

**Ludovic Lebart**

**Alain Morineau**

**Marie Piron**

---



# **Statistique exploratoire multidimensionnelle**



DUNOD

**Ludovic Lebart**

C.N.R.S.,  
École nationale supérieure  
des télécommunications

**Alain Morineau**

Centre international de statistique  
et d'informatique appliquées  
CISIA

**Marie Piron**

Institut français de recherche scientifique  
pour le développement en coopération,  
ORSTOM

# **Statistique exploratoire multidimensionnelle**

DUNOD

Ce pictogramme mérite une explication. Son objet est d'alerter le lecteur sur la menace que représente pour l'avenir de l'écrit, particulièrement dans le domaine de l'édition technique et universitaire, le développement massif du **photocopillage**.

Le Code de la propriété intellectuelle du 1er juillet 1992 interdit en effet expressément la photocopie à usage collectif sans autorisation des ayants droit. Or, cette pratique s'est généralisée dans les établisse-

ments d'enseignement supérieur, provoquant une baisse brutale des achats de livres et de revues, au point que la possibilité même pour les auteurs de créer des œuvres nouvelles et de les faire éditer correctement est aujourd'hui menacée.



Nous rappelons donc que toute reproduction, partielle ou totale, de la présente publication est interdite sans autorisation du Centre français d'exploitation du droit de copie (CFC, 3 rue Hauteleville, 75006 Paris).

© Dunod, Paris, 1995  
ISBN 2 10 002886 3

Toute représentation ou reproduction, intégrale ou partielle, faite sans le consentement de l'auteur, ou de ses ayants droit, ou ayants cause, est illicite (loi du 11 mars 1957, alinéa 1er de l'article 40). Cette représentation ou reproduction, par quelque procédé que ce soit, constituerait une contrefaçon sanctionnée par les articles 425 et suivants du Code pénal. La loi du 11 mars 1957 n'autorise, aux termes des alinéas 2 et 3 de l'article 41, que les copies ou reproductions strictement réservées à l'usage privé du copiste et non destinées à une utilisation collective d'une part, et d'autre part, que les analyses et les courtes citations dans un but d'exemple et d'illustration.

*Cet ouvrage est dédié à la mémoire de Brigitte Escofier*

## AVANT-PROPOS

Cet ouvrage s'adresse aux praticiens, scientifiques et étudiants de toutes disciplines qui ont à analyser et traiter de grands ensembles de données multidimensionnelles, c'est-à-dire finalement des recueils de données statistiques se présentant, totalement ou partiellement, sous forme de tableaux rectangulaires.

Le domaine d'application, limité au départ aux sciences de la vie (biométrie, agronomie, écologie) et aux sciences humaines (psychométrie, socio-économie), ne cesse de s'étendre car les possibilités offertes par les outils de traitement suscitent de nouveaux recueils de mesures. Les applications industrielles se développent rapidement et le contrôle de qualité, l'analyse des processus de production, la veille technologique, la recherche documentaire font de plus en plus appel à des ensembles de mesures multidimensionnelles.

On a tenté de faire le point sur les développements récents de la *statistique exploratoire multidimensionnelle* en continuité avec un ouvrage précédent<sup>1</sup> dont on reprend d'ailleurs, en les développant, certains chapitres. On s'est ainsi efforcé d'intégrer la substance de plusieurs centaines de publications (dont celles des auteurs) sur le thème de ce précédent travail.

Comme toujours pour ce type d'ouvrage qui s'adresse simultanément à des praticiens et des chercheurs de disciplines diverses, plusieurs lectures devraient être possibles selon les connaissances du lecteur notamment en mathématique et statistique : une lecture pratique, d'utilisateur, pour les personnes spécialisées dans les divers domaines d'application actuels et potentiels ; une lecture plus technique, complète, pour une personne ayant une formation en mathématiques appliquées et en statistique.

---

<sup>1</sup> *Technique de la description statistique*, (L. Lebart, A. Morineau, N. Tabard) Dunod, 1977.

La statistique exploratoire multidimensionnelle se prolonge naturellement et se diversifie en des outils et des modèles évidemment plus complexes que les méthodes de base. Mais l'essentiel des applications relèvent en fait de la partie la plus accessible. On a fait preuve d'une grande parcimonie dans l'utilisation de l'outil mathématique : le niveau d'abstraction choisi est toujours le niveau minimal compatible avec une présentation exacte, et la communication a été favorisée au détriment de la généralisation. Les lecteurs mathématiciens sauront sans difficulté introduire les notions qui permettent des formulations plus élégantes.

L'ensemble doit beaucoup à des collaborations et des cadres de travail divers : au sein du département Economie et Management, de l'Ecole Nationale Supérieure des Télécommunications et de l'URA820 du Centre National de la Recherche Scientifique (URA: Traitement et Communication de l'Information, dirigée par Jean-Pierre Tubach) ; au sein du Centre International de Statistique et d'Informatique Appliquées (CISIA), du Centre d'Etude de l'Emploi et de l'Institut français de recherche scientifique pour le développement en coopération (ORSTOM).

Nous remercions également les autres collègues, chercheurs ou professeurs auprès desquels nous avons puisé collaboration et soutien, ou simplement eu d'intéressants débats ou discussions, ou encore accès à des documents. Citons, sans être exhaustif, Mireille Bardos, Laurent Benzoni, Abdelhalim Bouamaine, Bernard Burtschy, Pierre Cazes, Frédéric Chateau, Jean-Pierre Fénelon, Christian Mullon, Jérôme Pagès, André Salem, Michel Tenenhaus, Gilbert Saporta et Wenhua Zhu.

Nous sommes heureux d'adresser ici nos remerciements à Gisèle Maïus et à Jean-Marc Quilbé, des éditions Dunod, pour l'accueil qu'ils ont réservé à cet ouvrage.

L. L., A. M., M. P.

Paris, Juillet 1995

# Sommaire

Introduction générale	1
-----------------------	---

## Chapitre 1

### MÉTHODES FACTORIELLES

Introduction	13
<b>Section 1.1 : Analyse générale, décomposition aux valeurs singulières</b>	15
1.1.1 Notions élémentaires et principe d'ajustement	15
1.1.2 Ajustement du nuage des individus	17
a- Droites d'ajustement	17
b- Caractéristiques du sous-espace d'ajustement	19
1.1.3 Ajustement du nuage des variables	19
1.1.4 Relation entre les ajustements dans les deux espaces	20
1.1.5 Reconstitution des données de départ	22
a- Reconstitution exacte	22
b- Reconstitution approchée	23
c- Qualité de l'approximation	24
1.1.6 Diversification de l'analyse générale	24
a- Analyse générale avec des métriques et des critères quelconques	24
b- Principe des éléments supplémentaires	27
c- Autres approches	28
1.1.7 Annexe 1 - Démonstration sur les extrema de formes quadratiques sous contraintes quadratiques	29
<b>Section 1.2 : Analyse en Composantes Principales</b>	32
1.2.1 Domaine d'application	32
1.2.2 Interprétations géométriques	33
a- Pour les $n$ individus	33
b- Pour les $p$ variables	34
1.2.3 Analyse du nuage des individus	34
a- Principe d'ajustement	34
b- Distance entre individus	36
c- Matrice à diagonaliser	37
d- Axes factoriels	37
1.2.4 Analyse du nuage des points-variables	38
a- distances entre points-variables	38
b- Distance à l'origine	39
c- Axes factoriels ou composantes principales	41

<b>1.2.5</b>	<b>Individus et variables supplémentaires</b>	42
	a- Individus supplémentaires	43
	b- Variables continues supplémentaires	43
	c- Variables nominales supplémentaires	43
<b>1.2.6</b>	<b>Représentation simultanée</b>	45
	a- Représentation séparée des deux nuages	45
	b- Justification d'une représentation simultanée	46
<b>1.2.7</b>	<b>Analyse en composantes principales non normée</b>	48
	a- Principe de l'analyse et nuage des individus	49
	b- Nuage des variables	50
<b>1.2.8</b>	<b>Analyses non-paramétriques</b>	51
	a- Analyse des rangs	51
	b- Analyse en composantes robustes	52
<b>1.2.9</b>	<b>Aperçu sur les autres méthodes dérivées</b>	53
<b>1.2.10</b>	<b>Éléments pour l'interprétation</b>	53
	a- Inertie liée aux facteurs	54
	b- Aides à l'interprétation	55
<b>1.2.11</b>	<b>Exemple d'application</b>	57
<b>Section 1.3 : Analyse des correspondances</b>		67
<b>1.3.1</b>	<b>Domaine d'application</b>	67
<b>1.3.2</b>	<b>Démarche et principe : une introduction élémentaire</b>	68
	a- Transformations du tableau de contingence	69
	b- Hypothèse d'indépendance	70
	c- Construction des nuages	71
	d- Critère d'ajustement	72
	e- Choix des distances	73
	f- Equivalence distributionnelle	74
	g- Relations de transition ou pseudo-barycentriques	75
	h- Justification de la représentation simultanée	78
<b>1.3.3</b>	<b>Schéma général de l'analyse des correspondances</b>	79
	a- Géométrie des nuages et éléments de base	79
	b- Critère à maximiser et matrice à diagonaliser	82
	c- Axes factoriels et facteurs	84
	d- Relation entre les deux espaces	84
	e- Relations de transition	85
	f- Représentation simultanée	86
	g- Autre présentation de l'analyse des correspondances	87
	h- Formule de reconstitution des données	89
<b>1.3.4</b>	<b>Règles d'interprétation : inertie, formes de nuages</b>	89
	a- Inertie et test d'indépendance	89
	b- Quelques formes caractéristiques de nuages de points	92
<b>1.3.5</b>	<b>Règles d'interprétation : contributions et cosinus</b>	94
	a- Contributions	94
	b- Cosinus carrés	95
	c- Exemple numérique	97
<b>1.3.6</b>	<b>Éléments supplémentaires</b>	99



1.3.7	<b>Mise en œuvre des calculs</b>	100
	a- Analyse par rapport à l'origine ou au centre de gravité du nuage	100
	b- Symétrisation de la matrice à diagonaliser	102
1.3.8	<b>Exemple d'application</b>	103
<b>Section 1.4 : Analyse des correspondances multiples</b>		108
1.4.1	<b>Domaine d'application</b>	108
1.4.2	<b>Notations et définitions</b>	109
	a- Hypercube de contingence	110
	b- Tableau disjonctif complet	110
	c- Tableau des faces de l'hypercube de contingence ou tableau de contingence de Burt	111
1.4.3	<b>Principes de l'analyse des correspondances multiples</b>	113
	a- Critère d'ajustement et distance du $\chi^2$	114
	b- Axes factoriels et facteurs	114
	c- Facteurs et relations pseudo-barycentriques	115
	d- Sous-nuage des modalités d'une même variable	117
	e- Support du nuage des modalités	117
	f- Meilleure représentation simultanée	118
	g- Inertie du nuage des modalités et conséquences pratiques	119
	h- Règles d'interprétation	120
	i- Principes du découpage en classes	121
1.4.4	<b>Éléments supplémentaires</b>	122
	a- Valeurs-test pour les modalités supplémentaires	123
	b- Variables continues supplémentaires	125
1.4.5	<b>Analyse du tableau de contingence de Burt : équivalence avec l'analyse du tableau disjonctif complet</b>	126
1.4.6	<b>Cas de deux questions</b>	127
1.4.7	<b>Cas particuliers</b>	130
	a- Toutes les questions ont deux modalités	131
	b- Sous-tableau d'un tableau de correspondances multiples	132
1.4.8	<b>Exemple d'application numérique</b>	135

## Chapitre 2

### QUELQUES MÉTHODES DE CLASSIFICATION

<b>Introduction</b>	145	
<b>Section 2.1 : Agrégation autour des centres mobiles</b>	148	
2.1.1	<b>Bases théoriques de l'algorithme</b>	148
2.1.2	<b>Justification élémentaire de l'algorithme</b>	150
2.1.3	<b>Techniques connexes</b>	151
2.1.4	<b>Formes fortes et groupements stables</b>	152

<b>Section 2.2 : Classification hiérarchique</b>	155
<b>2.2.1 Principe</b>	155
a- Distances entre éléments et entre groupes	156
b- Algorithme de classification	156
c- Éléments de vocabulaire	157
<b>2.2.2 Classification ascendante selon le saut minimal et arbre de longueur minimale</b>	159
a- Définition d'une ultramétrie	159
b- Équivalence entre ultramétrie et hiérarchie indicée	159
c- L'ultramétrie sous dominante	161
d- Arbre de longueur minimale: définition et généralités	163
e- Arbre de longueur minimale: algorithme de Kruskal	164
f- Arbre de longueur minimale: algorithme de Prim	165
g- Arbre de longueur minimale: algorithme de Florek et Sollin	165
h- Lien entre l'arbre et le saut minimal	166
<b>2.2.3 Critère d'agrégation selon la variance</b>	167
a- Notations et principe	168
b- Perte d'inertie par agrégation de deux éléments : le critère de Ward généralisé	170
<b>2.2.4 Algorithme de recherche en chaîne des voisins réciproques</b>	171
a- Algorithme	172
b- Critère de la médiane	173
<b>2.2.5 Exemple numérique d'application</b>	173
a- Classification des lignes (professions)	174
b- Classification des colonnes (médias)	175
<b>Section 2.3 : Classification mixte et description statistique des classes</b>	177
<b>2.3.1 Stratégie de classification mixte</b>	177
a- Les étapes de l'algorithme	177
b- Choix du nombre de classes par coupure de l'arbre	179
c- Procédure de consolidation	180
<b>2.3.2 Description statistique des classes</b>	181
a- Valeurs-test pour les variables continues	181
b- Valeurs-test pour les variables nominales	182
c- Variables caractéristiques d'une classe	184
<b>Section 2.4 : Complémentarité entre analyse factorielle et classification</b>	185
<b>2.4.1 Utilisation conjointe de l'analyse factorielle</b>	185
a- Nécessité... et insuffisance des méthodes factorielles	185
b- Mise en œuvre pratique dans le cas de la classification mixte	187
c- Autres aspects de la complémentarité	189
<b>2.4.2 Aspects techniques et théoriques de la complémentarité</b>	189
a- Classification des lignes ou colonnes d'un tableau de contingence	189
b- Un exemple de coïncidence entre les deux approches	190

<b>2.4.3 Valeurs propres et indices de niveau</b>	194
a- Quelques inégalités	194
b- Le cas des tables de contingence structurées par blocs	195
c- Une étude empirique du lien entre valeurs propres et indices	195
<b>2.4.4 La complémentarité en pratique : un exemple</b>	199
a- Les étapes	200
b- L'espace des variables actives	202
c- Exemples de description automatique de trois classes	202
d- Projection de variables signalétiques (en supplémentaires)	205

## Chapitre 3

### LIENS AVEC LES MÉTHODES EXPLICATIVES USUELLES, MÉTHODES DÉRIVÉES

<b>Introduction</b>	209
<b>Section 3.1 : Analyse canonique</b>	213
<b>3.1.1 Formulation du problème et notations</b>	213
<b>3.1.2 Les variables canoniques</b>	215
a- Calcul des variables canoniques	215
b- Interprétation géométrique	217
c- Cas de matrices non inversibles	218
<b>3.1.3 Liens avec l'analyse des correspondances</b>	219
a- Le cas de l'analyse des correspondances simples	219
b- L'analyse des correspondances multiples	220
<b>Section 3.2 : Régression multiple, modèle linéaire</b>	223
<b>3.2.1 Formulation du problème : le modèle linéaire</b>	223
<b>3.2.2 Ajustement par la méthode des moindres-carrés</b>	225
a- Calcul et propriétés de l'ajustement des moindres-carrés	226
b- Approche géométrique	227
c- Le coefficient de corrélation multiple	228
<b>3.2.3 Lien avec l'analyse canonique</b>	229
<b>3.2.4 Qualité de l'ajustement</b>	230
a- Spécification du modèle	230
b- Moyenne et variance des coefficients	231
c- Tests sous l'hypothèse de normalité des résidus	231
<b>3.2.5 Régression régularisée</b>	233
a- Principe de la régression régularisée	234
b- Variables supplémentaires et régression	236
c- Expression des coefficients dans la nouvelle base	236
<b>3.2.6 Régression sur variables nominales : analyse de la variance</b>	237
a- Codage des variables nominales	238

b-	Modèle linéaire sans interaction	239
c-	Modèle linéaire avec interaction	240
<b>3.2.7</b>	<b>Régression sur variables mixtes : analyse de la covariance</b>	<b>241</b>
a-	Modèles d'analyse de la covariance	242
b-	Test d'un effet différencié de $x$ dans chaque classe $k$	243
c-	Test de l'effet de la variable $u$	243
d-	Test d'un "effet classe global"	243
e-	Généralisation de l'analyse de la covariance	244
<b>3.2.8</b>	<b>Choix des variables, généralisations du modèle</b>	<b>244</b>
a-	Sélection et choix des variables explicatives	244
b-	Modèles linéaires généralisés	245
<b>3.2.9</b>	<b>Modèles de variables latentes</b>	<b>246</b>
a-	Le modèle	247
b-	Estimation des paramètres inconnus	249
<b>Section 3.3 : Analyse factorielle discriminante</b>		<b>251</b>
<b>3.3.1</b>	<b>Formulation du problème et notations</b>	<b>251</b>
<b>3.3.2</b>	<b>Fonctions linéaires discriminantes</b>	<b>253</b>
a-	Décomposition de la matrice de covariance	254
b-	Calcul des fonctions linéaires discriminantes	255
c-	Diagonalisation d'une matrice symétrique	256
<b>3.3.3</b>	<b>Cas de deux classes : équivalence avec la régression multiple</b>	<b>257</b>
<b>3.3.4</b>	<b>Lien avec d'autres méthodes</b>	<b>259</b>
a-	L'analyse canonique	259
b-	L'analyse des correspondances	260
c-	Une analyse en axes principaux avec une métrique particulière	262
<b>3.3.5</b>	<b>Principes des règles d'affectation (ou de classement)</b>	<b>263</b>
a-	Le modèle bayésien d'affectation	264
b-	Le modèle bayésien dans le cas normal	265
c-	Autres règles d'affectation	266
d-	Qualité des règles de classement	268
<b>3.3.6</b>	<b>Régularisation en analyse discriminante</b>	<b>269</b>
a-	Analyse régularisée de Friedman	270
b-	Analyse régularisée par axes principaux	270
<b>3.3.7</b>	<b>Discrimination sur variables nominales</b>	<b>275</b>
a-	Analyse factorielle discriminante qualitative	275
b-	Analyse discriminante barycentrique	276
c-	Note sur le "scoring"	276
<b>3.3.8</b>	<b>Discrimination et réseaux de neurones</b>	<b>277</b>
a-	Schéma et modèle du perceptron multi-couches	278
b-	Modèles non-supervisés ou auto-organisés	280
c-	Statistique et méthodes neuronales	282
<b>Section 3.4 : Modèles log-linéaires</b>		<b>284</b>
<b>3.4.1</b>	<b>Formulation du problème et principes de base</b>	<b>284</b>
<b>3.4.2</b>	<b>Ajustement d'un modèle log-linéaire</b>	<b>285</b>

a- Tableau de contingence à deux entrées	285
b- Tableau de contingence à $p$ entrées	286
c- modèles hiérarchiques	287
<b>3.4.3 Estimation et tests d'ajustement du modèle</b>	288
a- Estimation des paramètres	288
b- Tests d'ajustement	289
c- Choix du modèle	289
<b>3.4.4 La régression logistique</b>	290
a- Le modèle logistique	291
b- Estimation et tests des coefficients	293
c- Comparaison de deux modèles	294
d- Modèle avec interaction	294
<b>3.4.5 Modèles log-linéaire et analyse des correspondances</b>	295
a- Des champs d'application différents	295
b- Liens théoriques entre l'analyse des correspondances et les modèles log-linéaires	298
c- Difficultés de l'articulation exploration-inférence	300
 <b>Section 3.5 : Segmentation</b>	 302
<b>3.5.1 Formulation du problème, principe et vocabulaire</b>	303
<b>3.5.2 Construction d'un arbre de décision binaire</b>	304
a- Algorithme général de segmentation	304
b- Cas de la régression	306
c- Cas de la discrimination	309
<b>3.5.3 Sélection du "meilleur sous-arbre"</b>	312
a- Procédures de sélection	313
b- Estimation de l'Erreur Théorique de Prévision	314
c- Estimation du Taux d'Erreur Théorique de classement	314
<b>3.5.4 Divisions équi-réductrices et équi-divisantes</b>	316
a- Divisions équi-réductrices	316
b- Divisions équi-divisantes	316
<b>3.5.5 Lien avec les méthodes de classement</b>	317
 <b>Section 3.6 : Analyses partielles et projetées</b>	 319
<b>3.6.1 Définition du coefficient de corrélation partielle</b>	319
<b>3.6.2 Calcul des covariances et corrélations partielles</b>	320
a- Cas de deux variables	320
b- Cas de $p$ variables ( $X$ ) et de $q$ variables ( $Z$ )	321
<b>3.6.3 Analyse du nuage résiduel ou analyse partielle</b>	322
<b>3.6.4 Autres analyses partielles ou projetées</b>	323
a- Analyse canonique des correspondances	324
b- Analyse non-symétrique des correspondances	325
 <b>Section 3.7 : Structures de graphe, analyses locales</b>	 327
<b>3.7.1 Variance locale et covariance locale d'une variable</b>	328

<b>3.7.2</b>	<b>Coefficient de contiguïté de Geary</b>	329
<b>3.7.3</b>	<b>Analyse locale</b>	331
<b>3.7.4</b>	<b>Analyse de contiguïté et projections révélatrices</b>	331
	a- Analyse de contiguïté	331
	b- Représentation de groupes par projection	332
	c- Liens avec les analyses partielles	333
<b>3.7.5</b>	<b>Extensions, généralisations, applications</b>	334
<b>3.7.6</b>	<b>Cas particuliers : Structure de partition</b>	335
	a- Analyse inter-classes	335
	b- Analyse intra-classes	336
<b>Section 3.8 : Tableaux multiples, groupes de variables</b>		<b>337</b>
<b>3.8.1</b>	<b>Quelques travaux de référence</b>	337
<b>3.8.2</b>	<b>Analyses procrustéennes</b>	339
	a- Analyse procrustéenne orthogonale	339
	b- Analyse procrustéenne sans contrainte	341
	c- Formulaire de quelques méthodes d'analyse	341
<b>3.8.3</b>	<b>Méthode STATIS</b>	342
	a- Notations	342
	b- Comparaison globale entre les tableaux : l'interstructure	342
	c- Le nuage moyen ou compromis : l'intrastructure	343
	d- Représentation simultanée des nuages partiels : les trajectoires	344
<b>3.8.4</b>	<b>Analyse factorielle multiple</b>	344
	a- Une analyse en composantes principales pondérée	344
	b- Recherche de facteurs communs (intrastructures)	345
	c- Représentation des groupes de variables (interstructure)	346
	d- Représentations superposées des nuages partiels des groupes actifs (trajectoires)	346
<b>3.8.5</b>	<b>Analyse canonique généralisée</b>	347
	a- Formulation générale	348
	b- Propriétés de l'Analyse Canonique Généralisée	349
	c- Utilisation en pratique de l'analyse canonique généralisée	352

## Chapitre 4

### VALIDITÉ ET PORTÉE DES RÉSULTATS

<b>Introduction</b>	357
<b>Section 4.1 : Signification des valeurs propres et des taux d'inertie</b>	<b>359</b>
<b>4.1.1 Travaux sur la loi des valeurs propres</b>	<b>359</b>
<b>4.1.2 Approximation de la distribution des valeurs propres en analyse des correspondances</b>	<b>361</b>

---

<b>4.1.3</b>	<b>Indépendance des taux d'inertie et de la trace</b>	364
<b>4.1.4</b>	<b>Exemples d'abaques et tables statistiques</b>	366
<b>4.1.5</b>	<b>Taux d'inertie et information</b>	368
	a- Caractère partiel des taux d'inertie	368
	b- Quelle information?	371
<b>4.1.6</b>	<b>Choix du nombre d'axes : quelques résultats utiles</b>	373
	a- Règles empiriques	374
	b- Procédures externes	374
	c- Critères de choix statistiques, résultats asymptotiques	375
<b>Section 4.2 : Stabilité des axes, des formes, des classes</b>		379
<b>4.2.1</b>	<b>Méthodes de validation empiriques</b>	379
	a- Calculs de stabilité et de sensibilité	379
	b- Epreuves empiriques de stabilité	382
<b>4.2.2</b>	<b>Méthodes de rééchantillonnage</b>	385
	a- Technique de Jackknife	385
	b- Technique de Bootstrap	387
	c- Validation croisée	388
<b>4.2.3</b>	<b>Zones de confiance, zones de garde</b>	389
	a- Zones de confiance établies par bootstrap	389
	b- Autres types de simulation bootstrap	392
	c- Zones de garde en analyse des correspondances	396
	d - Autres régions de confiances	398
<b>4.2.4</b>	<b>Nombre de classes et validation des classifications</b>	399
	a- L'hypothèse d'absence de structure, les modèles	400
	b- Combien de classes retenir ?	402
	c- Les critères externes	404
<b>Références bibliographiques</b>		405
<b>Index des auteurs</b>		429
<b>Index des matières</b>		434

# Introduction

La statistique descriptive permet de représenter de façon vivante et assimilable des informations statistiques en les simplifiant et les schématisant. La statistique descriptive multidimensionnelle en est la généralisation naturelle lorsque ces informations concernent plusieurs variables ou dimensions.

Mais le passage au multidimensionnel induit un changement qualitatif important. On ne dit pas en effet que des microscopes ou des appareils radiographiques sont des instruments de description, mais bien des instruments d'observation ou d'exploration, et aussi de recherche. La réalité multidimensionnelle n'est pas seulement simplifiée parce que complexe, mais aussi explorée parce que cachée.

Le travail de préparation et de codage des données, les règles d'interprétation et de validation des représentations fournies par les techniques utilisées dans le cas multidimensionnel n'ont pas la simplicité rencontrée avec la statistique descriptive élémentaire. Il ne s'agit pas seulement de présenter mais d'analyser, de découvrir, parfois de vérifier et prouver, éventuellement de mettre à l'épreuve certaines hypothèses.

C'est pourquoi nous avons choisi de parler dans cet ouvrage de statistique exploratoire multidimensionnelle.

## La statistique et l'informatique

Née au tout début du vingtième siècle, notamment à la suite des travaux du précurseur l'astronome Quételet et des démographes et biométriciens Galton, Pearson, puis Fisher, la science statistique aura manipulé des chiffres pendant un demi-siècle sans disposer de véritables outils de calcul. Les appareils que l'on trouve maintenant dans la poche des écoliers et dans tous les bureaux auraient comblé les aspirations les plus insensées des statisticiens jusqu'en 1960. "Il est impensable d'utiliser des méthodes conçues avant l'avènement de l'ordinateur, il faut complètement réécrire la statistique", écrivait en substance Jean-Paul Benzécri dès 1965 dans son cours à la Sorbonne sur *l'Analyse des données et la reconnaissance des formes*.

Cet auteur, qui a profondément marqué le développement des recherches statistiques au cours des années récentes, préconise aussi, de manière un peu provocante pour une discipline où la notion de modèle a joué un rôle central : "le modèle doit suivre les données et non l'inverse".



Aux États-Unis, John Tukey, le fondateur du courant désigné par *Exploratory Data Analysis (EDA)*, a une attitude aussi radicale (cf. Mallows et Tukey, 1982). Il s'en faut cependant de beaucoup que ces deux pionniers aient été unanimement entendus. À défaut d'être repensée, la statistique s'est cependant considérablement enrichie. La période récente a connu des changements tout à fait notables du fait de la diffusion des moyens de calcul : les outils existants ont été améliorés, de nouveaux outils sont apparus, de nouveaux domaines d'application ont été explorés.

### *Meilleurs graphiques*

L'informatique, surtout la micro-informatique, a rendu familiers tous les outils graphiques de la statistique descriptive élémentaire. Autrefois fruits d'un travail laborieux et coûteux, ces représentations sont immédiatement accessibles dans pratiquement tous les logiciels intégrés. Les techniques de statistique exploratoire multidimensionnelle mettent à profit ces interfaces graphiques pour représenter, par exemple, les espaces factoriels et les arbres de classification : c'est là l'une de leurs fonctions iconographiques qui généralise effectivement la statistique descriptive usuelle au cas de variables nombreuses.

### *Désuétude des tables statistiques*

Classiquement, pour savoir si une quantité, dont la distribution est connue, ne dépasse pas les limites que lui assignent certaines hypothèses, on consultait la table donnant les valeurs que cette quantité ne dépassera que dans 5% ou 1% des cas. Le choix de seuils était imposé par la nécessité de limiter le volume des tables. À partir du moment où la quantité à tester est elle-même calculée sur ordinateur, il est facile d'ajouter au programme une procédure de calcul de la probabilité de dépassement de la valeur calculée. On gagne en confort, mais aussi en performance, car on pourra désormais comparer et trier des statistiques différentes grâce aux probabilités de dépassement, comme celles liées aux tests fishériens, évoqués au paragraphe suivant (au delà de la théorie classique des tests).

### *Emphase sur la robustesse, le non-paramétrique*

La mise en œuvre de la plupart des procédures inférentielles classiques est hypothéquée par la pertinence des hypothèses techniques<sup>1</sup> et par la sensibilité éventuelle des résultats à la non-vérification de ces hypothèses.

---

<sup>1</sup> Contrairement aux hypothèses générales qui sont les hypothèses d'ordre scientifique qui régissent l'étude d'un phénomène, et qui précèdent la phase d'observation ou d'expérimentation statistique, les hypothèses techniques interviennent dans la mise en œuvre pratique des méthodes statistiques. Elles concernent principalement la spécification des modèles et des distributions statistiques impliquées dans ces modèles. Certaines hypothèses techniques n'ont aucun lien avec les hypothèses générales, mais sont au contraire des exigences du modèle utilisé (exemple: les résidus sont indépendants et suivent une loi normale dont la matrice des covariances doit être spécifiée dans le cas de la régression linéaire multiple).

L'un des principaux obstacles à l'utilisation d'estimateurs robustes, c'est-à-dire peu sensibles à la présence de points aberrants (vis-à-vis des distributions étudiées), était la difficulté des calculs à mettre en œuvre. La plupart des panoplies existantes se sont donc enrichies de procédures plus robustes dès l'apparition de moyens de calcul plus puissants. Pour des raisons analogues, les techniques non-paramétriques qui s'affranchissent des hypothèses les plus lourdes ont connu un regain d'intérêt, comme ce fût le cas des techniques non-paramétriques de discrimination.

Les test "Fishériens", ou tests de permutation<sup>1</sup>, connaissent également un renouveau important. Les hypothèses statistiques sont éprouvées par permutations aléatoires de l'ensemble fini des observations effectivement disponibles : il y aura donc coïncidence entre les distributions marginales théoriques et observées. Seul l'obstacle du calcul pouvait faire écarter des techniques fondées sur des hypothèses qui épousent aussi étroitement la réalité. Mais les habitudes des praticiens (et aussi le coût de formation, la maîtrise des méthodes) sont telles qu'on ne peut attendre une substitution rapide des outils.

### *Taille et complexité des problèmes*

Il n'est pas rare maintenant de traiter des tableaux correspondant à des milliers d'observations et des centaines de variables. Bien sûr, les données les plus volumineuses et les plus complexes ont pu être abordées à l'aide d'outils préexistants. Mais très vite, l'adage: "c'est l'échelle qui fait le phénomène" s'est trouvé vérifié. Le changement d'échelle du volume des données a rapidement conduit à modifier les outils eux-mêmes et à imaginer de nouveaux outils dans le cadre de nouvelles approches.

### *Méthodes algorithmiques*

La levée de l'obstacle du calcul a eu pour effet de diffuser l'emploi des techniques de type algorithmique, au premier rang desquelles se trouvent les techniques de classification automatique et les méthodes impliquant des algorithmes coûteux (comme les diagonalisations de matrices par exemple). D'autres techniques, comme les techniques de sélection pas-à-pas, les techniques d'estimation par la méthode du maximum de vraisemblance, de programmation dynamique, connaissent des utilisations de plus en plus fréquentes.

### *Traitement des variables qualitatives*

L'étude statistique des variables qualitatives est par nature plus complexe que celle des variables numériques continues, qui s'appuie généralement sur la loi normale et sur les formalismes simples qui en dérivent (maximum de vraisemblance, moindres carrés, par exemple). Il n'est donc pas étonnant que les possibilités de calcul aient permis de fortes avancées

---

<sup>1</sup> Cf. sur les tests dits "exacts" : Mehta *et al.* (1991), Agresti (1992), Good (1994).

dans ce domaine : analyse des correspondances simples et multiples dans le cas descriptif, modèles log-linéaires, modèles logistiques dans le cas inférentiel.

### *Méthodes de validation*

Les techniques de simulation (ou de Monte-Carlo) connaissent des applications à grande échelle dans tous les domaines où les hypothèses distributionnelles usuelles sont inadaptées. La simulation permet de construire de l'inférence "sur-mesure" en combinant des sources, des formes et des niveaux de variabilité dans des processus complexes dont la formalisation est rigoureusement impossible. Mais le sur-mesure est plus coûteux que le prêt-à-porter.

Les techniques de rééchantillonnage telle que les techniques de "Jackknife" (la variabilité est étudiée en procédant à des prélèvements sans remise dans l'échantillon) et de "Bootstrap" (la variabilité est étudiée en procédant à des tirages pseudo-aléatoires avec remise dans l'échantillon) ont le mérite d'avoir donné lieu à des développements théoriques. A l'heure actuelle, le Bootstrap, qui présente de notables avantages (taille d'échantillon inchangée, facilité de mise en œuvre, propriétés théoriques satisfaisantes) est assez largement utilisé.

Les techniques de *validation croisée* sont surtout utilisées en analyse discriminante : pour estimer un vrai taux d'erreur, il convient de tester la méthode sur des individus ne faisant pas partie de l'échantillon d'apprentissage. D'où l'idée de procéder à  $n$  analyses discriminantes sur  $(n-1)$  individus, en retirant à chaque fois un individu de l'échantillon d'apprentissage, puis en notant le succès ou l'échec de son affectation. Ces principes de base peuvent être réaménagés et adaptés, notamment au cas des grands tableaux, mais on devine que le gain d'information réalisé a sa contrepartie en volume de calcul.

### *Réseaux neuronaux*

Les techniques neuronales ou connexionnistes ont une large intersection avec les méthodes classiques d'analyse des données<sup>1</sup>, intersection peu visible de prime abord en raison d'une terminologie et d'un cadre conceptuel tout à fait spécifiques. Inspirées à l'origine par des modèles de fonctionnement du cerveau, les méthodes connexionnistes peuvent être considérées comme des méthodes d'analyse non-linéaire des données. L'analyse en composantes principales, les méthodes de classification du type k-means ou nuées dynamiques sont des méthodes neuronales non supervisées ; la régression, l'analyse discriminante linéaire, des cas particuliers de méthodes neuronales supervisées.

---

<sup>1</sup> L'expression anglaise *data analysis* a un sens très général de statistique appliquée (avec une connotation d'approche pragmatique et informatisée). L'équivalent anglais de l'analyse des données serait à peu près *multivariate data analysis*.

### *Les logiciels*

Une des innovations de forme, sinon de fond, de ces dernières années aura été la matérialisation des méthodes et des techniques sous forme de "produits", les logiciels, développés avec des contraintes économiques et commerciales de conception, de production, de distribution. Comme tout produit fini, le logiciel a l'avantage de diffuser et l'inconvénient de figer. Comme tout produit coûteux, il introduit une discrimination par les moyens financiers disponibles. Comme tout produit à l'usage de spécialistes, il introduit de nouvelles divisions du travail, parfois peu souhaitables dans un processus de connaissance. Enfin, si cette division du travail se fait à l'échelle internationale, de nouvelles dépendances sont créées dans des secteurs sensibles : l'acquisition de connaissances, la recherche fondamentale.

Ces avantages et inconvénients sont indissolublement liés dans les logiciels statistiques. Les logiciels accessibles et faciles à utiliser permettront une large diffusion des méthodes, mais donneront parfois lieu à des utilisations inconsidérées dans des domaines où une réflexion minutieuse et une grande prudence seraient de mise. La médiation des logiciels est un nouveau paramètre dont il faut tenir compte<sup>1</sup>.

### *Nouveaux domaines d'application*

L'informatisation et les outils qu'elle a suscité ou dont elle a stimulé le développement (gestionnaires de base de données relationnelles, systèmes d'informations géographiques par exemple) ont pour effet le plus évident de permettre le traitement statistique de recueils de données plus grands et plus complexes, donnant lieu à de véritables systèmes d'information. Les méthodes d'analyse des données peuvent être des outils performants pour exploiter au mieux la structure organisée de ces systèmes.

On peut citer parmi les domaines récemment abordés: les analyses d'images, les analyses de séquences d'images (données de télédétection par exemple); les analyses de signaux, de processus, de systèmes; la recherche documentaire; les analyses de données textuelles; les analyses de grandes enquêtes.

---

<sup>1</sup> Les activités d'un club comme MODULAD (domicilié à l'INRIA) doivent pallier certains des inconvénients cités. Rassemblant des créateurs, des développeurs, des utilisateurs de logiciels, il doit faciliter certains types de communications et de diffusions. Les étudiants ou chercheurs ont ainsi accès, dans la bibliothèque de programme MODULAD, au "source" des programmes. Naturellement, les faibles moyens mis en oeuvre ne permettent pas de mener à bien les coûteuses opérations d'habillage, d'assurer les qualités de convivialité nécessaires et des mises à jour en fonction des nouveaux matériels et langages. Cette bibliothèque, ainsi que les listages de programmes publiés dans les ouvrages "Techniques de la description statistique" (L. Lebart, A. Morineau, N. Tabard. Dunod, 1977) et "Traitement des données statistiques" (L. Lebart, A. Morineau, J.-P. Fénelon. Dunod, 1979) peuvent donner accès à la plupart des traitements proposés dans cet ouvrage. Les traitements correspondant aux exemples ont été réalisés à l'aide du logiciel SPAD.N (Lebart *et al.*, 1991), actuellement développé et distribué par le CISIA.

## Panorama du contenu de ce manuel

Les avancées et innovations qui viennent d'être évoquées se retrouvent à des degrés divers dans le développement et la mise en œuvre de la statistique exploratoire multidimensionnelle, à laquelle est consacrée le présent ouvrage.

La gamme des méthodes qui permettent de décrire et d'explorer des tableaux de données statistiques (tableaux mesures-observations, tableaux de contingence ou tableaux croisés, tableaux de présence-absence ou tableaux d'incidence) est assez étendue.

Celles que nous retiendrons seront choisies en fonction de leur aptitude à traiter de tableaux volumineux, de la transparence de leur fonctionnement, de leur bonne insertion dans l'éventail des méthodes réellement applicables et appliquées.

Deux grandes familles de méthodes répondent à ces exigences :

- [chapitre 1] : *les méthodes factorielles*<sup>1</sup>, fondées sur des recherches d'axes principaux (l'analyse en composantes principales et les analyses des correspondances simples et multiples sont les méthodes factorielles les plus utilisées) qui produisent essentiellement des visualisations graphiques planes ou tridimensionnelles des éléments à décrire.

- [chapitre 2] : *les méthodes de classification* qui produisent des groupements en classes d'objets (ou en familles de classes hiérarchisées), obtenus à la suite de calculs algorithmiques. Les éléments à décrire sont groupés de la manière la moins arbitraire possible à partir de leurs vecteurs de description.

Les points de vue fournis par ces deux types de méthodes sont en fait très complémentaires. On insistera sur cette complémentarité qui se manifeste d'ailleurs à plusieurs niveaux, qu'il s'agisse de la possibilité d'appréhender des structures très diverses, ou d'aider à la lecture des résultats.

Lorsqu'on a peu d'information *a priori* sur les données (on parlera alors de données non structurées ou amorphes) l'application des techniques exploratoires multidimensionnelles est gratifiante. Mais il est plus difficile d'utiliser ce que l'on sait pour essayer d'en savoir plus. Et si l'information *a priori* sur les données est considérable, d'autres techniques faisant appel à des modèles qui utilisent effectivement cette information sont alors compétitives.

---

<sup>1</sup> Les techniques d'analyse factorielle comprennent dans la littérature statistique française des vingt dernières années toutes les techniques de représentation utilisant des "axes principaux": analyse en composantes principales, des correspondances simples et multiples, analyse factorielle dite classique ou des psychologues — alors que l'expression correspondante en anglais (factor analysis) ne désigne de façon assez stricte que cette dernière technique : analyse en facteurs communs et spécifiques de Spearman, Thurstone, utilisée principalement par les psychologues et les psychométriciens.

- [chapitre 3] : les liens avec les méthodes explicatives usuelles, éclaireront les utilisateurs sur la vocation spécifique de chacune de ces méthodes. Les cinq premières sections de ce chapitre présentent successivement l'analyse canonique, la régression multiple et le modèle linéaire, l'analyse discriminante, les modèles log-linéaires et logistiques, les méthodes de segmentation. Cet éventail de techniques recouvre une part très importante des applications potentielles de la statistique.

Il n'existe cependant pas de méthodologie générale de mise en œuvre des méthodes exploratoires de base impliquant une articulation et une synergie avec les méthodes dites explicatives. Chaque application demande un travail original de codage, de sélection et d'agencement d'outils particuliers en fonction des domaines et des problèmes.

Les méthodes d'analyse de tableaux ayant une structure *a priori* présentées dans les trois sections suivantes du chapitre 3 constituent le complément naturel ou le prolongement des analyses exploratoires. Elles présentent les techniques qui tentent d'intégrer en leur sein même une éventuelle information externe : les analyses partielles ou conditionnelles permettent de prendre en compte l'effet de certaines variables ; les analyses de contiguïté mettent à profit des structures de graphes sur les observations (contenant comme cas particulier les partitions et les séries chronologiques) ; enfin les analyses de tableaux multiples étudient le cas de tableaux comportant plusieurs groupes de variables.

- [chapitre 4] : la validité et la portée des résultats sont deux thèmes d'études qui ont donné lieu à des recherches nombreuses au cours des années récentes. Dans une première section, on fait le point sur les résultats théoriques disponibles (difficilement acquis et peu utilisables en pratique) puis, dans la seconde section, on présente quelques procédures plus empiriques, plus souples, incluant les techniques de rééchantillonnage, parmi lesquelles le *Bootstrap* jouera un rôle prédominant.

## Les méthodes descriptives et exploratoires de base

Les méthodes étudiées dans les deux premiers chapitres sont destinées à fournir des représentations et des réductions, complémentaires, de l'information contenue dans de volumineux tableaux de données numériques. D'autres méthodes de description qui ne rentrent pas dans les deux familles étudiées ici ne seront évoquées que brièvement, comme les méthodes purement graphiques<sup>1</sup>, dévolues à la représentation de tableaux

---

<sup>1</sup> Parmi les méthodes purement graphiques, citons la méthode des visages de Chernoff (1973), pour laquelle chaque visage correspond à un individu et chaque trait du visage à une variable; la méthode des courbes d'Andrews (1972), où les différents paramètres des courbes sont les variables; la méthode des constellations de Wakimoto et Taguri (1978), dans laquelle, après conversion de chaque  $x_{ij}$  (valeur de la variable  $j$  pour l'individu  $i$ ) en un  $\cos\theta_{ij}$ , chaque individu  $i$  est représenté par un point du plan complexe comme une somme de variables de modules constants et d'arguments  $\theta_{ij}$ .

de petites dimensions, les méthodes de sériation<sup>1</sup>, les méthodes de *multidimensional scaling*<sup>2</sup>.

Elles interviennent souvent dans des contextes particuliers d'application et sont moins adaptées aux traitements des grands tableaux.

Le tableau de données sur lequel sont effectuées les réductions ne sera pas en général un tableau de valeurs numériques quelconques. Il doit en particulier présenter une certaine homogénéité de forme et de contenu.

### *Représentation géométrique élémentaire d'un tableau de données*

Le tableau de données dispose la masse d'information sous forme rectangulaire.

Pour fixer les idées, les lignes ( $i=1,\dots,n$ ) peuvent représenter les  $n$  *individus* ou *observations*, appelés plus généralement *unités statistiques*; les colonnes ( $j=1,\dots,p$ ) sont alors les  $p$  *variables*, qui peuvent être des *mesures* (numériques) ou des *attributs* ou *caractères* observés sur les individus (cas de variables nominales)<sup>3</sup>.

Afin de comprendre le principe des méthodes de statistique exploratoire multidimensionnelle, il est utile de représenter géométriquement les  $n$  lignes et les  $p$  colonnes du tableau de données par des points dont les coordonnées sont précisément les éléments de ce tableau (figure 1).

Deux nuages de points sont alors construits :

- le nuage des  $n$  individus (le nuage des points-lignes) situé dans l'espace à  $p$  dimensions  $\mathbb{R}^p$  des variables (des colonnes); chacune des  $n$  lignes est représentée par un point à  $p$  coordonnées.
- le nuage des  $p$  variables (le nuage des points-colonnes) situé dans l'espace à  $n$  dimensions  $\mathbb{R}^n$  des individus (des lignes); chacune des  $p$  colonnes est représentée par un point à  $n$  coordonnées.

Le tableau de données noté  $X$  est donc une matrice dans laquelle chaque vecteur, ligne ou colonne, représente un point soit dans  $\mathbb{R}^p$  soit  $\mathbb{R}^n$ .

<sup>1</sup> Les méthodes de sériations visent à faire apparaître des structures particulières de tableaux par simple réordonnement de lignes et de colonnes. Pour des exposés de synthèse sur ce sujet, cf. par exemple Arabie (1978), Caraux (1984), Marcotorchino (1987).

<sup>2</sup> Cf. Shepard (1974), Kruskal et Wish (1978), Schiffman *et al.* (1981).

<sup>3</sup> Cette distinction entre variables et individus est commode parce qu'elle se réfère à une situation classique en statistique. Elle correspond au contexte de l'analyse en composantes principales (section 1.2) qui précède historiquement l'analyse des correspondances et ses variantes. Cette distinction n'a évidemment pas de sens dans le cas de tables de contingence pour lesquelles lignes et colonnes jouent des rôles symétriques.

Chacune des deux dimensions du tableau de données permet de définir des distances (ou des proximités) entre les éléments définissant l'autre dimension.

L'ensemble des colonnes permet de définir, à l'aide de formules appropriées, des distances entre lignes. De la même façon, l'ensemble des lignes permet de calculer des distances entre colonnes.

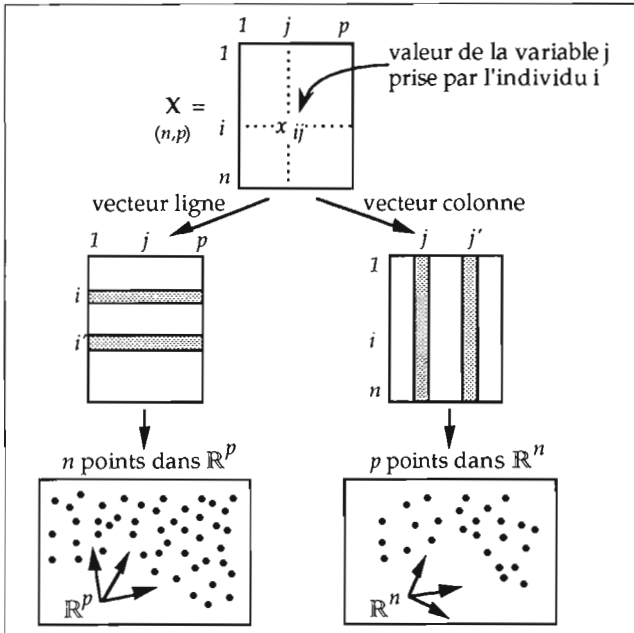


Figure 1  
Principe de représentation géométrique

Les *proximités géométriques* usuelles entre points-lignes et entre points-colonnes traduisent en fait des *associations statistiques* soit entre les individus, soit entre les variables. Les tableaux de distances associés à ces représentations géométriques (simples dans leur principe, mais complexes en raison du grand nombre de dimensions des espaces concernés) pourront alors être décrits par les deux grandes familles de méthodes que sont les méthodes factorielles et la classification (figure 2).

Ces représentations géométriques du tableau de données nous conduisent naturellement à utiliser les notions d'espaces vectoriels, de nuages de points, de métriques (permettant de calculer des distances entre points-lignes ou entre points-colonnes) mais aussi de masses affectées aux points si l'on ne leur accorde pas la même importance dans le nuage.

Les développements théoriques des méthodes de statistique exploratoire multidimensionnelle vont reposer sur ces notions.



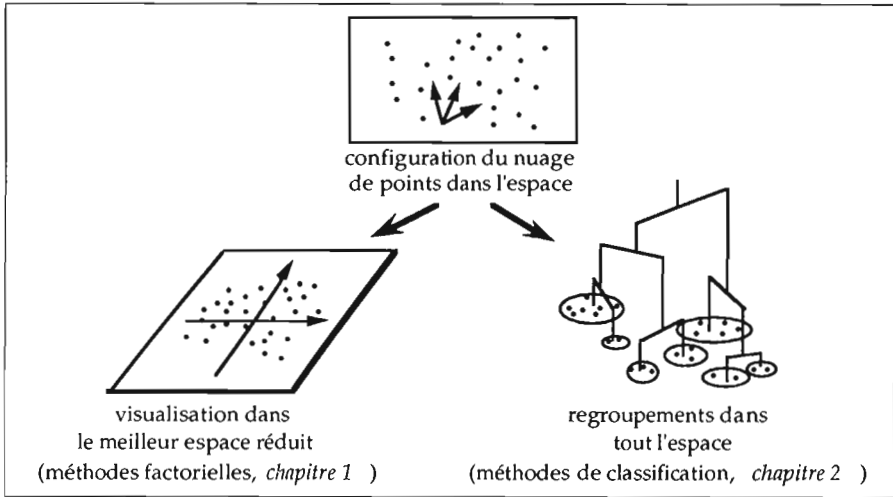


Figure 2  
Les deux grandes familles de méthodes

Ces méthodes impliquent souvent de la même manière les individus (lignes) et les variables (colonnes). Les individus ne sont plus de simples intermédiaires utilisés pour calculer des moyennes ou des corrélations sur les variables, suivant le schéma de la statistique traditionnelle où ils ne sont que des réalisations d'épreuves indépendantes. La confrontation des espaces d'individus et de variables enrichira les interprétations.

#### Notations de base

Malgré leur partielle inadaptation aux éléments mathématiques que l'on va traiter, les notations matricielles seront souvent utilisées par souci de cohérence et volonté de communication avec l'essentiel de la littérature statistique disponible.

Le tableau des données soumis à l'analyse est désigné par la lettre majuscule grasse  $X$ . La matrice  $X$  est d'ordre  $(n,p)$ , autrement dit, elle a  $n$  lignes et  $p$  colonnes. Son terme générique est  $x_{ij}$  ( $i^{\text{ème}}$  observation de la  $j^{\text{ème}}$  variable). Une colonne de  $X$  sera désignée par la lettre minuscule grasse  $x_j$ .

La transposée de  $X$  est notée  $X'$ ; cette matrice a donc  $p$  lignes et  $n$  colonnes.

Sauf mention contraire, pour les notations utilisant des caractères latins, les matrices sont représentées par des lettres majuscules grasses; les vecteurs par des lettres minuscules grasses; et les scalaires par des lettres minuscules en italique.

Chapitre 1

---

# **MÉTHODES FACTORIELLES**



## Introduction

Les méthodes factorielles se proposent de fournir des représentations synthétiques de vastes ensembles de valeurs numériques, en général sous forme de visualisations graphiques.

Pour cela, on cherche à réduire les dimensions du tableau de données en représentant les associations entre individus et entre variables dans des espaces de faibles dimensions.

Il est toujours possible de calculer des distances entre les lignes et entre les colonnes d'un tableau rectangulaire de valeurs numériques, mais il n'est pas possible de visualiser ces distances de façon immédiate (les représentations géométriques associées impliquant en général des espaces à plus de deux ou trois dimensions) : il est nécessaire de procéder à des transformations et des approximations pour en obtenir une représentation plane.

C'est une des tâches dévolues à l'analyse factorielle au sens large : opérer une réduction de certaines représentations "multidimensionnelles".

On recherchera donc des sous-espaces de faibles dimensions (une, deux ou trois par exemple) qui ajustent au mieux le nuage de points-individus et celui des points-variables, de façon à ce que les proximités mesurées dans ces sous-espaces reflètent autant que possible les proximités réelles. On obtient ainsi un espace de représentation, l'espace factoriel.

Mais la géométrie des nuages de points et les calculs de proximités ou de distances qui en découlent diffèrent selon la nature des lignes et des colonnes du tableau analysé.

Les colonnes peuvent être des variables continues ou des variables nominales ou des catégories dans le cas des tables de contingences. Les lignes peuvent être des individus ou des catégories.

La nature des informations, leur codage, les spécificités du domaine d'application vont introduire des variantes au sein des méthodes factorielles.

On présente ici trois techniques fondamentales :

- *l'analyse en composantes principales* (section 1.2) s'applique aux tableaux de type "variables-individus", dont les colonnes sont des variables à valeurs numériques continues et dont les lignes sont des individus, des observations, des objets, etc. Les proximités entre variables s'interprètent en termes de corrélation ; les proximités entre individus s'interprètent en termes de similitudes globales des valeurs observées. Elle peut donner lieu à de nombreuses variantes en s'appliquant par exemple à un tableau

de rangs (diagonalisation de la matrice de corrélation des rangs de Spearman), ou encore après l'élimination de l'effet de certaines variables (analyses locales ou partielles).

- *l'analyse des correspondances* (section 1.3) s'applique aux tableaux de contingences, c'est-à-dire aux tableaux de comptages obtenus par le croisement de deux variables nominales. Ces tableaux ont la particularité de faire jouer un rôle identique aux lignes et aux colonnes. L'analyse fournit des représentations des associations entre lignes et colonnes de ces tableaux, fondées sur une distance entre profils (qui sont des vecteurs de fréquences conditionnelles) désignée sous le nom de distance du  $\chi^2$ .
- *l'analyse des correspondances multiples* (section 1.4) est une extension du domaine d'application de l'analyse des correspondances, avec cependant des procédures de calcul et des règles d'interprétation spécifiques. Elle fait l'objet d'une mention particulière en raison de l'étendue de son champ d'application. Elle est particulièrement adaptée à la description de grands tableaux de variables nominales dont les fichiers d'enquêtes socio-économiques ou médicales constituent des exemples privilégiés. Les lignes de ces tableaux sont en général des individus ou observations (il peut en exister plusieurs milliers); les colonnes sont des modalités de variables nominales, le plus souvent des modalités de réponses à des questions.

Les techniques les plus utilisées dérivent des deux techniques fondamentales que sont l'analyse en composantes principales et l'analyse des correspondances. Quelle que soit la constitution du tableau de données, toutes les techniques d'analyse factorielle ont un noyau commun que nous désignons sous le nom d'*analyse générale* (section 1.1) et que nous allons présenter maintenant.

## Analyse générale, décomposition aux valeurs singulières

Considérons un tableau de valeurs numériques  $X$  ayant  $n$  lignes et  $p$  colonnes. Pour prendre un exemple, le tableau  $X$  a 1000 lignes et 100 colonnes. Il représente les 100 variables observées sur 1000 individus constituant un échantillon statistique.

Le tableau  $X$  possède donc 100 000 éléments. Pour des raisons diverses, il peut exister des liaisons fonctionnelles ou stochastiques entre certaines variables. Peut-on résumer ces 100 000 données par un nombre inférieur de valeurs sans perte notable d'information compte tenu des liaisons et interrelations entre les valeurs ?

Nous recherchons en fait une technique de réduction s'appliquant de façon systématique à divers types de tableaux et conduisant à une reconstitution rapide mais approximative du tableau de départ.

### 1.1.1 Notions élémentaires et principe d'ajustement

On a vu précédemment comment les lignes et les colonnes d'un tableau rectangulaire permettaient de définir des nuages de points.

La position des points dans le nuage est donnée par l'ensemble des distances entre tous les points et détermine la *forme du nuage*. C'est elle qui caractérise la nature et l'intensité des relations entre les individus (lignes) et entre les variables (colonnes) et révèle les structures de l'information contenues dans les données.

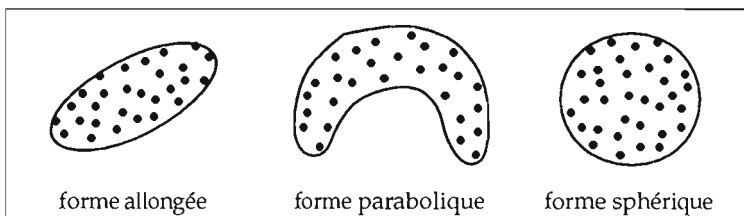


Figure 1.1 - 1  
Différentes formes de nuages

Par exemple, si le nuage de points est uniformément allongé le long d'une droite, il existe un support linéaire dominant pour les points. Une forme parabolique traduira une relation non linéaire tandis qu'un nuage de forme

sphérique marquera plutôt une absence de relation (cf. figure 1.1 - 1). On peut également rencontrer, parmi les formes classiques de nuages, des formes triangulaires ou un nuage composé de quelques amas de points (figure 1.1 - 2).

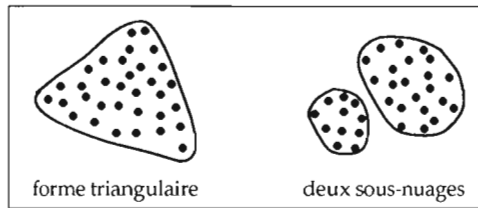


Figure 1.1 - 2  
Autres formes de nuages

Une façon simple de rendre compte visuellement de la forme d'un nuage est de le projeter sur des droites, ou mieux sur des plans, en minimisant les déformations que la projection implique. Pour cela, on peut chercher le sous-espace à une dimension  $H$  qui maximise la somme des carrés des distances entre les projections sur  $H$  de tous les couples de points  $(k, k')$  :

$$\text{Max}_{(H)} \left\{ \sum_k \sum_{k'} d^2(k, k') \right\} \quad [1.1 - 1]$$

Si chaque point est muni d'un masse, c'est la somme pondérée que l'on pourra chercher à maximiser :

$$\text{Max}_{(H)} \left\{ \sum_k \sum_{k'} p_k p_{k'} d^2(k, k') \right\}$$

On calcule ainsi le sous-espace vectoriel qui ajuste au mieux le nuage de points. Nous verrons plus loin, à propos de l'analyse en composantes principales, que ce dernier critère est équivalent au critère ci-dessous (où  $G$  désigne le point moyen ou centre de gravité des projections) :

$$\text{Max}_{(H)} \left\{ \sum_k p_k d^2(k, G) \right\}$$

Toutefois, on ne s'intéresse pas toujours à la forme d'un nuage, mais quelques fois à sa position par rapport à l'origine. Ainsi, en analyse en composantes principales, on s'intéresse bien à la forme du nuage des points-observations dans un espace, mais c'est la position par rapport à l'origine des points-variables qui aura du sens dans l'autre espace.

Le modèle d'analyse par rapport à l'origine désigné ici sous le nom d'analyse générale permet de rendre compte de ces diverses situations. Il n'est qu'une présentation sous forme géométrique de la *décomposition aux valeurs singulières* présentée pour la première fois par Eckart et Young (1936, 1939) pour les tableaux rectangulaires, généralisant les travaux de Sylvester (1889) relatifs aux matrices carrées. Gifi (1990) mentionne

également les travaux antérieurs et indépendants de Beltrami (1873) et Jordan (1874). Cf. également Gower (1966), Gabriel (1971).

Le problème que l'on se propose de résoudre est alors un problème de réduction purement numérique, autrement dit, un problème de compression de données.

Pour exposer cette technique de réduction factorielle, nous nous plaçons successivement dans les espaces vectoriels  $\mathbb{R}^p$  et  $\mathbb{R}^n$ , avec pour notre exemple :  $p = 100$ ,  $n = 1000$ .

### 1.1.2 Ajustement du nuage des individus dans l'espace des variables

On envisage ici le nuage de  $n$  points-individus définis dans l'espace des variables  $\mathbb{R}^p$  et qui sont non pondérés (pour simplifier la formulation). Chacune des  $n$  lignes du tableau  $X$  est considérée comme un vecteur ou encore un point de  $\mathbb{R}^p$ .

Si ce nuage est contenu dans un sous-espace vectoriel à  $q$  dimensions de  $\mathbb{R}^p$  et si  $q$  est notablement inférieur à  $p$ , autrement dit, si le tableau  $X$  est de rang  $q$ , le problème d'approximation est pratiquement résolu<sup>1</sup>.

#### a – Droites d'ajustement

Commençons par chercher un sous-espace vectoriel à *une dimension*, c'est-à-dire une droite passant par l'origine, qui réalise le meilleur ajustement possible du nuage de points.

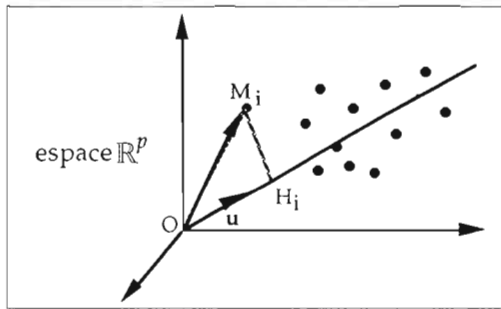


Figure 1.1 - 3  
Meilleur ajustement du nuage de points

<sup>1</sup> Par exemple, si les 1000 points-individus se trouvent dans un sous-espace à 10 dimensions (ou plus généralement si leurs positions sont reconstituées de façon satisfaisante à partir de leurs positions dans ce sous-espace) il suffit, pour retrouver les positions relatives de ces points dans  $\mathbb{R}^p$ , de connaître la nouvelle base (soit 10 vecteurs à 100 dimensions) et les nouvelles coordonnées des points dans cette base (soit 1000 vecteurs à 10 dimensions). On pourrait dans ce cas reconstituer les 100 000 nombres à partir des 11 000 nombres ainsi définis ( $10 \times 100 + 1000 \times 10 = 11\ 000$ ).



Il faut pour cela définir le vecteur directeur unitaire de cette droite. Soit  $\mathbf{u}$  ce vecteur. On désignera également par  $\mathbf{u}$  la matrice colonne associée, et par  $\mathbf{u}'$  sa transposée. On exprime que  $\mathbf{u}$  est unitaire par la relation  $\mathbf{u}'\mathbf{u} = 1$ .

La longueur de la projection  $\text{OH}_i$  d'un vecteur  $\text{OM}_i$  sur le sous-espace à une dimension porté par  $\mathbf{u}$  (figure 1.1 - 3) n'est autre que le produit scalaire de  $\text{OM}_i$  par  $\mathbf{u}$ , somme des produits terme à terme<sup>1</sup> des composantes de  $\text{OM}_i$  et de  $\mathbf{u}$  :

$$\text{OH}_i = \mathbf{x}'_i \mathbf{u} = \sum_j^p x_{ij} u_j$$

Chacune des  $n$  lignes du tableau  $\mathbf{X}$  est un vecteur-individu  $\mathbf{x}_i$  dans  $\mathbb{R}^p$ . Or le produit matriciel  $\mathbf{X}\mathbf{u}$  est la matrice-colonne à  $n$  éléments, dont chaque terme est le produit scalaire d'une ligne de  $\mathbf{X}$  par  $\mathbf{u}$  :

$$\mathbf{X}\mathbf{u} = \begin{bmatrix} x_{11} & \dots & x_{1p} \\ \dots & x_{ij} & \dots \\ \dots & \dots & \dots \\ x_{n1} & \dots & x_{np} \end{bmatrix} \begin{bmatrix} u_1 \\ \dots \\ u_j \\ \dots \\ u_p \end{bmatrix} = \begin{bmatrix} \dots \\ \sum_j x_{ij} u_j \\ \dots \end{bmatrix}$$

Ce sont les  $n$  composantes de la matrice colonne  $\mathbf{X}\mathbf{u}$  qui repèrent sur  $\mathbf{u}$  les  $n$  projections  $\text{OH}_i$  des points du nuage.

Parmi les critères d'ajustement d'un sous-espace à un nuage de  $n$  points, celui que l'on retient et qui conduit aux calculs analytiques sans doute les plus simples, est le critère classique *des moindres carrés*. Il consiste à rechercher la *droite d'allongement maximum* du nuage de points et donc à rendre minimale la somme des carrés des écarts :

$$\sum_{i=1}^n \text{M}_i \text{H}_i^2$$

Le théorème de Pythagore appliqué à chacun des  $n$  triangles rectangles du type  $\text{H}_i \text{OM}_i$  conduit à la relation :

$$\sum_i \text{M}_i \text{H}_i^2 = \sum_i \text{OM}_i^2 - \sum_i \text{OH}_i^2$$

Comme  $\sum_i \text{OM}_i^2$  est une quantité fixe, indépendante du vecteur  $\mathbf{u}$  cherché, il est équivalent de rendre maximale la quantité :

$$\sum_i \text{OH}_i^2$$

<sup>1</sup> On suppose implicitement (et provisoirement) que la métrique dont est muni cet espace est la métrique euclidienne usuelle.

qui s'exprime en fonction de  $X$  et  $u$  par :

$$\sum_i OH_i^2 = (Xu)'Xu = u'X'Xu$$

Pour trouver  $u$ , on est donc conduit à chercher le maximum de la forme quadratique  $u'X'Xu$  :

$$\begin{cases} \text{Max}_{(u)} \{u'X'Xu\} \\ \text{sous la contrainte : } u'u = 1 \end{cases}$$

Soit  $u_1$  le vecteur qui réalise ce maximum. Le sous-espace à deux dimensions s'ajustant au mieux au nuage contient nécessairement le sous-espace engendré par  $u_1$ <sup>1</sup>. On cherche ensuite  $u_2$ , le second vecteur de base de ce sous-espace, orthogonal à  $u_1$  et rendant maximal  $u_2'X'Xu_2$ .

On recherche de façon analogue le meilleur sous-espace au sens des moindres carrés à  $q$  dimensions (pour  $q \leq p$ ).

## b – Caractéristiques du sous-espace d'ajustement

Les démonstrations qui figurent en annexe (§ 1.1.7 ci-après) conduisent à l'énoncé suivant :

"le vecteur unitaire  $u_1$  qui caractérise le sous-espace à une dimension ajustant au mieux le nuage des  $n$  points-individus dans  $\mathbb{R}^p$ , est le *vecteur propre* de la matrice  $X'X$  correspondant à la plus grande *valeur propre*  $\lambda_1$ ".

Plus généralement, le sous-espace à  $q$  dimensions qui ajuste au mieux (au sens des moindres carrés) le nuage dans  $\mathbb{R}^p$  est engendré par les  $q$  premiers vecteurs propres de la matrice symétrique  $X'X$  correspondant aux  $q$  plus grandes valeurs propres. On diagonalisera, par conséquent, la matrice  $X'X$  d'ordre  $(p,p)$ .

L'analyse générale effectue donc une rotation du repère autour de l'origine  $O$  et fournit un système de vecteurs orthonormés dont  $u_1$  puis  $(u_1, u_2), \dots, (u_1, u_2, \dots, u_q, \dots, u_p)$  passent "au plus près" du nuage.

### 1.1.3 Ajustement du nuage des variables dans l'espace des individus

Plaçons-nous maintenant dans l'espace des individus  $\mathbb{R}^n$ , où le tableau  $X$  peut être représenté par un nuage de  $p$  points-variables dont les  $n$  coordonnées représentent les colonnes de  $X$ .

<sup>1</sup> Le raisonnement par l'absurde prouve que s'il ne contenait pas  $u_1$ , il en existerait un meilleur contenant  $u_1$ .

La démarche pour ajuster le nuage des  $p$  points-variables dans cet espace est exactement la même que pour le nuage des points-individus et consiste à rechercher le vecteur unitaire  $\mathbf{v}$ , puis le sous-espace à  $q$  dimensions dans  $\mathbb{R}^n$  qui ajuste au mieux le nuage de points.

Cela conduit à rendre maximale la somme des carrés des  $p$  projections sur  $\mathbf{v}$ , qui sont les  $p$  composantes du vecteur  $\mathbf{X}'\mathbf{v}$ . On maximise la quantité :

$$(\mathbf{X}'\mathbf{v})'\mathbf{X}'\mathbf{v} = \mathbf{v}'\mathbf{X}\mathbf{X}'\mathbf{v} \quad \text{avec la contrainte} \quad \mathbf{v}'\mathbf{v} = 1$$

Comme précédemment, nous sommes amenés à retenir les  $q$  vecteurs propres de  $\mathbf{X}\mathbf{X}'$  correspondant aux  $q$  plus grandes valeurs propres. La matrice à diagonaliser sera cette fois la matrice  $\mathbf{X}\mathbf{X}'$  d'ordre  $(n,n)$ .

On notera  $\mathbf{v}_\alpha$  le vecteur propre de  $\mathbf{X}\mathbf{X}'$  correspondant à la valeur propre  $\mu_\alpha$ .

### 1.1.4 Relation entre les ajustements dans les deux espaces

Recherchons les relations dites de transition entre les deux espaces.

Dans  $\mathbb{R}^p$ , nous avons :

$$\mathbf{X}'\mathbf{X}\mathbf{u}_\alpha = \lambda_\alpha\mathbf{u}_\alpha \quad [1.1 - 2]$$

et dans  $\mathbb{R}^n$  :

$$\mathbf{X}\mathbf{X}'\mathbf{v}_\alpha = \mu_\alpha\mathbf{v}_\alpha \quad [1.1 - 3]$$

En prémultipliant les deux membres de [1.1 - 2] par  $\mathbf{X}$ , on obtient :

$$(\mathbf{X}\mathbf{X}')\mathbf{X}\mathbf{u}_\alpha = \lambda_\alpha(\mathbf{X}\mathbf{u}_\alpha)$$

Cette relation montre qu'à tout vecteur propre  $\mathbf{u}_\alpha$  de  $\mathbf{X}'\mathbf{X}$  relatif à une valeur propre  $\lambda_\alpha$  non nulle, correspond un vecteur propre  $\mathbf{X}\mathbf{u}_\alpha$  de  $\mathbf{X}\mathbf{X}'$ , relatif à la même valeur propre  $\lambda_\alpha$ . Comme on a appelé  $\mu_1$  la plus grande valeur propre de  $\mathbf{X}\mathbf{X}'$ , on a nécessairement  $\lambda_1 \leq \mu_1$ .

En prémultipliant les deux membres de [1.1 - 3] (pour  $\alpha = 1$ ) par  $\mathbf{X}'$ , on voit de même  $\mathbf{X}'\mathbf{v}_1$  est vecteur propre de  $\mathbf{X}'\mathbf{X}$  relativement à la valeur propre  $\mu_1$ , d'où la relation  $\mu_1 \leq \lambda_1$ , ce qui prouve finalement que  $\lambda_1 = \mu_1$ .

On verrait de la même façon que toutes les valeurs propres non nulles des deux matrices  $\mathbf{X}'\mathbf{X}$  et  $\mathbf{X}\mathbf{X}'$  sont égales<sup>1</sup> (avec le même ordre de multiplicité éventuellement) :

$$\lambda_\alpha = \mu_\alpha$$

<sup>1</sup> Il est donc inutile de refaire les calculs de diagonalisation sur  $\mathbf{X}\mathbf{X}'$ , puisqu'une simple transformation linéaire, associée à la matrice  $\mathbf{X}$  de départ, nous permet d'obtenir les directions propres  $\mathbf{X}\mathbf{u}_\alpha$  cherchées dans  $\mathbb{R}^n$ . Il suffit de diagonaliser la matrice  $\mathbf{X}'\mathbf{X}$  ( $p,p$ ) ou  $\mathbf{X}\mathbf{X}'$  ( $n,n$ ) ayant la plus petite dimension.

Remarquons que le vecteur  $Xu_\alpha$  a pour norme  $\lambda_\alpha$  (on a  $u'_\alpha X'Xu_\alpha = \lambda_\alpha$ ) et donc le vecteur  $v_\alpha$  unitaire correspondant à la même valeur propre  $\lambda_\alpha$  est facilement calculable en fonction de  $u_\alpha$ . On obtient ainsi, pour  $\lambda_\alpha \neq 0$ , les formules de transition entre les deux espaces,  $\mathbb{R}^p$  et  $\mathbb{R}^n$  :

$$\begin{cases} v_\alpha = \frac{1}{\sqrt{\lambda_\alpha}} Xu_\alpha \\ u_\alpha = \frac{1}{\sqrt{\lambda_\alpha}} X'v_\alpha \end{cases} \quad [1.1-4]$$

$$\begin{cases} u_\alpha = \frac{1}{\sqrt{\lambda_\alpha}} X'v_\alpha \end{cases} \quad [1.1-5]$$

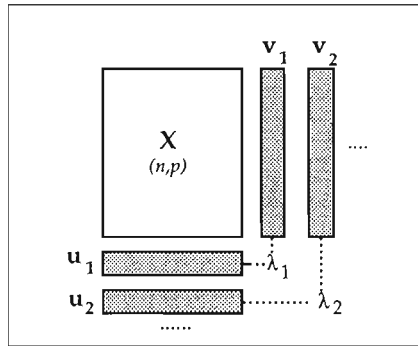


Figure 1.1 - 4  
Relations de transitions

Dans  $\mathbb{R}^p$ ,  $u_\alpha$  est le  $\alpha$ <sup>ième</sup> axe factoriel et l'on calcule le vecteur  $\psi_\alpha$  des coordonnées sur cet axe par :

$$\psi_\alpha = Xu_\alpha$$

De même dans  $\mathbb{R}^n$ ,  $v_\alpha$  est le  $\alpha$ <sup>ième</sup> axe factoriel et l'on construit les coordonnées  $\varphi_\alpha$  par :

$$\varphi_\alpha = X'v_\alpha$$

Compte tenu de [1.1 - 4] et [1.1 - 5], les facteurs peuvent se calculer par :

$$\begin{cases} \psi_\alpha = v_\alpha \sqrt{\lambda_\alpha} \\ \varphi_\alpha = u_\alpha \sqrt{\lambda_\alpha} \end{cases}$$

Sur le sous-espace de  $\mathbb{R}^p$  engendré par  $u_\alpha$  les coordonnées des points du nuage des individus sont les composantes de  $Xu_\alpha$ . Ce sont aussi les composantes de  $v_\alpha \sqrt{\lambda_\alpha}$ .

Les coordonnées des points sur un axe factoriel dans  $\mathbb{R}^p$  sont donc proportionnelles aux composantes de l'axe factoriel dans  $\mathbb{R}^n$  correspondant à la même valeur propre. Il en est de même pour les coordonnées des points du nuage des variables où l'on échangera  $\mathbb{R}^p$  et  $\mathbb{R}^n$ .

**Remarques**

1) L'orientation des axes est arbitraire. En effet, les vecteurs propres sont définis au signe près. La figure 1.1 - 5, concernant trois points, montre que toutes les images, obtenues suivant des orientations différentes des facteurs, respectent la forme du nuage c'est-à-dire les distances entre les points.

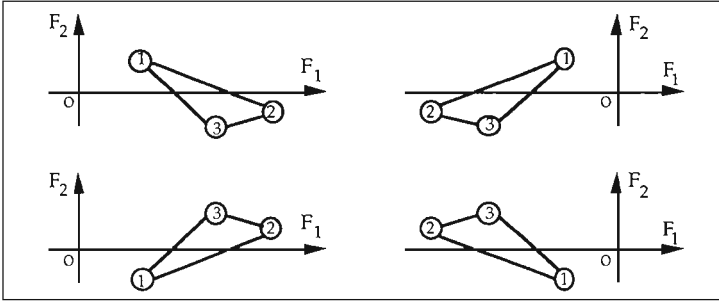


Figure 1.1 - 5  
Orientation arbitraire des axes

2) Les vecteurs de coordonnées dans  $\mathbb{R}^p$  et  $\mathbb{R}^n$  ont pour norme :

$$\Psi'_\alpha \Psi_\alpha = \sum_i^n \varphi_{\alpha i}^2 = \lambda_\alpha$$

et

$$\Phi'_\alpha \Phi_\alpha = \sum_j^p \varphi_{\alpha j}^2 = \lambda_\alpha$$

### 1.1.5 Reconstitution des données de départ

Nous désignons toujours par  $\mathbf{u}_\alpha$  le  $\alpha^{\text{ième}}$  vecteur propre de norme 1 de la matrice  $\mathbf{X}'\mathbf{X}$ , correspondant à la valeur propre  $\lambda_\alpha$ ;  $\mathbf{v}_\alpha$  le  $\alpha^{\text{ième}}$  vecteur propre de norme 1 de  $\mathbf{X}\mathbf{X}'$ . Nous avons :

$$\Psi_\alpha = \mathbf{X}\mathbf{u}_\alpha = \mathbf{v}_\alpha \sqrt{\lambda_\alpha}$$

#### a – Reconstitution exacte

Postmultiplions les deux membres de cette relation par  $\mathbf{u}'_\alpha$  et sommons sur l'ensemble des axes<sup>1</sup> :

$$\mathbf{X} \left\{ \sum_{\alpha=1}^p \mathbf{u}_\alpha \mathbf{u}'_\alpha \right\} = \sum_{\alpha=1}^p \sqrt{\lambda_\alpha} \mathbf{v}_\alpha \mathbf{u}'_\alpha$$

<sup>1</sup> Certains d'entre eux peuvent correspondre à une valeur propre nulle; ils sont alors choisis de façon à compléter la base orthonormée formée par les axes précédents.

Désignons par  $U$  la matrice d'ordre  $(p,p)$  ayant en colonne les vecteurs propres  $u_\alpha$  de  $X'X$ . Ces vecteurs étant orthogonaux et de norme 1, on a :

$$UU' = I \quad \text{et donc} \quad U'U = I$$

où  $I$  est la matrice unité. Mais :

$$\sum_{\alpha=1}^p u_\alpha u'_\alpha = UU'$$

Les valeurs propres  $\lambda_\alpha$  étant toujours rangées par ordre décroissant, la formule précédente devient :

$$X = \sum_{\alpha=1}^p \sqrt{\lambda_\alpha} v_\alpha u'_\alpha \quad [1.1 - 6]$$

et apparaît comme une formule de reconstitution du tableau  $X$ , à partir des  $\lambda_\alpha$  et des vecteurs  $u_\alpha$  et  $v_\alpha$  associés (figure 1.1 - 6).

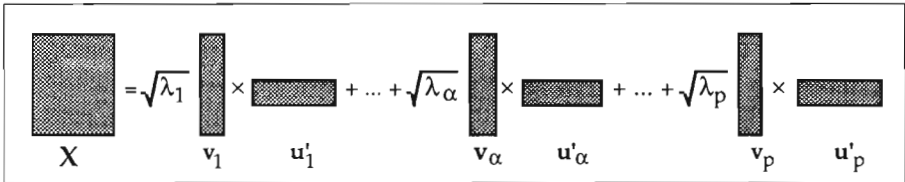


Figure 1.1 - 6  
Reconstitution exacte du tableau de données;  
décomposition aux valeurs singulières.

### Remarque

Les méthodes d'analyse factorielle reposent toutes sur une propriété mathématique des tableaux (ou matrices) rectangulaires : la décomposition aux valeurs singulières [Eckart et Young, 1936]. Cela signifie principalement que, sous des conditions assez générales, une matrice rectangulaire peut être écrite de façon unique comme une "somme optimale" de matrices de rang 1 (produits d'une matrice ligne par une matrice colonne). Que veut-on dire par somme optimale? que la première matrice de rang 1 constitue la meilleure approximation de rang 1 de la matrice initiale (au sens des moindres carrés), que la somme des deux premières constitue la meilleure approximation de rang 2, etc<sup>1</sup>.

### b – Reconstitution approchée

Si les  $p-q$  plus petites valeurs propres sont très faibles et jugées "négligeables", on peut limiter la sommation aux  $q$  premiers termes correspondant aux valeurs propres  $(\lambda_1, \lambda_2, \dots, \lambda_q)$ :

<sup>1</sup> Cette propriété qui concerne le tableau de données lui-même, et non pas seulement la matrice de corrélation ou un tableau de distances construit à partir des données, a ceci de remarquable qu'elle implique de façon similaire les lignes et les colonnes du tableau.

$$\mathbf{X} \approx \mathbf{X}^* = \sum_{\alpha=1}^q \sqrt{\lambda_{\alpha}} \mathbf{v}_{\alpha} \mathbf{u}'_{\alpha} \quad [1.1-7]$$

Si  $q$  est notablement inférieur à  $p$ , on apprécie le gain réalisé en comparant les deux membres de cette relation : le vecteur  $\sqrt{\lambda_{\alpha}} \mathbf{v}_{\alpha}$  a  $n$  composantes et le vecteur  $\mathbf{u}_{\alpha}$  a  $p$  composantes.

Les  $np$  termes de  $\mathbf{X}$  sont donc approchés par des termes construits à partir des  $q(n+p)$  valeurs contenues dans le membre de droite.

### c – Qualité de l'approximation

La qualité de la reconstitution peut être évaluée par la quantité :

$$\tau_q = \frac{\sum_i \sum_j x_{ij}^{*2}}{\sum_i \sum_j x_{ij}^2}$$

On a encore :

$$\tau_q = \frac{\text{tr } \mathbf{X}^* \mathbf{X}^*}{\text{tr } \mathbf{X}' \mathbf{X}}$$

où  $\text{tr}$  désigne l'opérateur trace.

Remplaçant  $\mathbf{X}$  et  $\mathbf{X}^*$  par leurs valeurs tirées de [1.1 - 6] et [1.1 - 7], on obtient immédiatement :

$$\tau_q = \frac{\sum_{\alpha \leq q} \lambda_{\alpha}}{\sum_{\alpha=1}^p \lambda_{\alpha}}$$

Le coefficient  $\tau_q$ , inférieur ou égal à 1, sera appelé *taux d'inertie* ou encore *pourcentage de variance* relatif aux  $q$  premiers facteurs. Son interprétation comme mesure de la qualité numérique de la reconstitution est assez claire, mais nous verrons plus loin que le problème de sa signification statistique est délicat.

## 1.1.6 Diversification de l'analyse générale

La métrique (c'est-à-dire la formule de distance) et le critère d'ajustement (c'est-à-dire la pondération des points) varient suivant le problème et donc suivant la nature des variables.

### a – Analyse générale avec des métriques et des critères quelconques

Jusqu'à présent, nous avons considéré les espaces munis de la métrique  $\mathbf{I}$  (matrice identité) et nous avons supposé que tous les points du nuage avaient la même importance.

Cependant il arrive que l'on ait à travailler avec une métrique plus générale et avec des individus dont les masses sont différentes (pondérations calculées après un redressement d'échantillon, regroupements divers d'individus, etc.). Ces masses vont intervenir dans les calculs de moyennes et lors de l'ajustement des sous-espaces.

Généralisons le principe d'analyse factorielle présenté ci-dessus à des métriques et des critères quelconques.

Plaçons-nous dans l'espace  $\mathbb{R}^p$  et considérons le nuage de  $n$  points-lignes pesants.

Soit  $X$  la matrice d'ordre  $(n,p)$  des coordonnées c'est-à-dire le tableau de données,  $M$  la matrice symétrique définie positive d'ordre  $(p,p)$  définissant la métrique dans  $\mathbb{R}^p$ , et  $N$  la matrice diagonale d'ordre  $(n,n)$  dont les éléments diagonaux sont les masses  $m_i$  des  $n$  points.

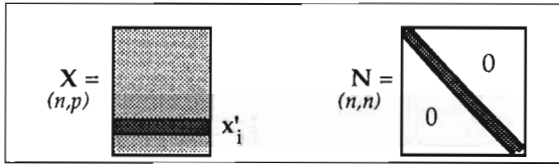


Figure 1.1 - 7  
 $X$ , tableau de coordonnées et  $N$ , matrice diagonale des masses

Un vecteur unitaire  $u$  de  $\mathbb{R}^p$  vérifie maintenant la relation de normalisation  $u'Mu = 1$ .

La coordonnée de la projection  $H_i$  du point  $i$  sur l'axe  $u$  vaut :

$$OH_i = X_i M u$$

et l'ensemble  $F$  des coordonnées des projections sur l'axe  $u$  des  $n$  points-lignes s'exprime par :

$$F = X M U$$

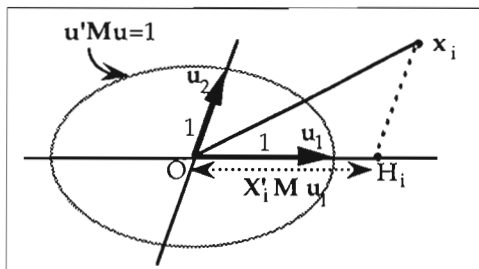


Figure 1.1 - 8  
 Métrique  $M$  dans  $\mathbb{R}^p$



Compte tenu du critère d'ajustement, on veut trouver le vecteur  $\mathbf{u}$  qui rende maximale la somme pondérée des carrés des projections :

$$\text{Max}_{(\mathbf{u})} \left\{ \sum_i m_i OH_i^2 \right\} = \text{Max}_{(\mathbf{u})} \{ \mathbf{u}' \mathbf{M} \mathbf{X}' \mathbf{N} \mathbf{X} \mathbf{M} \mathbf{u} \}$$

sous la contrainte :

$$\mathbf{u}' \mathbf{M} \mathbf{u} = 1$$

Les résultats de l'annexe de cette section nous montrent que  $\mathbf{u}$  est le vecteur propre de la matrice  $\mathbf{A} = \mathbf{X}' \mathbf{N} \mathbf{X} \mathbf{M}$  correspondant à la plus grande valeur propre  $\lambda$ .

L'équation de l'axe factoriel  $\mathbf{u}$  dans  $\mathbb{R}^p$  s'écrit :

$$\mathbf{X}' \mathbf{N} \mathbf{X} \mathbf{M} \mathbf{u} = \lambda \mathbf{u}$$

et les coordonnées factorielles des  $n$  points sont données par la relation :

$$\boldsymbol{\psi} = \mathbf{X} \mathbf{M} \mathbf{u}$$

#### - Relation entre $\mathbb{R}^p$ et $\mathbb{R}^n$

Si les masses et les métriques dans  $\mathbb{R}^p$  ( $\mathbf{N}$  et  $\mathbf{M}$ ) et dans  $\mathbb{R}^n$  ( $\mathbf{P}$ , matrice des masses des  $p$  points-colonnes et  $\mathbf{Q}$ , métrique dans  $\mathbb{R}^n$ ) n'ont pas de relations privilégiées entre elles, on perd les relations de transition et la formule de reconstitution.

En analyse en composantes principales, on utilise la même métrique dans les deux espaces. En analyse des correspondances, on verra que la matrice des masses dans un espace est liée à la métrique de l'autre espace, ce qui permettra de conserver les relations de transition.

#### - Axes d'inertie

La quantité :

$$\mathbf{u}' \mathbf{M} \mathbf{X}' \mathbf{N} \mathbf{X} \mathbf{M} \mathbf{u} = \boldsymbol{\psi}' \mathbf{N} \boldsymbol{\psi} = \sum_i m_i \psi_i^2$$

représente l'inertie du nuage de points pesants le long de l'axe d'allongement maximal, l'axe factoriel  $\mathbf{u}$ . Elle est égale à la valeur propre  $\lambda$  associée au vecteur propre  $\mathbf{u}$ .

Les  $p$  vecteurs propres définissent donc des axes d'inertie du nuage de points et on les obtient par ordre d'inerties décroissantes.

La somme de toutes les valeurs propres donne l'inertie totale du nuage. C'est la trace de la matrice diagonalisée  $\mathbf{A} = \mathbf{X}' \mathbf{N} \mathbf{X} \mathbf{M}$  :

$$\text{Trace}(\mathbf{A}) = \sum_{\alpha=1}^p \lambda_{\alpha}$$

$\mathbf{A}$  est appelée *matrice d'inertie*.

## b – Principe des éléments supplémentaires

L'analyse factorielle permet de trouver des sous-espaces de représentation des proximités entre vecteurs de description d'observations. Elle s'appuie, pour cela, sur des éléments (variables et individus) appelés *éléments actifs*.

Mais elle permet aussi de positionner, dans ce sous-espace, des éléments (points-lignes ou points-colonnes du tableau de données) n'ayant pas participé à l'analyse qui sont appelés *éléments supplémentaires* ou *illustratifs*.

Les éléments supplémentaires interviennent *a posteriori* pour caractériser les axes. Leur introduction dans l'analyse factorielle constitue un apport fondamental car elle permettra de conforter et d'enrichir l'interprétation des facteurs.

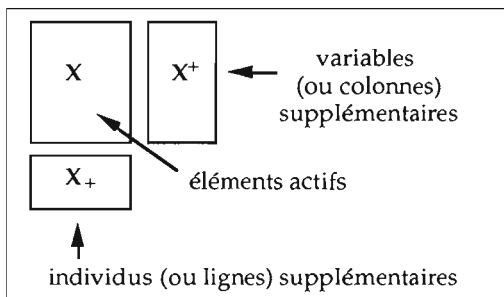


Figure 1.1 - 9  
Représentation des éléments supplémentaires

En effet, il est fréquent, dans la pratique, que l'on dispose d'informations complémentaires élargissant le tableau de données. Ce peut être de nouveaux individus (lignes supplémentaires), par exemple un groupe témoin extérieur à l'échantillon, et il est intéressant alors de positionner ces témoins dans le nuage des individus analysés.

Très souvent dans les applications, ce ne sont pas les individus par eux-mêmes qui sont intéressants mais certaines de leurs caractéristiques connues par ailleurs; on cherchera alors à représenter comme "individus" supplémentaires les centres de gravité des classes d'individus appartenant à une même catégorie. Ce peut être aussi de nouvelles variables (colonnes supplémentaires); on peut disposer d'un ensemble de variables nominales qu'il est intéressant de faire apparaître dans l'analyse réalisée sur des variables continues (et réciproquement). Par ailleurs de nouvelles variables observées sur l'échantillon initial peuvent être disponibles alors qu'on les a volontairement écartées de l'analyse pour ne conserver qu'un corpus homogène de caractéristiques.

Les éléments supplémentaires n'interviennent pas dans les calculs d'ajustement et ne participent donc pas à la formation des axes factoriels. On

cherche uniquement à les positionner dans le nuage des individus ou dans celui des variables en calculant *a posteriori* leurs coordonnées sur les axes factoriels.

Les coordonnées des nouvelles variables sur l'axe  $\alpha$  sont les composantes du vecteur :

$$(X^+)'v_\alpha$$

et les coordonnées des nouveaux individus sur l'axe  $\alpha$  sont :

$$(X_+)u_\alpha$$

Les éléments actifs, définis dans un espace et servant à calculer les plans factoriels, doivent former un ensemble homogène en texture (c'est-à-dire doivent être de même nature, continues ou nominales) pour que les distances entre éléments aient un sens. Mais pour interpréter les similitudes entre ces éléments, ils doivent aussi être homogènes en contenu c'est-à-dire relatifs à un même thème ; on compare les objets selon un certain point de vue et non pas en utilisant sans différenciation tous les attributs connus et souvent disparates. Les variables supplémentaires, quant à elles, ne sont pas soumises à cette condition d'homogénéité.

Cette dichotomie entre variables actives et variables illustratives est analogue à la distinction établie entre les variables explicatives (exogènes) et les variables à expliquer (endogènes) dans les modèles de régression multiple (cf. section 3.2).

D'un point de vue géométrique, nous verrons que les deux situations sont d'ailleurs très similaires. Notons que les points supplémentaires peuvent être considérés comme des points actifs affectés d'une masse nulle.

### c – Autres approches

La décomposition aux valeurs singulières est une propriété de tous les tableaux rectangulaires. Elle fait appel à des distances euclidiennes, c'est-à-dire à des formes quadratiques définies positives, et à des ajustements de sous-espaces vectoriels par minimisation d'un critère lié à ces distances. D'autres approches sont possibles, qui modifient le type de distance, ou la nature des sous-espaces, ou les deux. Il faut s'attendre à perdre beaucoup des propriétés mathématiques simples de l'analyse générale : unicité de la décomposition, symétrie des rôles joués par les lignes et les colonnes, simplicité de la formule de reconstitution, positionnement aisé de variables supplémentaires.

D'autres critères d'ajustements peuvent tout d'abord être utilisés. A la méthode des moindres carrés  $\min\{\sum e_i^2\}$  (norme dite "L<sub>2</sub>"), on peut par exemple substituer celle des moindres valeurs absolues  $\min\{\sum |e_i|\}$  (norme dite "L<sub>1</sub>"). Nous évoquerons à nouveau ces normes à propos de la régression, chapitre 3, § 3.2.1. Sur les méthodes d'analyse des données

utilisant la norme  $L_1$  (dite aussi *city-block distance*) on consultera les contributions et points de vue de Fichet (1987, 1988, ainsi que dans Van Cutsem *et al.*, 1994), Arabie (1991) et le recueil édité par Dodge (1987).

Dans un esprit un peu différent, Meyer (1994) donne un algorithme pour ajuster (au sens des moindres carrés, c'est-à-dire de  $L_2$ ) une matrice de distances de type  $L_p$  à une matrice de dissimilarité donnée.

Pour étudier certaines tables de contingence, notamment les tableaux d'échanges, Domenges et Volle (1979) proposent d'utiliser la distance de Hellinger :  $d^2(x, y) = \sum_i (\sqrt{x_i} - \sqrt{y_i})^2$  ("analyse factorielle sphérique").

Enfin, sans changer la métrique ni le critère d'ajustement, on peut songer à ajuster d'autres surfaces que des hyperplans. Ainsi, dans le cas de l'analyse en composantes principales normée qui est, dans l'espace  $\mathbb{R}^n$ , l'analyse générale de points situés sur une sphère (cf. § 1.2.4), Falissard (1995) propose d'ajuster une hypersphère.

### 1.1.7 Annexe 1 - Démonstration sur les extrema de formes quadratiques sous contraintes quadratiques

Le problème est la recherche du vecteur  $\mathbf{u}$  qui rend maximale la quantité  $\mathbf{u}'\mathbf{A}\mathbf{u}$ , avec la contrainte  $\mathbf{u}'\mathbf{M}\mathbf{u} = 1$ , expression où  $\mathbf{A}$  et  $\mathbf{M}$  sont des matrices symétriques;  $\mathbf{M}$  est de plus définie non-négative et définit la métrique dans  $\mathbb{R}^p$ .

On donnera deux démonstrations élémentaires pour la solution de ce problème. L'une fait appel aux multiplicateurs de Lagrange (calcul classique d'extremum sous contrainte), l'autre suppose connues certaines propriétés spectrales des matrices symétriques<sup>1</sup>.

#### - Démonstration directe

La forme quadratique  $\mathbf{u}'\mathbf{A}\mathbf{u}$  s'écrit :

$$\mathbf{u}'\mathbf{A}\mathbf{u} = \sum_{ij} a_{ij}u_iu_j$$

En dérivant cette quantité successivement par rapport aux  $p$  composantes du vecteur  $\mathbf{u}$ , on voit que le vecteur des dérivées partielles de  $\mathbf{u}'\mathbf{A}\mathbf{u}$  s'écrit sous forme matricielle :

$$\frac{\partial(\mathbf{u}'\mathbf{A}\mathbf{u})}{\partial\mathbf{u}} = 2\mathbf{A}\mathbf{u}$$

---

<sup>1</sup> Le problème est ici un peu plus général que celui rencontré précédemment, pour lequel  $\mathbf{A} = \mathbf{X}'\mathbf{X}$  et  $\mathbf{M} = \mathbf{I}$  où  $\mathbf{I}$  est la matrice unité. Mais cette formulation plus large, avec une métrique et des critères quelconques tels que des masses affectées aux points, sera utile à propos de l'analyse des correspondances et de l'analyse discriminante. Elle n'introduit guère de difficulté supplémentaire au niveau des démonstrations.

De même :

$$\frac{\partial(\mathbf{u}'\mathbf{M}\mathbf{u})}{\partial\mathbf{u}} = 2\mathbf{M}\mathbf{u}$$

La recherche d'un maximum lié implique que s'annulent les dérivées du Lagrangien :

$$\mathcal{L} = \mathbf{u}'\mathbf{A}\mathbf{u} - \lambda(\mathbf{u}'\mathbf{M}\mathbf{u} - 1)$$

$\lambda$  étant un multiplicateur de Lagrange. Par suite :

$$\frac{\partial\mathcal{L}}{\partial\mathbf{u}} = 2\mathbf{A}\mathbf{u} - 2\lambda\mathbf{M}\mathbf{u} = 0$$

exprime la condition d'extremum. On en déduit la relation :

$$\mathbf{A}\mathbf{u} = \lambda\mathbf{M}\mathbf{u} \quad [1.1 - 8]$$

Prémultipliant les deux membres de cette relation par  $\mathbf{u}'$ , et tenant compte du fait que  $\mathbf{u}'\mathbf{M}\mathbf{u} = 1$ , il vient :

$$\lambda = \mathbf{u}'\mathbf{A}\mathbf{u}$$

La valeur du paramètre  $\lambda$  est donc le maximum cherché.

Lorsque la matrice  $\mathbf{M}$  est définie positive, donc inversible, la relation [1.1 - 8] s'écrit alors :

$$\mathbf{M}^{-1}\mathbf{A}\mathbf{u} = \lambda\mathbf{u}$$

$\mathbf{u}$  est le *vecteur propre* de la matrice  $\mathbf{M}^{-1}\mathbf{A}$  correspondant à la plus grande *valeur propre*  $\lambda$  (si celle-ci est unique, ce qui sera le cas général).

Appelons désormais  $\mathbf{u}_1$ , le vecteur  $\mathbf{u}$  correspondant à la plus grande valeur  $\lambda_1$  telle que la relation [1.1 - 8] soit vérifiée. Cherchons le vecteur  $\mathbf{u}_2$ , unitaire et  $\mathbf{M}$ -orthogonal à  $\mathbf{u}_1$  (c'est-à-dire tel que  $\mathbf{u}_2'\mathbf{M}\mathbf{u}_2 = 1$  et  $\mathbf{u}_1'\mathbf{M}\mathbf{u}_2 = 0$ ), qui rend maximale la forme quadratique  $\mathbf{u}_2'\mathbf{A}\mathbf{u}_2$ .

On est conduit à annuler les dérivées du Lagrangien :

$$\mathcal{L} = \mathbf{u}_2'\mathbf{A}\mathbf{u}_2 - \lambda_2(\mathbf{u}_2'\mathbf{M}\mathbf{u}_2 - 1) - \mu_2\mathbf{u}_2'\mathbf{M}\mathbf{u}_1$$

où  $\lambda_2$  et  $\mu_2$  sont deux multiplicateurs de Lagrange.

La condition d'extremum s'écrit pour  $\mathbf{u}_2$  :

$$\frac{\partial\mathcal{L}}{\partial\mathbf{u}_2} = 2\mathbf{A}\mathbf{u}_2 - 2\lambda_2\mathbf{M}\mathbf{u}_2 - \mu_2\mathbf{M}\mathbf{u}_1 = 0$$

En multipliant les divers membres de cette relation par  $\mathbf{u}_1'$ , on voit que  $\mu_2 = 0$  (puisque  $\mathbf{u}_1'\mathbf{A}\mathbf{u}_2 = \lambda_1\mathbf{u}_1'\mathbf{M}\mathbf{u}_2 = 0$ ).

Il reste donc comme précédemment :

$$\mathbf{A}\mathbf{u}_2 = \lambda_2\mathbf{M}\mathbf{u}_2$$

Quand  $\mathbf{M}$  est inversible,  $\mathbf{u}_2$  est le second vecteur propre de  $\mathbf{M}^{-1}\mathbf{A}$ , relatif à la seconde plus grande valeur propre  $\lambda_2$  si celle-ci est unique.

La démonstration s'étend aisément au cas d'un vecteur unitaire  $\mathbf{u}_\alpha$  pour  $\alpha \leq p$  (i.e. :  $\mathbf{u}'_\alpha \mathbf{M} \mathbf{u}_\alpha = 1$ ),  $\mathbf{M}$ -orthogonal aux vecteurs  $\mathbf{u}_\beta$  trouvés précédemment ( $\mathbf{u}'_\alpha \mathbf{M} \mathbf{u}_\beta = 0$  pour  $\beta < \alpha$ ) et rendant maximale la forme  $\mathbf{u}'_\alpha \mathbf{A} \mathbf{u}_\alpha$ . On a alors :

$$\mathbf{A} \mathbf{u}_\alpha = \lambda_\alpha \mathbf{M} \mathbf{u}_\alpha$$

Et si  $\mathbf{M}$  est inversible :

$$\mathbf{M}^{-1}\mathbf{A} \mathbf{u}_\alpha = \lambda_\alpha \mathbf{u}_\alpha$$

### - Seconde démonstration

Nous ne ferons qu'esquisser cette démonstration, dans le cas où  $\mathbf{M}$  est définie positive. On peut alors décomposer cette matrice sous la forme classique  $\mathbf{M} = \mathbf{L}'\mathbf{L}$ , où  $\mathbf{L}$  est inversible puisque  $\mathbf{M}$  est supposée définie positive.

Posant alors  $\mathbf{u} = \mathbf{L}^{-1}\mathbf{y}$ , la contrainte de normalisation  $\mathbf{u}'\mathbf{M}\mathbf{u} = 1$  s'écrit maintenant  $\mathbf{y}'\mathbf{y} = 1$ , et la quantité à rendre maximale  $\mathbf{u}'\mathbf{A}\mathbf{u}$  devient  $\mathbf{y}'\mathbf{S}\mathbf{y}$ , avec  $\mathbf{S} = \mathbf{L}^{-1}\mathbf{A}\mathbf{L}^{-1}$ .

Soit  $\mathbf{T}$  la matrice orthogonale ( $p,p$ ) dont les colonnes sont les vecteurs propres  $\mathbf{t}_\alpha$  de  $\mathbf{S}$ , normés et ordonnés suivant les valeurs propres  $\lambda_\alpha$  décroissantes, et soit  $\Lambda$  la matrice diagonale dont le  $\alpha^{\text{ième}}$  élément vaut  $\lambda_\alpha$ .

Posons encore  $\mathbf{z} = \mathbf{T}\mathbf{y}$  (ce qui implique  $\mathbf{y} = \mathbf{T}\mathbf{z}$  car  $\mathbf{T}' = \mathbf{T}^{-1}$ ). On a alors :

$$\mathbf{y}'\mathbf{S}\mathbf{y} = \mathbf{y}'\mathbf{T}\Lambda\mathbf{T}'\mathbf{y} = \mathbf{z}'\Lambda\mathbf{z}$$

avec la contrainte  $\mathbf{z}'\mathbf{z} = 1$ .

La solution est alors proche. On remarque que  $\lambda_1 \geq \mathbf{z}'\Lambda\mathbf{z}$ ; en effet :

$$\lambda_1 - \mathbf{z}'\Lambda\mathbf{z} = \mathbf{z}'(\lambda_1\mathbf{I} - \Lambda)\mathbf{z} \geq 0$$

Le maximum  $\lambda_1$  est effectivement atteint pour  $\mathbf{z}' = (1,0,0,0,\dots,0)$ , donc pour  $\mathbf{y} = \mathbf{t}_1$  et pour  $\mathbf{u}_1 = \mathbf{L}^{-1}\mathbf{t}_1$ . De la relation  $\mathbf{S}\mathbf{t}_1 = \lambda_1\mathbf{t}_1$ , on tire :

$$\mathbf{L}^{-1}\mathbf{A}\mathbf{L}^{-1}\mathbf{t}_1 = \lambda_1\mathbf{t}_1$$

D'où, finalement<sup>1</sup> :

$$\mathbf{M}^{-1}\mathbf{A}\mathbf{u}_1 = \lambda_1\mathbf{u}_1$$

<sup>1</sup> On note au passage qu'il suffit ici de procéder à la diagonalisation d'une matrice symétrique  $\mathbf{S}$  (après avoir décomposé  $\mathbf{M}$  sous la forme :  $\mathbf{M} = \mathbf{L}'\mathbf{L}$ ), alors que la matrice précédente  $\mathbf{M}^{-1}\mathbf{A}$  est en général non-symétrique. Cette propriété est utilisée dans les programmes de calcul (en particulier en analyse des correspondances), car la recherche des éléments spectraux est notablement plus rapide et fiable dans le cas des matrices symétriques.

## Section 1.2

---

# Analyse en Composantes Principales

Conçue pour la première fois par Karl Pearson en 1901, intégrée à la statistique mathématique par Harold Hotelling en 1933, l'analyse en composantes principales n'est vraiment utilisée que depuis l'avènement et la diffusion des moyens de calculs actuels.

La technique d'analyse en composantes principales peut être présentée de divers points de vue. Pour le statisticien classique, il s'agit de la recherche des axes principaux de l'ellipsoïde indicateur d'une distribution normale multidimensionnelle, ces axes étant estimés à partir d'un échantillon. C'est la présentation initiale de Hotelling (1933), puis celle des manuels classiques d'analyse multivariée, comme l'ouvrage fondamental d'Anderson (1958).

Pour les factorialistes classiques, il s'agit d'un cas particulier de la méthode d'analyse factorielle des psychométriciens (cas de variances spécifiques nulles ou égales ; cf. Horst, 1965; Harman, 1967 ; cf. également § 3.2.9).

Enfin, du point de vue plus récent des analystes de données, il s'agit d'une technique de représentation des données, ayant un caractère optimal selon certains critères algébriques et géométriques et que l'on utilise en général sans référence à des hypothèses de nature statistique ni à un modèle particulier. Ce point de vue, fort répandu actuellement est peut-être le plus ancien. C'est celui qui avait été adopté par Pearson (1901). Bien entendu, il ne s'agissait pas de l'analyse en composantes principales telle que nous la présentons, mais les idées essentielles de la méthode étaient déjà entrevues par cet auteur. On trouvera une présentation plus proche de nos préoccupations dans l'article de synthèse de Rao (1964).

L'analyse en composantes principales présente de nombreuses variantes selon les transformations apportées au tableau de données : le nuage des points-individus peut être centré ou non, réduit ou non. Parmi ces variantes, l'analyse en composantes principales normée (nuage centré-réduit) est certainement la plus utilisée et c'est celle-ci que nous choisirons pour présenter les principes de l'analyse en composantes principales.

### 1.2.1 Domaine d'application

L'utilisateur éventuel de l'analyse en composantes principales se trouve dans la situation suivante : il possède un tableau rectangulaire de mesures, dont les colonnes figurent des variables à valeurs numériques continues (des mensurations, des taux, etc.) et dont les lignes représentent les individus sur lesquels ces variables sont mesurées.

En biométrie, il est fréquent de procéder à de nombreuses mensurations sur certains organes ou certains animaux. En micro-économie, on aura par exemple à relever les dépenses des ménages en divers postes.

D'une manière générale, la condition que doivent remplir ces tableaux numériques pour être l'objet d'une description par l'analyse en composantes principales est la suivante : l'une au moins des dimensions du tableau (les lignes en général) est formée d'unités ayant un caractère répétitif, l'autre pouvant être éventuellement plus hétérogène.

Dans les exemples cités, les lignes ont ce caractère répétitif : on les désignera en général sous le nom d'individus ou d'observations, les colonnes étant désignées sous le nom de variables. Quelquefois, ces lignes pourront être considérées comme des réalisations indépendantes de vecteurs aléatoires, dont les composantes correspondent aux différentes variables.

Pour fixer les idées, nous considérons le tableau  $\mathbf{R}$  des mesures prises sur quelques milliers d'hommes actifs concernant leurs temps d'activités quotidiennes. On dispose de 16 variables décrivant des temps d'activités, en minutes par jour (sommeil, repos, repas chez soi, etc.). Les personnes enquêtées sont regroupées en 27 groupes selon l'âge, le niveau d'éducation et le type d'agglomération. Ce sont ces groupes qui sont observés et sont ici considérés comme des "individus" (cf. tableau 1.2 - 1, au § 1.2 - 11). Il s'agit de disposer d'un tableau de dimensions raisonnables dans le cadre d'un exposé pédagogique, et non pas d'un exemple ayant une portée méthodologique générale, une des attitudes de base en analyse descriptive des données étant au contraire "*de ne pas réduire a priori le champ de l'observable*".

Le tableau  $\mathbf{R}$  aura en colonne les 16 mesures caractérisant les 27 observations. Le terme général  $r_{ij}$  de ce tableau décrit la durée moyenne de l'activité  $j$  de l'observation  $i$  (constituant un groupe d'individus).

Nous voulons avoir une idée de la structure de l'ensemble des 16 activités, ainsi que des similitudes éventuelles de comportement entre les groupes d'individus retenus.

### 1.2.2 Interprétations géométriques

Les représentations géométriques entre les lignes et entre les colonnes du tableau de données permettent de représenter visuellement les proximités entre les individus et entre les variables.

#### a – Pour les $n$ individus

Dans  $\mathbb{R}^p$ , les  $n(n-1)$  distances attachées aux couples de points qui représentent des individus ont une interprétation directe pour l'utilisateur :

$$d^2(i, i') = \sum_{j=1}^p (r_{ij} - r_{i'j})^2 \quad [1.2 - 1]$$



Il s'agit ici de la distance euclidienne classique. Deux points sont très voisins si, dans l'ensemble, leurs  $p$  coordonnées sont très proches. Les deux individus concernés sont alors caractérisés par des valeurs presque égales pour chaque variable. Dans l'exemple évoqué ci-dessus, deux individus représentés par des points proches consacrent les mêmes temps aux mêmes activités.

## b – Pour les $p$ variables

Si les valeurs prises par deux variables particulières sont très voisines pour tous les individus, ces variables seront représentées par deux points très proches dans  $\mathbb{R}^n$ . Cela peut vouloir dire que ces variables mesurent une même chose ou encore qu'elles sont liées par une relation particulière.

Toutefois la définition de ces proximités dans les deux espaces est assez fruste. Des problèmes d'échelle de mesure se posent d'emblée : le temps consacré au sommeil est toujours beaucoup plus important que le temps passé à la lecture.

Par ailleurs, dans un cadre plus général, comment calculer la distance entre deux variables si l'une est exprimée en centimètre et l'autre en kilogramme? Comment interpréter un éloignement moyen dans  $\mathbb{R}^p$ ? Est-ce que deux individus assez proches dans  $\mathbb{R}^p$  ont des valeurs assez voisines pour chacune des variables, ou au contraire très proches pour certaines et éloignées pour d'autres?

L'analyse en composantes principales normée permet de donner des éléments de réponses à ces questions.

### 1.2.3 Analyse du nuage des individus

Nous considérons tout d'abord ici le nuage des  $n$  individus non pondérés. Nous voulons, dans l'espace des variables, ajuster le nuage de  $n$  points par un sous-espace à une, puis deux dimensions, de façon à obtenir sur un graphique une représentation visuelle la plus fidèle possible des proximités existant entre les  $n$  individus vis-à-vis des  $p$  variables.

#### a – Principe d'ajustement

Ce n'est donc plus la somme des carrés des distances à l'origine en projection qu'il faut rendre maximum (cf. formule [1.1 - 1]), mais la somme des carrés des distances entre *tous les couples d'individus* :

$$\underset{(H)}{\text{Max}} \left\{ \sum_i^n \sum_{i'}^n d_H^2(i, i') \right\}$$

Autrement dit, la droite d'ajustement  $H_1$  ne doit pas être astreinte à passer par l'origine, comme  $H_0$  dans l'analyse générale (figure 1.2 - 1).

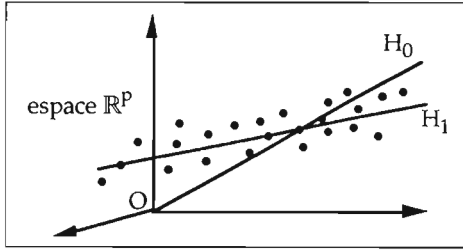


Figure 1.2 - 1  
Droite d'ajustement du nuage de  $n$  points

Si  $h_i$  et  $h_{i'}$  désignent les valeurs des projections de deux points-individus  $i$  et  $i'$  sur  $H_1$ , on a la relation classique :

$$\begin{aligned} \sum_{i,i'}^n d^2(i,i') &= \sum_{i,i'}^n (h_i - h_{i'})^2 = n \sum_{i,i'}^n h_i^2 + n \sum_{i,i'}^n h_{i'}^2 - 2 \sum_i^n h_i \sum_{i'}^n h_{i'} \\ &= 2n^2 \left( \frac{1}{n} \sum_i^n h_i^2 - \bar{h}^2 \right) = 2n \sum_i^n (h_i - \bar{h})^2 \end{aligned}$$

où  $\bar{h}$  désigne la moyenne des projections des  $n$  individus :

$$\bar{h} = \frac{1}{n} \sum_i^n h_i$$

et correspond à la projection sur  $H_1$  du centre de gravité  $G$  du nuage dont la  $j^{\text{ème}}$  coordonnée vaut :

$$\bar{r}_j = \frac{1}{n} \sum_i^n r_{ij}$$

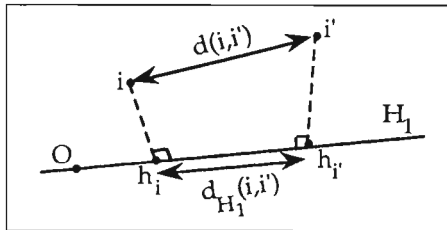


Figure 1.2 - 2  
Projections sur  $H_1$

Par conséquent, on a :

$$\sum_{i,i'}^n d^2(i,i') = 2n \sum_i^n d^2(i,G)$$

Rendre maximum la somme des carrés des distances entre tous les couples d'individus revient à maximiser la somme des carrés des distances entre les points et le centre de gravité du nuage G :

$$\text{Max}_{(H)} \left\{ \sum_{i,i'}^n d_H^2(i, i') \right\}$$

est équivalent à :

$$\text{Max}_{(H)} \left\{ \sum_i^n d_H^2(i, G) \right\}$$

Si l'origine est prise en G, la quantité à maximiser sera à nouveau la somme des carrés des distances à l'origine, ce qui correspond au problème de l'analyse générale dans  $\mathbb{R}^p$  (cf. § 1.1.2).

Le sous-espace cherché résulte de l'analyse générale du tableau transformé X, de terme général :

$$x_{ij} = r_{ij} - \bar{r}_j$$

## b – Distance entre individus

La distance entre deux individus  $i$  et  $i'$  est la distance euclidienne usuelle donnée par la formule [1.2 - 1].

Il peut exister des valeurs de  $j$  pour lesquelles les variables correspondantes sont d'échelles très diverses, (exemple : temps passé au sommeil, temps passé à la lecture) ; on veut que la distance entre deux points soit indépendante des unités sur les variables. On peut parfois désirer, surtout lorsque les unités de mesures ne sont pas les mêmes, faire jouer à chaque variable un rôle identique dans la définition des proximités entre individus : on parle alors d'analyse en composantes principales normée. Pour cela on corrige les échelles en adoptant la distance :

$$d^2(i, i') = \sum_{j=1}^p \left( \frac{r_{ij} - r_{i'j}}{s_j \sqrt{n}} \right)^2$$

$s_j$  désignant l'écart-type de la variable  $j$  :

$$s_j^2 = \frac{1}{n} \sum_{i=1}^n (r_{ij} - \bar{r}_j)^2$$

Finalement, nous retiendrons que l'analyse normée dans  $\mathbb{R}^p$  du tableau brut R est l'analyse générale de X, de terme général :

$$x_{ij} = \frac{r_{ij} - \bar{r}_j}{s_j \sqrt{n}} \quad [1.2 - 2]$$

Toutes les variables ainsi transformées sont "comparables" et ont même dispersion :

$$s^2(x_j) = 1$$

Les variables sont *centrées réduites*. On mesure l'écart à la moyenne en nombre d'écart-types de la variable  $j$ .

### c – Matrice à diagonaliser

En résumé, l'analyse du nuage des points-individus dans  $\mathbb{R}^p$  nous a amené à effectuer une translation de l'origine au centre de gravité de ce nuage et à changer, dans le cas de l'analyse normée, les échelles sur les différents axes.

L'analyse du tableau transformé  $\mathbf{X}$  nous conduit à diagonaliser la matrice  $\mathbf{C} = \mathbf{X}'\mathbf{X}$ .

Le terme général  $c_{jj'}$  de cette matrice s'écrit :

$$c_{jj'} = \sum_i^n x_{ij}x_{ij'}$$

soit :

$$c_{jj'} = \frac{1}{n} \sum_i^n \frac{(r_{ij} - \bar{r}_j)(r_{ij'} - \bar{r}_{j'})}{s_j s_{j'}}$$

c'est-à-dire :

$$c_{jj'} = \text{cor}(j, j')$$

$c_{jj'}$  n'est autre que le coefficient de corrélation empirique entre les variables  $j$  et  $j'$  (d'où l'utilité du coefficient  $\sqrt{n}$  introduit au dénominateur de la relation [1.2 - 2]).

La matrice à diagonaliser est donc la *matrice de corrélations*.

### d – Axes factoriels

Les coordonnées des  $n$  points-individus sur l'axe factoriel  $\mathbf{u}_\alpha$  ( $\alpha^{\text{ième}}$  vecteur propre de la matrice  $\mathbf{C}$  associé à la valeur propre  $\lambda_\alpha$ ) sont les  $n$  composantes du vecteur :

$$\Psi_\alpha = \mathbf{X}\mathbf{u}_\alpha$$

Le facteur  $\Psi_\alpha$  est une combinaison linéaire des variables initiales.

Puisque le nuage des individus est centré sur le centre de gravité, la moyenne du facteur est nulle :

$$\sum_i^n \psi_{\alpha i} = 0$$

et sa variance vaut :

$$\text{var}(\Psi_\alpha) = \lambda_\alpha$$

La coordonnée du point-individu  $i$  sur cet axe s'écrit explicitement :

$$\psi_{\alpha i} = \sum_{j=1}^p u_{\alpha j} x_{ij} = \sum_{j=1}^p u_{\alpha j} \frac{r_{ij} - \bar{r}_j}{s_j \sqrt{n}}$$

### 1.2.4 Analyse du nuage des points-variables

L'analyse générale développée dans la section précédente nous a montré qu'en effectuant un ajustement dans un espace, on effectuait implicitement un ajustement dans l'autre espace. Nous avons volontairement choisi de commencer en travaillant dans  $\mathbb{R}^p$ . Dans cet espace, la transformation du tableau  $\mathbf{R}$  initial selon la relation [1.2 - 2] avait deux objectifs :

- d'une part obtenir un ajustement qui respecte dans la mesure du possible les distances entre points-individus ;
- d'autre part, faire jouer des rôles similaires à toutes les variables dans la définition des distances entre individus.

Notons que la formule [1.2 - 2] ne fait pas intervenir de façon symétrique les lignes et les colonnes du tableau initial  $\mathbf{R}$ .

Que signifie, dans  $\mathbb{R}^n$ , la proximité entre deux points-variables  $j$  et  $j'$  si l'on prend comme coordonnées de ces variables les colonnes du tableau transformé  $\mathbf{X}$  ?

#### a – distances entre points-variables

La distance entre variables découle de l'analyse dans  $\mathbb{R}^p$ . Calculons la distance euclidienne usuelle entre deux variables  $j$  et  $j'$  :

$$d^2(j, j') = \sum_{i=1}^n (x_{ij} - x_{ij'})^2$$

soit :

$$d^2(j, j') = \sum_{i=1}^n x_{ij}^2 + \sum_{i=1}^n x_{ij'}^2 - 2 \sum_{i=1}^n x_{ij} x_{ij'}$$

Remplaçant  $x_{ij}$  par sa valeur tirée de [1.2 - 2] et tenant compte du fait que :

$$s_j^2 = \frac{1}{n} \sum_{i=1}^n (r_{ij} - \bar{r}_j)^2$$

on obtient :  $\sum_{i=1}^n x_{ij}^2 = \sum_{i=1}^n x_{ij'}^2 = 1$  et également :  $\sum_{i=1}^n x_{ij} x_{ij'} = c_{jj'}$

D'où la relation liant la distance dans  $\mathbb{R}^n$  entre deux points-variables  $j$  et  $j'$  et le coefficient de corrélation  $c_{jj'}$  entre ces variables :

$$d^2(j, j') = 2(1 - c_{jj'}) \quad [1.2 - 3]$$

ce qui implique :

$$0 \leq d^2(j, j') \leq 4$$

Dans l'espace  $\mathbb{R}^n$ , le cosinus de l'angle de deux vecteurs-variables est le coefficient de corrélation entre ces deux variables ( $c_{jj'} = \cos(j, j')$ ). Si ces deux

variables sont à la distance 1 de l'origine (i.e. si elles sont de variance unité), le cosinus n'est autre que leur produit scalaire.

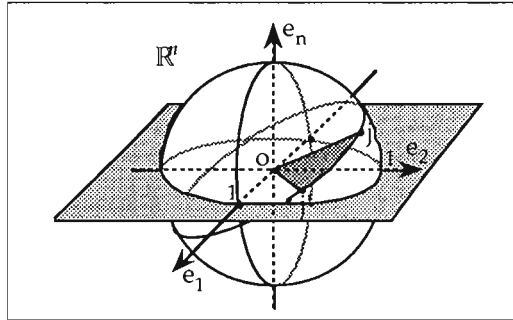


Figure 1.2 - 3  
Système de proximités entre deux points-variables

Le système de proximités entre points-variables induit par la relation [1.2 - 3] est familier au statisticien :

- Deux variables fortement corrélées sont très proches l'une de l'autre ( $c_{jj'} = 1$ ) ou au contraire les plus éloignées possible ( $c_{jj'} = -1$ ) selon que la relation linéaire qui les lie est directe ou inverse :
- Deux variables orthogonales ( $c_{jj'} = 0$ ) sont à distance moyenne.

Les proximités entre points-variables s'interprètent donc en termes de corrélations.

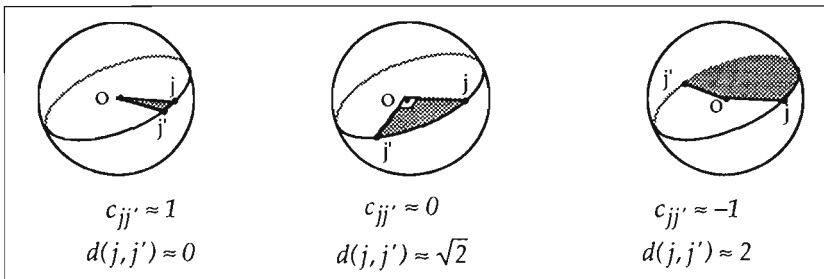


Figure 1.2 - 4  
Corrélations et distances entre points-variables

### b – Distance à l'origine

L'analyse dans  $\mathbb{R}^n$  ne se fait pas par rapport au centre de gravité du nuage de points-variables, contrairement à celui des points-individus, mais par rapport à l'origine.

La distance d'une variable  $j$  à l'origine  $O$  s'exprime par :

$$d^2(O, j) = \sum_{i=1}^n x_{ij}^2 = 1$$

Tous les points-variables sont sur une sphère de rayon 1 centrée à l'origine des axes, la *sphère des corrélations*.

Les plans d'ajustement couperont la sphère suivant de grands cercles (de rayon 1), les *cercles des corrélations*, à l'intérieur desquels se trouveront les points-variables.

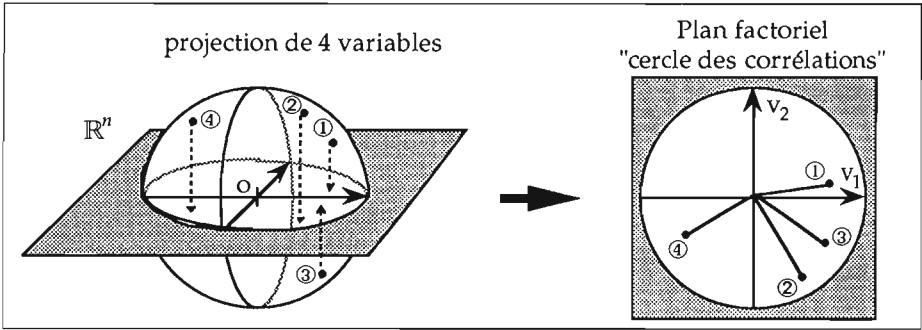


Figure 1.2 - 5  
Représentation de la sphère et du cercle des corrélations

**Remarque**

La transformation analytique simple [1.2 - 2] a dans les espaces  $\mathbb{R}^p$  et  $\mathbb{R}^n$  des interprétations géométriques différentes. Considérons par exemple l'opération de centrage des variables  $x_{ij} \rightarrow (x_{ij} - \bar{x}_j)$  :

- Dans  $\mathbb{R}^p$ , cette transformation équivaut à une translation de l'origine des axes au centre de gravité (ou point moyen) du nuage (cf. figure 1.2 - 6).
- Dans  $\mathbb{R}^n$ , cette transformation est une projection parallèlement à la première bissectrice des axes sur l'hyperplan qui lui est orthogonal<sup>1</sup> (cf. figure 1.2 - 7).

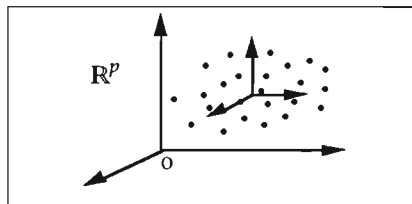


Figure 1.2 - 6  
Transformation dans  $\mathbb{R}^p$

<sup>1</sup> La matrice  $P$  d'ordre  $(n, n)$  associée à cette transformation a pour terme général  $p_{ii'} = \delta_{ii'} - \frac{1}{n}$  où  $\delta_{ii'} = 1$  si  $i = i'$ , et 0 sinon.  $P$  est idempotente :  $P^2 = P$ .

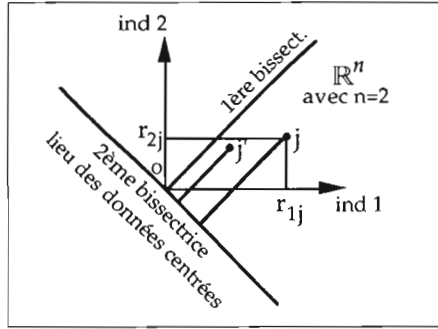


Figure 1.2 - 7  
Transformation dans  $\mathbb{R}^n$

### c – Axes factoriels ou composantes principales

Nous avons vu dans l'analyse générale (§ 1.1.4.) qu'il est inutile de procéder à la diagonalisation de la matrice  $\mathbf{X}\mathbf{X}'$  d'ordre  $(n,n)$  une fois connus les vecteurs propres  $\mathbf{u}_\alpha$  et les valeurs propres  $\lambda_\alpha$  de la matrice  $\mathbf{C} = \mathbf{X}'\mathbf{X}$ .

Le vecteur  $\mathbf{v}_\alpha = \frac{1}{\sqrt{\lambda_\alpha}} \mathbf{X}\mathbf{u}_\alpha$  est en effet un vecteur propre unitaire de  $\mathbf{X}\mathbf{X}'$ , relativement à la même valeur propre  $\lambda_\alpha$ . Le  $\alpha^{\text{ième}}$  facteur dans  $\mathbb{R}^n$  s'écrit :

$$\varphi_\alpha = \mathbf{X}' \mathbf{v}_\alpha = \frac{1}{\sqrt{\lambda_\alpha}} \mathbf{X}' \mathbf{X} \mathbf{u}_\alpha = \mathbf{u}_\alpha \sqrt{\lambda_\alpha}$$

comme  $\psi_\alpha = \mathbf{X}\mathbf{u}_\alpha$ , on a :

$$\varphi_\alpha = \frac{1}{\sqrt{\lambda_\alpha}} \mathbf{X}' \psi_\alpha$$

alors les coordonnées factorielles  $\varphi_{\alpha j}$  des points-variables sur l'axe  $\alpha$  sont les composantes de  $\mathbf{X}'\mathbf{v}_\alpha$  soit encore<sup>1</sup> de  $\mathbf{u}_\alpha \sqrt{\lambda_\alpha}$  :

$$\varphi_{\alpha j} = \sum_{i=1}^n \left( \frac{r_{ij} - \bar{r}_j}{s_j \sqrt{n}} \right) \left( \frac{\psi_{\alpha i}}{\sqrt{\lambda_\alpha}} \right)$$

et l'on a :

$$\varphi_{\alpha j} = \text{cor}(j, \psi_\alpha) \quad [1.2 - 4]$$

La coordonnée d'un point-variable sur un axe n'est autre que le *coefficient de corrélation* de cette variable avec le facteur  $\psi_\alpha$  (combinaison linéaire des variables initiales) considéré lui-même comme variable artificielle dont les coordonnées sont constituées par les  $n$  projections des individus sur cet axe.

<sup>1</sup> Ce sont en quelque sorte des sous-produits des calculs déjà effectués dans l'autre espace.



Les axes factoriels étant orthogonaux deux à deux, on obtient ainsi une série de variables artificielles non corrélées entre elles, appelées *composantes principales*, qui synthétisent les corrélations de l'ensemble des variables initiales.

### Remarques

1) L'analyse en composantes principales ne traduit que des liaisons linéaires entre les variables. Un coefficient de corrélation faible entre deux variables signifie donc que celles-ci sont indépendantes linéairement alors qu'il peut exister une relation de degré supérieur à 1 (liaison non linéaire).

2) La coordonnée d'un point-variable sur l'axe  $\alpha$  est nécessairement inférieure à 1 en valeur absolue :

$$|\varphi_{\alpha j}| \leq 1$$

et :

$$\sum_{\alpha=1}^p \text{cor}^2(j, \psi_{\alpha}) = 1$$

3) Le nuage de points-variables dans  $\mathbb{R}^n$  n'est pas centré sur l'origine.

### 1.2.5 Individus et variables supplémentaires

On dispose d'informations complémentaires que l'on veut rapporter à l'analyse des temps d'activités des hommes actifs regroupés en catégories. Par exemple, on veut enrichir cette analyse par une série d'indicateurs d'habitudes de fréquentation-média, constituant des variables continues et par le niveau d'éducation et l'âge qui sont des variables nominales. On désire également positionner, dans le nuage analysé, des groupes de femmes actives, que l'on va mettre en lignes supplémentaires.

Le tableau de données  $\mathbf{R}$  peut être ainsi complété en colonne par un tableau  $\mathbf{R}^+$  à  $n_s$  lignes et  $p$  colonnes. Il n'est pas nécessaire de connaître le tableau  $\mathbf{R}_+$  à  $n_s$  lignes et  $p_s$  colonnes croisant individus et variables supplémentaires (cf. figure 1.2 - 8).

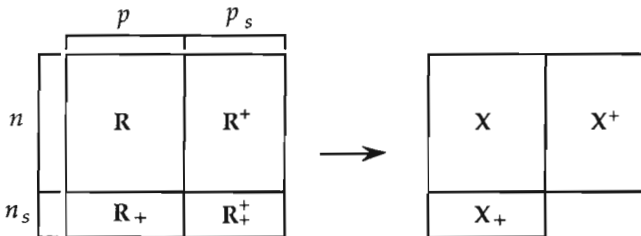


Figure 1.2 - 8  
Lignes et colonnes supplémentaires

Les tableaux  $\mathbf{R}^+$  et  $\mathbf{R}_+$  vont être respectivement transformés en tableaux  $\mathbf{X}^+$  et  $\mathbf{X}_+$  de façon à rendre ces nouvelles lignes et colonnes comparables à celles de  $\mathbf{X}$ .

Dans l'espace  $\mathbb{R}^n$  les  $p_s$  points-variables supplémentaires peuvent être continues ou nominales<sup>1</sup>.

### a – Individus supplémentaires

Pour situer les individus supplémentaires par rapport aux autres dans l'espace  $\mathbb{R}^p$  il est nécessaire de les positionner par rapport au centre de gravité du nuage (déjà calculé sur les  $n$  individus) et de diviser les coordonnées par les écarts-types des variables (déjà calculés sur les  $n$  individus). D'où la transformation :

$$x_{+ij} = \frac{r_{+ij} - \bar{r}_j}{s_j}$$

Les coordonnées des nouveaux points-individus sont donc les  $n_s$  lignes du vecteur  $X_+u_\alpha$ .

En appelant  $X_s$  le tableau  $\begin{bmatrix} X \\ X_+ \end{bmatrix}$  on obtient simultanément les  $n + n_s$  coordonnées des individus analysés et supplémentaires en effectuant le produit  $X_s u_\alpha$ .

### b – Variables continues supplémentaires

Dans  $\mathbb{R}^n$ , pour que les distances entre variables s'interprètent encore en termes de corrélations, ces variables doivent être à valeurs numériques continues et il est indispensable d'effectuer la transformation :

$$x_{ij}^+ = \frac{r_{ij}^+ - \bar{r}_j^+}{s_j^+}$$

On calcule donc les nouvelles moyennes et les nouveaux écarts-types correspondant aux variables supplémentaires, pour positionner celles-ci sur la sphère de rayon unité.

Les coordonnées des  $p_s$  variables supplémentaires sur cet axe sont donc les  $p_s$  lignes du vecteur  $X^+v_\alpha$  et correspondent chacune au coefficient de corrélation entre la variable et le facteur (cf. formule [1.2 - 4]).

### c – Variables nominales supplémentaires

Si la variable à mettre en supplémentaire est nominale, on ne peut plus effectuer la même transformation.

Dans ce cas, on ramène la variable nominale ayant  $m$  modalités, à  $m$  groupes d'individus définis par les modalités de la variable. On traite

---

<sup>1</sup> L'homogénéité de nature des variables supplémentaires n'est plus exigée sous réserve des transformations indiquées.

ensuite ces  $m$  groupes d'individus comme des individus supplémentaires. Ce sont les centres de gravité de ces groupes d'individus qui vont être positionnés dans l'espace  $\mathbb{R}^p$ .

Supposons, par exemple, que l'on mesure la taille et le poids de 10 individus et que l'on désire mettre en supplémentaire la variable sexe. Nous disposons du tableau de mesures représenté figure 1.2 - 9.

variables continues actives			variable nominale supplémentaire à 2 modalités	modalité 1 (homme)		modalité 2 (femme)	
taille	poids	sexe		taille	poids	taille	poids
1	150	45	2			150	45
	168	68	1	168	68		
	175	72	1	175	72		
	178	70	2			178	70
i	185	70	1	185	70		
	160	53	2			160	53
	165	49	2			165	49
	180	90	1	180	90		
	175	65	2			175	65
10	174	72	2			174	72
lignes supplém.				177	75	167	59

Figure 1.2 - 9  
Les modalités de la variable nominale supplémentaire sont des individus supplémentaires

On calcule alors la taille et le poids moyens des hommes (177; 75) et celui des femmes (167; 59). Ce sont ces points moyens qui vont être positionnés parmi les points-individus.

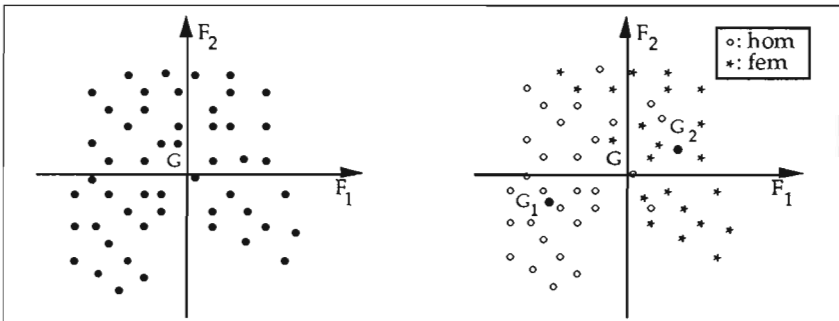


Figure 1.2 - 10  
Représentation d'une variable nominale supplémentaire

La représentation par deux points  $G_1$  et  $G_2$  d'une variable nominale à deux modalités est esquissée sur la figure 1.2 - 10.

L'analyse d'une variable nominale supplémentaire ne se fait donc pas dans  $\mathbb{R}^n$  mais dans  $\mathbb{R}^p$ .

La figure 1.2 - 11 schématise le positionnement des variables supplémentaires :

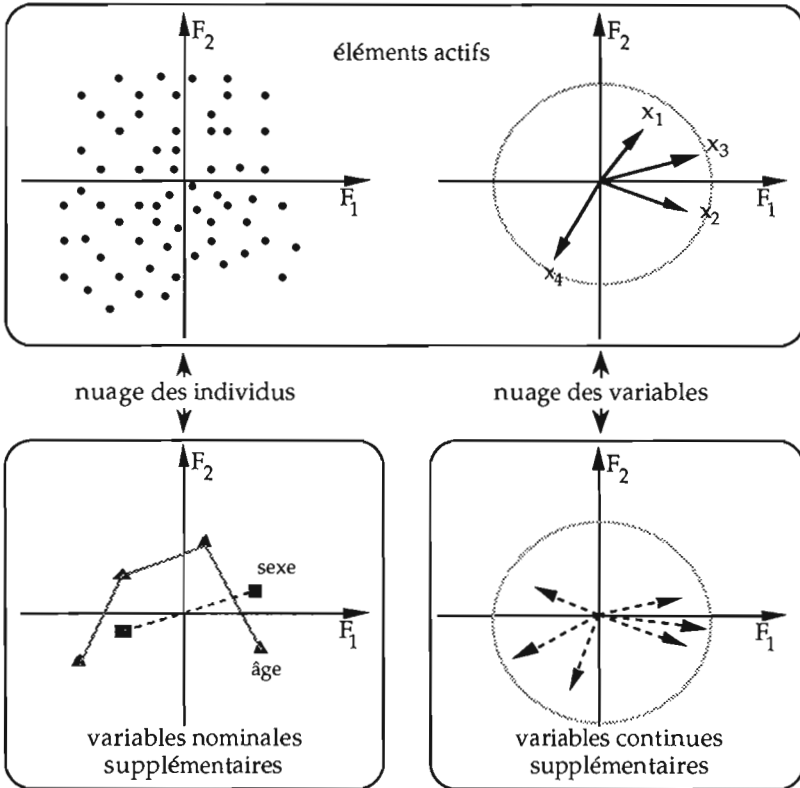


Figure 1.2 - 11  
Représentation des variables supplémentaires

### 1.2.6 Représentation simultanée

L'analyse du nuage des variables est déduite de celle du nuage des individus : la représentation des variables sur les axes factoriels dans  $\mathbb{R}^n$  aide l'interprétation des axes factoriels dans  $\mathbb{R}^p$  et réciproquement.

#### a – Représentation séparée des deux nuages

Mais les deux nuages ne sont pas dans le même repère, ce qui rend impossible la représentation simultanée des individus et des variables.

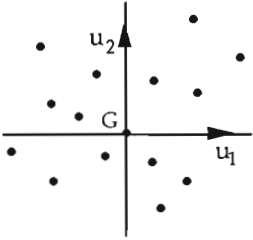
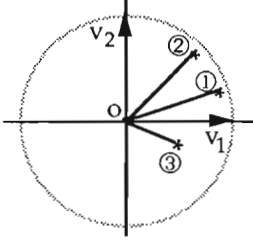
Dans l'espace $\mathbb{R}^p$	Dans l'espace $\mathbb{R}^n$
<p>L'analyse du nuage des <math>n</math> points-individus se fait dans le repère :</p> <p><math>\{G, u_1, \dots, u_\alpha, \dots, u_p\}</math></p>  <p>La représentation des individus sur les axes factoriels fournit la meilleure visualisation approchée des distances entre les individus.</p>	<p>L'analyse du nuage des <math>p</math> points-variables se fait dans le repère :</p> <p><math>\{O, v_1, \dots, v_\alpha, \dots, v_n\}</math></p>  <p>La représentation des variables sur les axes factoriels fournit une synthèse graphique de la matrice de corrélations.</p>

Figure 1.2 - 12

Nuage des individus dans  $\mathbb{R}^p$ 

Figure 1.2 - 13

Nuage des variables dans  $\mathbb{R}^n$ 

Les proximités entre individus s'interprètent en termes de similitudes de comportement vis-à-vis des variables et les proximités entre variables en termes de corrélations. Il faut bien se garder d'interpréter la distance séparant un point-variable d'un point-individu car ces deux points ne font pas partie d'un même nuage dans un même espace : la superposition de ces deux plans factoriels est dénuée de sens.

### b – Justification d'une représentation simultanée

Cependant si l'on considère non plus des points-variables mais des directions de variables dans  $\mathbb{R}^p$ , on peut alors envisager de représenter simultanément, dans cet espace, à la fois les points-individus et des vecteurs représentant les variables.

Dans l'espace  $\mathbb{R}^p$  des  $n$  points-individus, après transformation du tableau de données, on dispose de deux systèmes d'axes :

- les anciens axes unitaires  $(e_1, e_2, \dots, e_j, \dots, e_p)$  correspondant aux  $p$  variables avant l'analyse où :

$$e_j = (0, 0, \dots, 1, 0, \dots, 0)$$

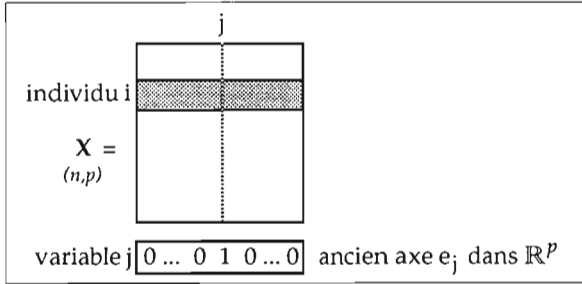
$\{e_j, (j = 1, \dots, p)\}$  est le système d'axes de référence pour les coordonnées initiales des individus.

- les nouveaux axes unitaires  $\{u_\alpha (\alpha = 1, \dots, p)\}$  constitués des axes factoriels.

La possibilité d'une représentation simultanée réside alors dans la projection (en ligne supplémentaire) de l'ancien axe  $e_j$  sur le nouvel axe  $u_\alpha$ .

La coordonnée de la projection de  $e_j$  sur  $u_\alpha$  vaut :

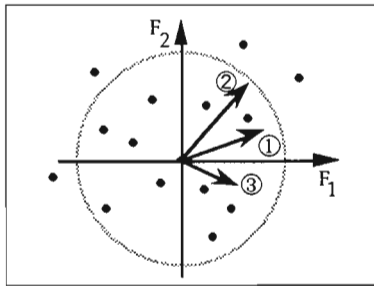
$$e_j \cdot u_\alpha = u_{\alpha j}$$



**Figure 1.2 - 14**  
Ancien axe dans  $\mathbb{R}^p$  en supplémentaire  
La variable  $j$  est un individu particulier

Il est ainsi possible de représenter dans  $\mathbb{R}^p$  les directions données par les variables d'origine sur le plan factoriel du nuage des individus ; ces directions peuvent être matérialisées par des vecteurs unitaires. Ces vecteurs constituent le repère d'origine dans lequel on a construit le nuage des individus. Ils sont donc orthogonaux deux à deux <sup>1</sup>.

Ce qui s'appellera *représentation simultanée* est donc "l'écrasement" du repère orthonormé des axes d'origine sur le plan factoriel du nuage des individus.



**Figure 1.2 - 15**  
Projection des anciens axes sur le plan factoriel  
du nuage des individus

Rappelons que, dans  $\mathbb{R}^n$ , la coordonnée de la variable  $j$  sur l'axe  $\alpha$  est égale au coefficient de corrélation (cf. formule [1.2 - 4]) entre la variable et le facteur et vaut :

<sup>1</sup> Il apparaît donc clairement que cette représentation des variables est distincte du nuage de variables décrit précédemment.

$$\varphi_{\alpha j} = \sqrt{\lambda_{\alpha}} u_{\alpha j}$$

Les deux nuages des variables ne coïncident donc pas. Ils diffèrent l'un de l'autre par une dilatation définie sur chaque axe par le coefficient  $\sqrt{\lambda_{\alpha}}$ .

Dans le cas de la représentation simultanée, qui est en fait une représentation dans  $\mathbb{R}^p$ , on n'interprète pas la distance entre deux variables en terme de corrélation, puisqu'il s'agit en réalité des extrémités de deux vecteurs unitaires orthogonaux<sup>1</sup>. L'interprétation de la distance entre deux variables (en terme de corrélation) ne peut se faire<sup>2</sup> que dans  $\mathbb{R}^n$ . En tenant compte de ces considérations, il est licite de comparer, sur la représentation simultanée, les positions respectives de deux individus vis-à-vis de l'ensemble des variables, ou de deux variables vis-à-vis de l'ensemble des individus.

On dispose ainsi d'une perspective déformée du système d'axes originel tenant compte des liaisons existant entre les variables initiales.

La direction d'une variable définit des zones pour les individus : d'un côté, ceux qui prennent des fortes valeurs pour cette variable et, à l'opposé, ceux qui prennent des valeurs faibles. On s'intéressera à l'éloignement des individus dans la direction de la variable. A l'intersection des axes se trouvent les valeurs moyennes de toutes les variables.

**Remarques:**

1) Si l'échelle des coordonnées des points-variables a une interprétation en termes de corrélations, il n'en est pas de même pour les points-individus. On appliquera à leurs coordonnées un coefficient de dilatation convenable. La valeur  $\sqrt{n/p}$  assure souvent un positionnement dans le plan compatible avec la répartition des points-variables et permet ainsi une représentation équilibrée des deux nuages.

2) Dans la représentation simultanée, il ne peut y avoir de variables continues supplémentaires (elles ne constituent pas des axes d'origine pour le positionnement des individus). Il peut y avoir des variables nominales supplémentaires car ce sont des individus supplémentaires.

## 1.2.7 Analyse en composantes principales non normée

L'analyse en composantes principales non normée revient à considérer le nuage de points centré et non réduit. On généralisera l'analyse en faisant jouer maintenant à chaque point-individu un rôle proportionnel à sa masse (ce que l'on aurait évidemment pu faire à propos de l'analyse normée).

<sup>1</sup> Toutes ces distances sont égales à  $\sqrt{2}$  dans l'espace complet.

<sup>2</sup> On note toutefois que le nuage projeté des extrémités des vecteurs unitaires dans  $\mathbb{R}^p$  et le nuage des extrémités des vecteurs variables dans  $\mathbb{R}^n$  ont généralement des allures voisines, surtout si les valeurs propres sont presque égales, car alors la dilatation est peu déformante.

### a – Principe de l'analyse et nuage des individus

Plaçons-nous dans l'espace  $\mathbb{R}^p$  et considérons le nuage des points-individus pesants, centré sur le centre de gravité G. L'analyse en composantes principales revient à effectuer une analyse générale de points pondérés avec comme origine le centre de gravité du nuage.

Le tableau de données initiales  $R$  subit plusieurs transformations : on construit le tableau  $X$  de données centrées et chaque individu  $i$  est affecté d'une masse ou d'un poids<sup>1</sup>  $p_i$  éléments diagonaux de la matrice diagonale  $N$ .

Le tableau  $Z$  soumis à l'analyse en composantes principales non normée est par conséquent de la forme :

$$Z = N^{1/2}X$$

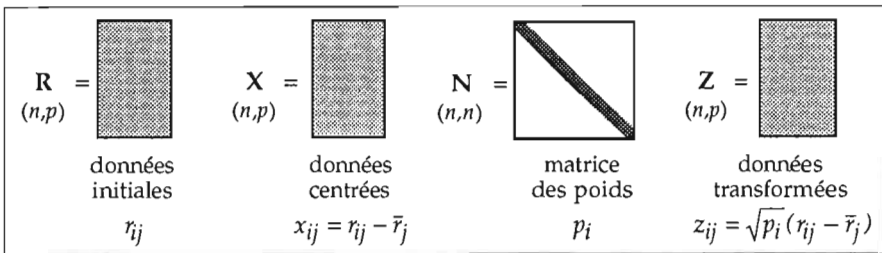


Figure 1.2 - 16  
Transformation du tableau de données  
en analyse en composantes principales non normée

La matrice à diagonaliser est la matrice d'inertie autour du centre de gravité du nuage G :

$$A = Z'Z = X'NX$$

de terme général :

$$a_{jj'} = \sum_{i=1}^n p_i (r_{ij} - \bar{r}_j)(r_{ij'} - \bar{r}_{j'})$$

Si les masses représentent des fréquences, alors la matrice à diagonaliser est la *matrice des covariances*.

A partir de là, on détermine les axes factoriels  $u_\alpha$  tels que  $u'_\alpha u_\alpha = 1$ . Les coordonnées factorielles sur ces axes sont données par :

$$\psi_\alpha = Xu_\alpha$$

dont les composantes s'écrivent :

$$\psi_{\alpha i} = \sum_{j=1}^p (r_{ij} - \bar{r}_j) u_{\alpha j}$$

<sup>1</sup> Les termes de masse et de poids sont utilisés indifféremment en statistique. Ils désignent souvent des fréquences relatives ou des probabilités *a priori*.



avec :

$$\sum_{i=1}^n p_i \psi_{i\alpha}^2 = \lambda_\alpha$$

### b – Nuage des variables

L'analyse du nuage des  $p$  variables dans  $\mathbb{R}^n$  revient à faire l'analyse générale du tableau  $\mathbf{Z}$  :

$$z_{ij} = \sqrt{p_i} (r_{ij} - \bar{r}_j)$$

avec :

$$\sum_{i=1}^n p_i = 1 \quad \text{et} \quad \bar{r}_j = \sum_{i=1}^n p_i r_{ij}$$

La distance induite entre deux variables s'exprime par :

$$d^2(j, j') = \sum_{i=1}^n (z_{ij} - z_{ij'})^2$$

soit :

$$d^2(j, j') = \sum_{i=1}^n z_{ij}^2 + \sum_{i=1}^n z_{ij'}^2 - 2 \sum_{i=1}^n z_{ij} z_{ij'}$$

Par conséquent<sup>1</sup> :

$$d^2(j, j') = \text{var}(j) + \text{var}(j') - 2\text{cov}(j, j') \quad [1.2 - 5]$$

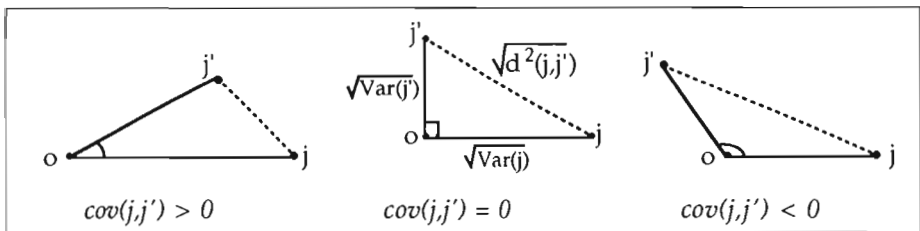


Figure 1.2 - 17  
Distance entre deux variables

La distance entre deux variables s'exprime en terme de covariance et augmente avec les variances. Elle diminue si la liaison est positive et augmente si la liaison est négative.

La distance d'une variable à l'origine des axes est sa variance :

$$d^2(O, j) = \text{var}(j) = \sum_{i=1}^n z_{ij}^2 = \sum_{i=1}^n p_i (r_{ij} - \bar{r}_j)^2$$

<sup>1</sup> La formule [1.2 - 3] est un cas particulier lorsque  $\text{var}(j) = \text{var}(j') = 1$ , c'est-à-dire lorsqu'il s'agit d'une analyse en composantes principales normée.

Par conséquent, pour l'analyse en composantes principales non normée, la sphère de corrélations n'est plus l'espace de départ<sup>1</sup>.

### 1.2.8 Analyses non-paramétriques

Ces méthodes ne diffèrent de la précédente que par une transformation préliminaire des données. Elles sont recommandées lorsque les données de base sont hétérogènes. Elles donnent des résultats d'une grande robustesse, se prêtant par ailleurs à des interprétations simples en termes statistiques.

#### a – Analyse des rangs

Le tableau initial des données est transformé en tableau de rangs. L'observation  $i$  de la variable  $j$  consiste alors en un classement  $q_{ij}$  : c'est le rang de l'observation  $i$  lorsque les  $n$  observations sont classées par ordre de grandeur. Dans ces conditions, la distance entre deux variables  $j$  et  $j'$  est définie par la formule<sup>2</sup> :

$$d^2(j, j') = \frac{6}{n(n-1)(n+1)} \sum_{i=1}^n (q_{ij} - q_{ij'})^2$$

L'utilisation des rangs sera justifiée dans les contextes suivants :

- Les données de base peuvent être elles-mêmes des classements, auquel cas ce type d'analyse s'impose.
- Les échelles de mesure des variables peuvent être si différentes que l'opération de réduction pratiquée par l'analyse en composantes principales normée reste insuffisante. De plus cette opération ne remédie pas par exemple à la dissymétrie des distributions. Il paraît enfin plus justifié de synthétiser une famille de classements qu'un ensemble très hétérogène de mesures.
- Les hypothèses *a priori* faites implicitement sur les mesures sont plus faibles et par conséquent moins arbitraires : la loi des distances est maintenant non-paramétrique; nous disposerons donc de seuils de confiance qui ne dépendront que de l'hypothèse de continuité des lois des observations, plus plausible que celle de normalité.
- Enfin, les représentations fournies sont robustes, très peu sensibles à l'existence de valeurs aberrantes, ce qui sera souvent une qualité appréciable.

Les règles d'interprétation se déduisent de celles de l'analyse en composantes principales puisque c'est cette analyse que l'on effectue après

<sup>1</sup> Dans une représentation simultanée, les anciens axes (distance 1 de l'origine) seront toujours dans un cercle de corrélations (cf. § 1.2.6.).

<sup>2</sup> On reconnaît dans cette formule le complément à 1 du coefficient de corrélation de Spearman (cf. Kendall, 1962).

l'opération de transformation en rangs<sup>1</sup>. La proximité entre deux variables s'interprète en terme de corrélation de rangs : deux variables seront très proches pour des classements voisins des observations ; au contraire, deux variables éloignées correspondront à des classements pratiquement inverses. Deux observations seront proches si elles ont des rangs similaires pour chacune des variables. Enfin, dans la représentation simultanée, on a une idée de l'ensemble du classement des observations pour une variable en examinant les positions respectives de cette variable et de l'ensemble des observations<sup>2</sup>.

### b – Analyse en composantes robustes

Le critère d'ajustement des moindres-carrés est particulièrement bien adapté à la distribution normale. Dans le cas d'une distribution uniforme (cas de l'analyse des rangs), il tend à donner une importance excessive aux observations extrêmes. On rendra donc plus robuste l'analyse par une transformation qui "normalise" la distribution uniforme des rangs.

Considérons la  $k^{\text{ième}}$  observation de  $n$  observations rangées et soit  $F$  la fonction de répartition de la loi Normale. On remplacera l'observation de rang  $k$  par la valeur  $y_k$  tirée de la fonction de répartition inverse de la loi Normale<sup>3</sup> :

$$y_k = F^{-1}\left(\frac{k}{n+1}\right)$$

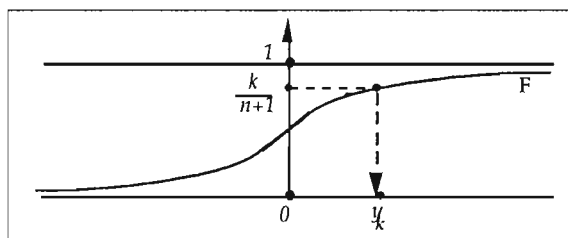


Figure 1.2 - 18

Transformation suivant la fonction de répartition inverse de la loi Normale

Pour  $n$  grand, la transformation est équivalente au remplacement de la  $k^{\text{ième}}$  observation par l'espérance de la  $k^{\text{ième}}$  observation dans un échantillon rangé de  $n$  valeurs normales.

<sup>1</sup> Notons qu'il n'est pas indispensable ici de réduire en terme d'écart-type car tous les rangs ont la même variance.

<sup>2</sup> Ajoutons enfin que le caractère non-paramétrique de la représentation obtenue permet de procéder à des tests de validité sur les valeurs propres. La loi des valeurs propres issues de l'analyse d'un tableau de rangs ne dépend en effet que des paramètres  $n$  et  $p$ , nombres de lignes et de colonnes du tableau. Il est donc possible de procéder à une tabulation permettant de connaître les seuils de signification des valeurs propres.

<sup>3</sup> On trouve déjà ce type de transformation dans Fisher et Yates (1949).

### 1.2.9 Aperçu sur les autres méthodes dérivées

De nombreuses techniques sont directement dérivées de l'analyse en composantes principales. Les variantes non-paramétriques du paragraphe précédent en sont des exemples.

Certaines présentations de l'analyse des correspondances (cf. section 1.3) considèrent cette méthode comme une analyse en composantes principale particulière. Cela est possible si l'on traite les deux espaces (lignes et colonnes) séparément, ce qui n'est pas l'optique choisie ici. Ce traitement séparé masque un des apports méthodologiques fondamentaux des analyses factorielles descriptives. L'analyse en composantes principales, qu'il s'agisse d'analyse normée ou non-normée, analyse les individus par rapport à leur *centre de gravité* et les variables par rapport à l'*origine des axes*. Cette dissymétrie de traitement des lignes et des colonnes correspond à des domaines d'applications spécifiques et induit des règles d'interprétation particulières. La décomposition aux valeurs singulières (ou encore analyse générale, ou théorème d'Eckart et Young) est bien le noyau théorique commun-aux deux méthodes.

Citons parmi les méthodes dérivées l'*analyse des corrélations partielles* ou *analyse avec variables instrumentales* (Rao, 1964), qui sera abordée au chapitre 3, section 3.6. Dans ce cas, on ne se contente plus d'éliminer les effets de l'hétérogénéité des variables (opérations de centrage et de réduction) mais on se propose d'éliminer également l'effet d'autres variables, en procédant à une régression multiple préalable. L'analyse logarithmique (Kazmierczak, 1985) est une analyse en composantes principales non-normée du tableau (doublement centré en lignes et en colonnes) des logarithmes des variables initiales. Cette variante possède d'intéressantes propriétés de stabilité et de robustesse.

D'autres techniques enfin, comme la régression sur composantes principales (§ 3.2.5) ou la classification sur facteurs (section 2.4) sont plutôt des techniques complémentaires que dérivées.

### 1.2.10 Eléments pour l'interprétation

Les axes factoriels permettent d'obtenir la meilleure visualisation approchée (au sens des moindres carrés) des distances entre les individus d'une part et entre les variables d'autre part. Pour interpréter les facteurs, il faut apprécier correctement cette approximation. On procédera dans un premier temps à un examen de l'inertie<sup>1</sup> de chaque facteur puis on s'intéressera aux éléments contribuant à construire et à définir les facteurs.

---

<sup>1</sup> Inertie, terme emprunté à la mécanique, est ici synonyme de variance, terme statistique.

### a – Inertie liée aux facteurs

Rappelons que la valeur propre (ou l'inertie liée à un facteur) est la variance des coordonnées des points-individus sur l'axe correspondant. C'est un indice de dispersion du nuage des individus dans la direction définie par l'axe.

Il n'existe pas de critères simples et définitifs qui permettent de se prononcer sur l'importance d'une valeur propre. Les problèmes de validité des résultats communs à l'ensemble des méthodes factorielles seront étudiés plus systématiquement dans le chapitre 4. On mentionnera ici simplement les règles pratiques les plus courantes.

Dans une analyse normée, la somme des inerties est égale au nombre de variables et donc l'inertie moyenne vaut 1. Chaque axe étant une combinaison particulière des variables d'origine, on s'intéresse en général aux axes ayant une inertie "notablement" supérieure à la moyenne<sup>1</sup>. On observe souvent une décroissance assez irrégulière des premières valeurs propres (Figure 1.2 - 19).

Si les données sont peu structurées (les variables ne sont pas fortement corrélées entre elles), le nuage a une forme "régulière". Dans ce cas, les valeurs propres sont régulièrement décroissantes (Figure 1.2 - 20) et l'analyse factorielle ne fournira pas des résultats intéressants.

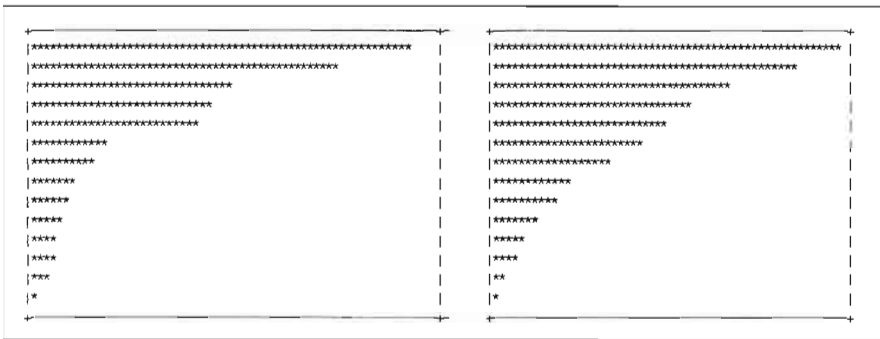


Figure 1.2 - 19  
Paliers dans la décroissance  
des valeurs propres

Figure 1.2 - 20  
Décroissance régulière  
des valeurs propres

Les pourcentages d'inertie des axes définissent les "pouvoirs explicatifs" des facteurs : ils représentent la part de la variance (ou inertie) totale prise en compte par chaque facteur. Son appréciation doit cependant tenir compte du nombre de variables et du nombre d'individus. Un taux d'inertie (relatif à un axe) égal à 10% peut être une valeur importante si le tableau possède 100 variables et faible s'il n'en a que 10. Comme nous le signalerons à propos de l'analyse des correspondances (section 1.3) et sur la validité des résultats

<sup>1</sup> Cette règle, toute empirique, est adoptée par certains utilisateurs.

(chapitre 4), l'inertie est une mesure pessimiste du pouvoir explicatif des facteurs, liée parfois de façon assez arbitraire au codage des données.

L'importance d'un facteur peut dépendre d'informations exogènes (variables supplémentaires par exemple).

Il existe d'autres aides à l'interprétation qui permettent d'apprécier les rôles respectifs des lignes, des colonnes, des axes en analyse en composantes principales.

### **b – Aides à l'interprétation**

On procède axe par axe pour définir les composantes principales. L'examen du plan factoriel permet de visualiser les corrélations entre les variables et d'identifier des groupes d'individus ayant pris les mêmes valeurs pour les mêmes variables.

Considérons le cas de l'analyse en composantes principales normée.

#### *- les variables*

Nous ne nous plaçons pas ici dans le cas de la représentation simultanée mais dans le nuage des variables ( $p$  points de  $\mathbb{R}^n$ ).

Les variables fortement corrélées avec un axe vont contribuer à la définition de cet axe. Cette corrélation se lit directement sur le graphique puisqu'il s'agit de la coordonnée du point-variable  $j$  sur l'axe  $\alpha$  (formule [1.2 - 4]).

On s'intéresse par conséquent aux variables présentant les plus fortes coordonnées (ce qui les situent proches du cercle de corrélations) et l'on interprétera les composantes principales en fonction des regroupements de certaines de ces variables et de l'opposition avec les autres.

Rappelons que le cosinus de l'angle sous lequel on voit deux points-variables actives dans  $\mathbb{R}^n$  n'est autre que le coefficient de corrélation de ces deux variables. Selon la qualité de l'ajustement, cette propriété sera plus ou moins bien conservée en projection. On se gardera d'interpréter la distance entre deux variables actives qui ne sont pas proches du cercle de corrélation.

Ainsi l'examen du plan factoriel permet de visualiser les distances réelles et donc les corrélations entre les variables actives et d'apprécier la qualité de leur représentation. La figure 1.2 - 22 du paragraphe suivant donne un exemple de cercle des corrélations dans le plan des deux premiers facteurs.

Dans le cas des variables continues supplémentaires, les corrélations n'étant pas transitives, il est prudent de ne pas interpréter abusivement les proximités entre variables en terme de corrélation, bien que celles-ci en soient souvent de bonnes images. Ceci sera commenté plus loin au § 1.2.11.

#### *- Les individus*

Si les points-individus ne sont pas anonymes pour l'étude, on s'intéresse à ceux qui participent à la formation des axes. On calcule la *contribution* de

chaque point  $i$  (de masse  $m_i$ ) à l'inertie de l'axe  $\alpha$ . Celle-ci s'exprime par la formule :

$$Cr_{\alpha}(i) = \frac{m_i \psi_{\alpha i}^2}{\lambda_{\alpha}}$$

où  $\lambda_{\alpha}$  est l'inertie de l'axe  $\alpha$  et  $m_i \psi_{\alpha i}^2$  est la contribution de l'individu  $i$  à l'inertie de cet axe. On a :

$$\sum_{i=1}^n Cr_{\alpha}(i) = 1$$

On s'intéressera surtout aux individus qui ont les plus fortes contributions relatives aux axes.

Lorsque les  $n$  individus sont affectés d'une même masse égale à  $1/n$ , l'inertie d'un point varie comme sa distance au centre de gravité. Les individus qui contribuent le plus à la détermination de l'axe sont les plus excentrés et l'examen des coordonnées factorielles ou la lecture du graphique suffisent à interpréter les facteurs dans ce cas. La représentation des individus sur le plan factoriel permet d'apprécier leur répartition et de repérer des zones de densités plus ou moins fortes.

#### - Possibilité d'apparition de facteur "taille"

L'analyse du nuage des variables se faisant à partir de l'origine, les variables peuvent être toutes situées du même côté d'un axe factoriel. Une telle disposition apparaît lorsque toutes les variables sont corrélées positivement entre elles. Si pour un individu, une variable prend une valeur forte, toutes les autres variables prennent également des valeurs fortes. Cette caractéristique apparaît le plus souvent sur le premier axe, que l'on appelle alors "facteur taille".

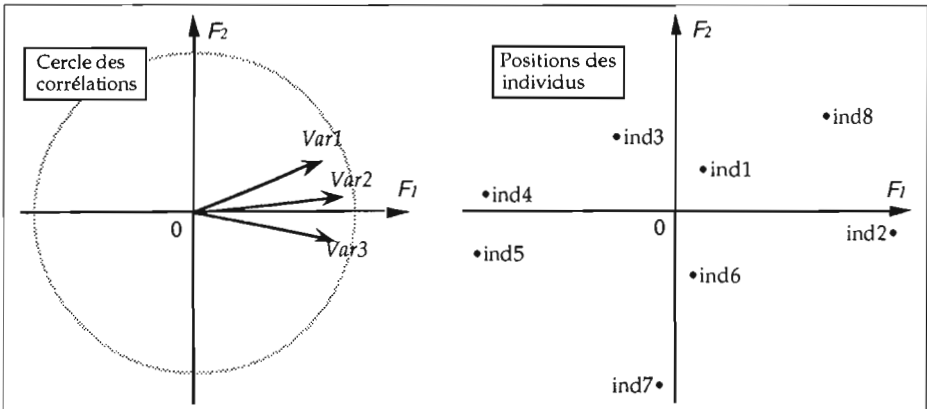


Figure 1.2 - 21  
Exemple de *Facteur taille*

On peut lire, par exemple sur la figure 1.2 - 21, que les individus 4 et 5 ont des comportements semblables caractérisés par des valeurs faibles pour les

trois variables, alors que les individus 2 et 8 ont au contraire simultanément des bons "scores" pour ces mêmes variables. L'orthogonalité des axes fait qu'il ne peut exister qu'un seul facteur taille.

### 1.2.11 Exemple d'application

Nous présentons ici l'exemple (cf. tableau 1.2 - 1) relatif aux temps d'activités quotidiennes évoqué au paragraphe 1.2.1.

Le CESP (Centre d'Étude des Supports de Publicité) a relevé, dans son *Enquête Budget-temps Multimédia* de 1991/1992 auprès de 17 665 personnes, des descripteurs de fréquentation de divers médias (radio, télévision, presse) et des temps d'activités quotidiennes (cf. Boeswillwald, 1992). Ont été également relevées de nombreuses caractéristiques socio-économiques, parmi lesquelles l'âge, le sexe, l'activité, le niveau d'éducation, et le lieu de résidence de ces personnes, ce qui a conduit à créer 96 catégories en croisant ces divers critères.

Nous nous intéressons seulement ici à la sous-population des hommes actifs, soit 27 groupes qui seront, pour cet exemple, les "individus". On cherche à connaître les associations entre les temps consacrés à différentes activités par les "individus" observés et à étudier les liens entre ces familles d'activités et les caractéristiques de base des individus.

Enfin, on se propose d'étudier le lien entre les activités quotidiennes et la fréquentation de divers médias (presse, radio, télévision, cinéma). Pour ce faire, on fera intervenir les caractéristiques socio-économiques (variables nominales) et les habitudes de fréquentation des médias (variables numériques continues) en tant que variables supplémentaires.

#### Lecture du tableau 1.2 - 1

(16 variables continues actives)

Les 27 "individus" (qui sont en réalité dans le cadre de cet exemple des groupes d'individus) sont repérés par un identificateur en 4 caractères :

- le 1er caractère est l'âge du groupe (1=jeune, 2=moyen, 3=âgé)
- le 2ème caractère est ici toujours égal à 1 (car il s'agit ici d'une sélection d'hommes actifs)
- le 3ème est le niveau d'éducation (1=primaire, 2=secondaire, 3=supérieur)
- le 4ème est le type d'agglomération (1=communes rurales; 2=villes moyennes; 3=villes importantes; 4=agglomération parisienne; 5,6,7 = groupes mixtes).

(On trouvera des libellés plus détaillés des variables dans le tableau 1.2 - 2 ci-après.)

On lit par exemple sur la première ligne du tableau 1.2 - 1 que le groupe '1111' (jeunes, actifs, peu instruits, ruraux) consacre en moyenne par jour 463.8 minutes au "sommeil", 23.8 minutes à des activités regroupées sous la rubrique "repos", 107.3 minutes pour les "repas chez soi", etc.



Tableau 1.2 - 1 : Budget-temps agrégé quotidien de 27 groupes d'hommes actifs

IDENT	Somm	Repo	Reps	Repr	Trar	Ména	Visi	Jard	Lois	Disq	Lect	Cour	Prom	A pi	Voit	Fréq
1111	463.8	23.8	107.3	4.8	300.0	21.3	51.0	82.3	10.0	1.2	.0	41.3	6.9	7.1	52.1	135.8
1115	515.6	58.5	102.7	10.4	208.8	41.9	30.0	32.9	2.1	4.6	.6	33.7	8.3	24.6	29.4	225.8
1121	463.3	34.2	84.8	17.1	298.3	18.1	37.8	55.8	18.4	5.9	2.6	30.7	5.9	8.8	56.7	135.8
1122	456.4	43.1	74.2	21.9	239.0	26.0	51.2	59.7	18.4	3.6	4.6	52.2	9.5	10.8	72.7	142.3
1123	478.0	44.2	76.7	15.2	212.3	22.3	42.0	43.7	18.4	2.3	6.4	48.3	14.7	15.5	72.8	167.7
1124	465.1	41.6	85.2	23.7	226.0	37.0	42.5	16.3	10.7	8.7	9.4	44.3	13.7	19.8	59.0	145.1
1136	458.4	47.4	94.7	15.1	314.3	25.3	39.1	42.4	16.9	.9	16.7	34.5	4.6	6.4	61.5	103.4
1133	457.2	30.7	82.0	26.2	269.8	52.1	37.6	35.6	25.6	6.0	8.0	42.8	10.4	12.0	81.4	107.6
1134	465.2	40.2	78.6	31.1	268.6	36.3	21.6	4.0	19.4	6.0	14.8	46.9	10.7	21.9	48.3	82.4
2111	449.0	42.1	86.2	7.9	312.5	15.1	16.1	112.9	15.4	.0	2.2	32.1	7.6	8.1	60.1	153.9
2112	450.2	63.1	86.7	9.8	249.6	40.4	55.6	83.3	3.0	2.2	.0	45.0	9.4	10.4	61.9	145.4
2117	455.2	47.4	95.6	9.0	250.8	30.4	13.5	57.3	7.9	2.9	7.0	52.2	15.1	15.7	49.1	194.8
2121	461.9	39.3	90.3	8.5	323.5	14.9	21.7	81.8	15.4	1.2	5.3	26.0	3.8	7.4	59.6	130.8
2122	453.7	44.7	97.5	18.7	269.0	23.1	39.6	93.5	3.1	3.4	12.1	42.0	12.1	10.6	62.4	129.1
2123	433.1	49.8	91.7	12.6	283.7	22.4	21.0	62.9	13.1	6.2	7.3	38.1	11.6	11.7	47.6	168.6
2124	438.3	32.8	102.3	11.1	338.3	28.0	6.5	64.8	13.8	1.4	19.8	34.9	7.4	14.1	53.2	130.5
2131	457.7	44.0	87.9	6.9	313.0	24.4	23.2	63.8	9.2	.6	11.8	30.0	7.3	7.5	69.7	108.3
2132	455.0	47.0	78.9	31.6	380.6	23.9	7.5	40.0	13.0	.0	10.3	23.3	1.4	9.4	59.4	100.0
2133	467.3	37.5	86.9	21.9	264.0	40.8	27.6	33.4	11.9	1.6	10.8	45.3	6.7	10.7	72.8	135.2
2134	433.5	35.6	76.1	17.1	355.0	34.1	13.4	31.7	12.6	3.2	13.2	37.5	8.5	22.3	57.5	96.5
3116	473.0	51.5	99.3	6.3	356.3	21.2	27.6	82.1	8.6	.0	1.5	35.7	13.4	7.1	40.6	107.7
3117	461.9	60.0	103.7	9.1	240.5	35.3	14.5	83.4	1.4	2.0	7.4	46.1	5.7	16.6	53.3	183.7
3121	453.4	45.6	86.2	7.8	358.7	12.9	18.5	54.4	4.2	.0	4.9	34.3	3.3	10.3	48.7	143.1
3122	485.1	53.5	86.0	.3	222.4	24.7	23.2	91.9	8.5	.0	3.7	52.9	7.1	9.9	75.3	166.3
3123	456.7	43.2	94.6	12.1	265.3	30.5	23.7	61.1	9.1	2.3	11.6	50.1	17.6	13.2	46.3	185.3
3136	444.2	53.6	90.7	7.2	302.4	31.7	16.4	97.6	4.7	2.4	4.3	38.8	13.6	11.4	61.8	127.2
3137	438.4	50.7	81.0	11.2	306.6	19.3	23.8	10.5	13.6	.0	18.4	67.6	8.3	18.6	63.1	143.3

L'analyse du tableau de données (tableau 1.2 - 1) nous conduit tout d'abord à calculer les paramètres descriptifs élémentaires regroupés dans le tableau 1.2 - 2. Les moyennes et écart-types vont servir à transformer les variables de base et n'interviendront plus directement dans la suite. Il importe donc de prendre connaissance de ces mesures de niveau et de dispersion. Les valeurs extrêmes sont également utiles pour apprécier la qualité de l'information recueillie.

Ce tableau donne les mêmes paramètres pour les variables continues supplémentaires. Pour le thème "budget-temps", trois variables seront projetées a posteriori : autres activités, total des activités à domicile, total des activités déclarées en déplacement, ces deux dernières étant des regroupements de variables actives ; pour le thème "fréquentation média" (qui donne lieu à une mesure de durée globale au niveau des variables actives) six variables décrivent les intensités de contacts avec le cinéma, la radio, la télévision, les presses quotidiennes et magazines, en isolant dans celle-ci les hebdomadaires dits "News".

**Tableau 1.2 - 2**  
**Statistiques sommaires des variables continues**  
**effectif total : 27**

IDEN - LIBELLE	MOYENNE	ECART- TYPE	MINIMUM	MAXIMUM
<b>Variables actives</b>				
Somm - Sommeil	458.91	16.47	433.10	515.60
Repo - Repos	44.63	8.90	23.80	63.10
Reps - Repas chez soi	89.18	8.90	74.20	107.30
Repr - Repas restaurant	13.87	7.82	.30	31.60
Trar - Travail rémunéré	286.27	46.75	208.80	380.60
Ména - Ménage	27.90	9.29	12.90	52.10
Visi - Visite à amis	27.64	13.26	6.50	55.60
Jard - Jardinage, Bricolage	58.49	27.39	4.00	112.90
Lois - Loisirs extérieur	11.42	5.95	1.40	25.60
Disq - Disque cassette	2.54	2.32	.00	8.70
Lect - Lecture livre	7.95	5.47	.00	19.80
Cour - Courses démarches	40.99	9.47	23.30	67.60
Prom - Promenade	9.06	3.88	1.40	17.60
A pi - Déplacement a pied	12.66	5.01	6.40	24.60
Voit - Déplacement en Voiture	58.38	11.29	29.40	81.40
Fréq - Fréquentation Média	140.58	32.56	82.40	225.80
<b>Variables continues supplémentaires</b>				
Autr - Autres activités	12.71	5.70	2.10	25.90
Domi - Total Domicile	928.73	49.92	826.00	1034.00
Tdep - Total Déplacement	88.45	14.65	67.50	122.10
Habitudes Cinema	.14	.14	.00	.60
Habitudes Radio.	1.92	.23	1.49	2.64
Habitudes Télévision	3.20	.37	2.13	3.90
Habitudes Presse Quotidienne	.18	.14	.03	.53
Habitudes Presse magazine	3.56	.74	2.00	5.31
Habitudes Hebdomadaires News	.31	.18	.00	.67

Tableau 1.2 - 3 : Matrice des corrélations, et valeurs propres correspondantes

Sommeil	1.00																	
Repos	.21	1.00																
Repas c.	.21	.10	1.00															
Repas r.	-.08	-.30	-.53	1.00														
Travail	-.52	-.28	-.02	-.01	1.00													
Ménage	.20	.08	-.01	.39	-.46	1.00												
Visites	.27	-.08	-.07	.10	-.47	.15	1.00											
Jardin.	-.09	.19	.43	-.64	.08	-.37	-.02	1.00										
Loisirs	-.17	-.61	-.55	.52	.10	-.01	.12	-.39	1.00									
Disques	.07	-.17	-.15	.52	-.46	.50	.30	-.42	.25	1.00								
Lecture	-.44	-.21	-.15	.38	.24	.08	-.36	-.51	.27	-.01	1.00							
Courses	-.04	.18	-.17	-.03	-.56	.23	-.24	-.24	-.01	.08	.18	1.00						
Promen.	.00	.09	.04	-.02	-.45	.27	.18	-.01	-.05	.40	-.03	.48	1.00					
A pied	.17	.15	-.14	.28	-.38	.49	-.18	-.62	-.09	.48	.27	.37	.30	1.00				
Voiture	-.19	-.22	-.55	.21	-.15	.10	.27	.03	.44	-.09	.15	.23	-.11	-.33	1.00			
Fréq.med	.40	.42	.37	-.44	-.62	.05	.01	.18	-.45	.07	-.38	.30	.28	.28	-.33	1.00		
		Somm	Repo	Reps	Repr	Trar	Ména	Visi	Jard	Lois	Disq	Lect	Cour	Prom	A pi	Voit	Fréq	

NUMER.	VALEUR PROPRE	POURCENTAGES	POURCENTAGES CUMULES	HISTOGRAMME DES 16 PREMIERES VALEURS PROPRES
1	3.871	24.20	24.20	*****
2	3.660	22.88	47.07	*****
3	2.006	12.54	59.61	*****
4	1.514	9.47	69.08	*****
5	1.126	7.04	76.12	*****
6	.837	5.23	81.35	*****
7	.766	4.79	86.15	*****
8	.596	3.73	89.87	*****
9	.444	2.78	92.65	*****
10	.374	2.34	94.99	*****
11	.246	1.54	96.53	*****
12	.222	1.39	97.92	*****
13	.161	1.01	98.93	*****
14	.114	.72	99.64	***
15	.037	.23	99.88	*
16	.019	.12	100.00	*

La matrice des corrélations (tableau 1.2 - 3) nous fournit des éléments de description des associations entre variables actives. Sa lecture nous donne une première idée du réseau d'interrelations existant entre les variables, mais l'analyse en composantes principales va permettre d'obtenir une synthèse de ces liaisons.

Le premier résultat est constitué par la liste des valeurs propres et des pourcentages de variance (cf. tableau 1.2 - 3). La somme des valeurs propres est égale au nombre de variables soit 16. Les deux premiers axes fournissent presque la moitié de l'inertie (47%) mais l'on sait que ces quantités sont d'interprétation délicate. On note cependant, à la vue de l'histogramme, qu'il existe une concentration nette du nuage dans un sous-espace à deux dimensions, le plan factoriel principal.

On lira sur le tableau 1.2 - 4 les coordonnées des points variables sur les trois premiers axes ainsi que les coordonnées des extrémités des axes unitaires (cf. § 1.2.6) destinés à une éventuelle représentation simultanée des individus et des variables. Les deux premières valeurs propres étant voisines (3.871 et 3.660), leurs racines carrées le sont également (1.97 et 1.91) et donc les nuages bidimensionnels des points variables et des anciens axes unitaires auront des allures très voisines (cf. § 1.2 - 6).

**Tableau 1.2 - 4**  
Coordonnées des variables actives sur les axes 1 à 3

VARIABLES	COORDONNEES			ANCIENS AXES UNIT.		
	1	2	3	1	2	3
Sommeil	.22	-.52	.18	.11	-.27	.13
Repos	.46	-.40	-.17	.23	-.21	-.12
Repas chez soi	.67	-.15	-.23	.34	-.08	-.17
Repas restaurant	-.84	.00	-.07	-.43	.00	-.05
Travail rémunéré	.05	.88	-.34	.03	.46	-.24
Ménage	-.40	-.57	-.08	-.20	-.30	-.06
Visite à amis	-.13	-.33	.73	-.07	-.17	.52
Jardinage, Bricolage	.76	.22	.35	.39	.11	.25
Loisirs extérieur	-.72	.30	.30	-.37	.16	.21
Disque cassette	-.53	-.53	.01	-.27	-.27	.01
Lecture livre	-.54	.24	-.50	-.27	.12	-.36
Courses démarches	-.21	-.54	.11	-.11	-.28	.08
Promenade	-.10	-.58	.04	-.05	-.30	.03
A pied	-.37	-.62	-.57	-.19	-.33	-.40
En Voiture	-.41	.22	.65	-.21	.11	.46
Fréquentation Média	.49	-.68	-.05	.25	-.36	-.03

La figure 1.2 - 22 donne une représentation des variables sur les deux premiers axes factoriels. Les données étant ici centrées réduites, les coordonnées des variables sur les axes sont les coefficients de corrélations entre ces variables et les facteurs.

Le premier axe oppose les activités extérieures ou d'ouverture (lecture, loisir extérieur, repas restaurant, déplacement en voiture) à des activités plus intérieures (jardinage, jeux, bricolage, repas chez soi). Le deuxième axe oppose essentiellement l'activité professionnelle (travail rémunéré) aux

activités de temps disponible ou libre (promenade, disque cassette, fréquentation média) mais aussi le temps passé au ménage et au sommeil.

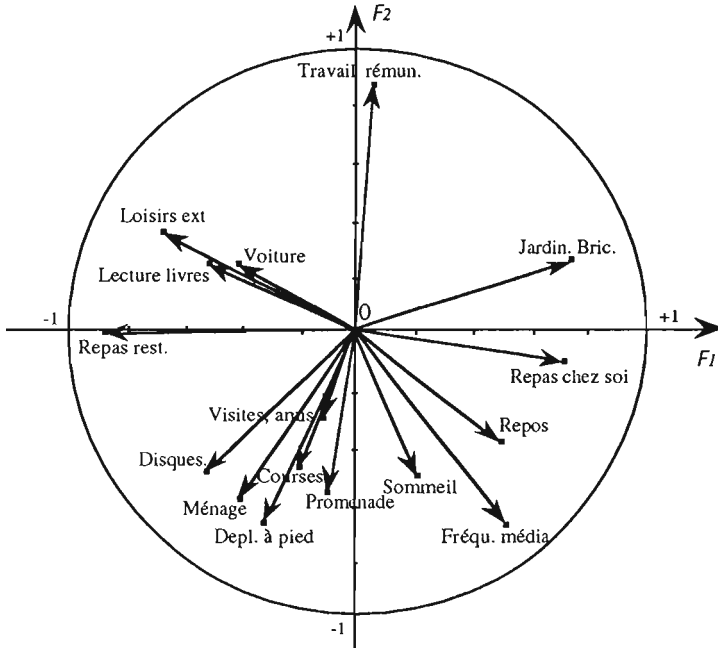


Figure 1.2 - 22  
Représentation des 16 variables actives dans le plan des facteurs 1 et 2

Les variables supplémentaires (tableau 1.2 - 5 et figure 1.2 - 23) relatives aux déplacements et aux médias illustrent ces propos. Les activités "total déplacement" et "total domicile" caractérisent bien le premier axe. La presse quotidienne et surtout le cinéma sont corrélés aux activités dites d'ouvertures, pour lesquelles le temps passé en déplacement est important. Le temps passé au domicile est pratiquement au centre de gravité des activités Repos, Jardinage-bricolage, Repas chez soi, Télévision, qui est le media dominant en durée.

Tableau 1.2 - 5  
Coordonnées des variables supplémentaires  
(ou illustratives) sur les axes 1 à 3

VARIABLES	COORDONNEES		
	1	2	3
Autres activités	.08	.16	.04
Total Domicile	.67	-.50	-.21
Total Déplacement	-.72	.05	.14
Habitudes Cinema	-.87	-.11	-.14
Habitudes Radio.	-.27	-.57	.07
Habitudes Télévision	.04	-.55	.34
Habitudes Presse Quot	-.39	.01	-.70
Habitudes Presse mag	-.24	-.38	-.26
Habitudes Hebdo-News	-.46	.20	-.48

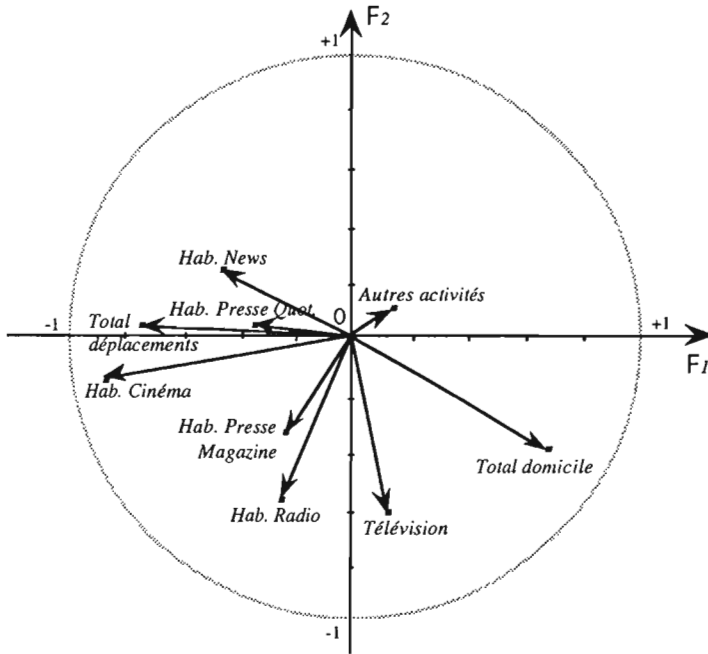


Figure 1.2 - 23  
Positionnement des variables supplémentaires  
(plan de la figure 1.2 - 22)

On présente le rôle de certaines caractéristiques socio-économiques, qui seront positionnées dans l'espace des individus. Les positions des individus dans le plan factoriel (tableau 1.2 - 6 et figure 1.2 - 24) vont permettre d'expliquer certaines des corrélations observées.

Ainsi, deux groupes (1133 et 1134) se distinguent à l'extrême gauche du premier axe : il s'agit de jeunes actifs instruits des grandes métropoles régionales ou de Paris, qui ont un profil d'activité typé (lecture, repas au restaurant, ...), expliquant à eux deux 35% de la variance le long de cet axe.

Le second groupe (1115 : jeunes peu instruits habitant dans des communes de profils variés) se distingue sur le deuxième axe (contribution de 26%). Remarquons aussi que ce même groupe a une distance à l'origine des axes (colonne DISTO, c'est-à-dire carré de la distance à l'origine) anormalement élevée (47.51) qui confirme son atypicité.

On vérifie sur le tableau de données 1.2 - 1 que ce groupe a un temps de travail moyen exceptionnellement faible (208.8, valeur qui est d'ailleurs le minimum de cette variable donné par le tableau 1.2 - 2) et des temps maxima pour "déplacement à pied" et "fréquentation média" (il s'agit essentiellement d'écoute télévision).

Souvent, dans les applications en vraie grandeur, les individus sont beaucoup plus nombreux et les identificateurs renvoient en général à un

numéro de questionnaire ou d'observation. Les variables nominales sont alors projetées selon la procédure indiquée au paragraphe 1.2.5.c.

**Tableau 1.2 - 6**  
**Coordonnées, contributions et cosinus carrés**  
**des individus sur les axes 1 et 2**

INDIVIDUS		COORDONNEES		CONTRIBUT.		COS. CARRE	
IDENTIF.	DISTO	1	2	1	2	1	2
1111	19.89	2.01	.85	3.8	.7	.20	.04
1115	47.51	2.26	-5.11	4.9	26.4	.11	.55
1121	10.55	-.71	1.01	.5	1.0	.05	.10
1122	13.29	-1.86	-.64	3.3	.4	.26	.03
1123	14.49	-1.28	-1.81	1.6	3.3	.11	.23
1124	19.06	-2.72	-2.93	7.1	8.7	.39	.45
1136	10.68	-.56	1.97	.3	3.9	.03	.36
1133	27.04	-4.21	-.30	17.0	.1	.66	.00
1134	25.35	-4.29	-.91	17.6	.8	.73	.03
2111	12.86	1.91	2.12	3.5	4.5	.28	.35
2112	17.27	1.43	-1.68	2.0	2.8	.12	.16
2117	10.89	1.03	-2.16	1.0	4.7	.10	.43
2121	10.96	1.27	2.55	1.5	6.6	.15	.59
2122	7.92	.62	-.21	.4	.0	.05	.01
2123	8.33	.30	-.33	.1	.1	.01	.01
2124	15.54	-.12	2.06	.0	4.3	.00	.27
2131	7.39	.55	2.03	.3	4.2	.04	.56
2132	24.45	-1.17	3.53	1.3	12.6	.06	.51
2133	7.85	-1.63	-.11	2.5	.0	.34	.00
2134	17.19	-2.54	1.36	6.2	1.9	.37	.11
3116	16.19	2.68	.96	6.9	.9	.45	.06
3117	15.96	2.43	-1.84	5.7	3.4	.37	.21
3121	13.00	1.90	2.11	3.4	4.5	.28	.34
3122	17.31	2.12	-.95	4.3	.9	.26	.05
3123	10.26	.56	-1.74	.3	3.1	.03	.30
3136	9.09	1.56	.09	2.3	.0	.27	.00
3137	21.68	-1.55	.08	2.3	.0	.11	.00

Le tableau 1.2 - 7 fournit les coordonnées des modalités (ou catégories) de ces variables qui sont, rappelons-le, les centres de gravité des individus concernés.

Ces centres de gravité ont été portés sur la figure 1.2 - 24 et les modalités contiguës d'une même variable nominale (il s'agit en fait de variables ordinales) ont été jointes par des lignes polygonales. Dans l'hypothèse où les groupes correspondant à une modalité particulière pourraient être considérés comme tirés au hasard parmi les 27 groupes, ces centres de gravité ne devraient pas s'éloigner beaucoup du centre de gravité du nuage (origine des axes factoriels).

On peut convertir cette distance au centre de gravité en "valeur-test"<sup>1</sup>, qui sera alors la réalisation d'une variable normale centrée réduite (deux premières colonnes du tableau 1.2 - 7).

<sup>1</sup> Ces aides à l'interprétation sont abordées dans un cadre plus général à l'occasion de l'analyse des correspondances multiples, au paragraphe 1.4.4.a.

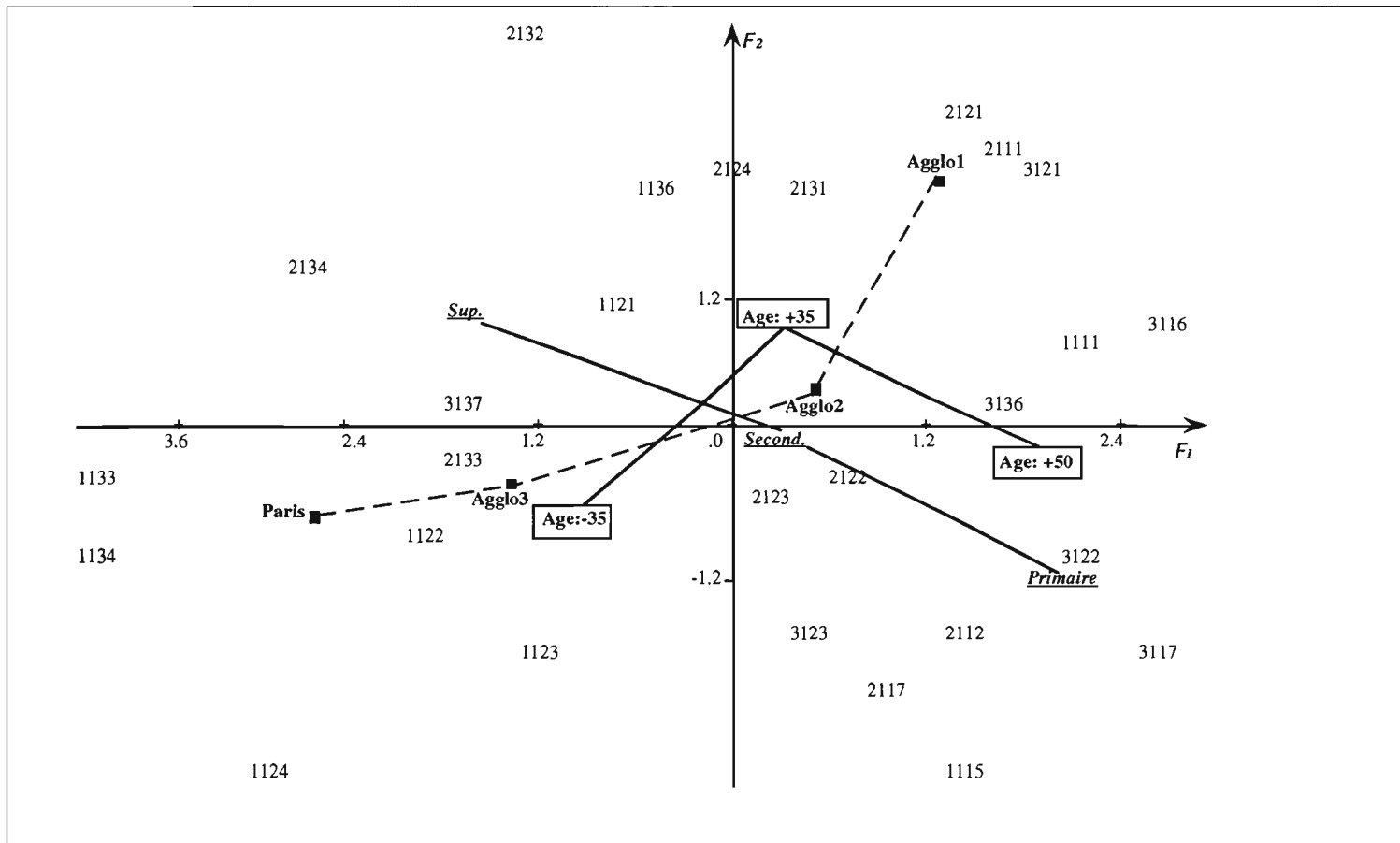


Figure 1.2 - 24 : Positionnement des individus (symboles à 4 chiffres) et des variables nominales (Age, niveau d'éducation, catégories d'agglomération)



**Tableau 1.2 - 7**  
**Valeurs-test et coordonnées des modalités**  
**supplémentaires sur les axes 1 et 2**

IDEN - LIBELLE	MODALITES EFFECT.	VALEURS-TEST		COORDONNEES	
		1	2	1	2
<b>. AGE</b>					
A-35 - Jeunes	9	-2.3	-1.6	-1.26	-.87
A+35 - Age-Moy	11	.3	1.8	.15	.83
A+50 - Ages	7	2.1	-.3	1.39	-.18
<b>. Education</b>					
prim - primaire	7	3.0	-1.5	1.96	-.98
seco - secondaire	11	.0	-.2	.01	-.08
supe - superieur	9	-2.8	1.6	-1.54	.86
<b>. Agglomération (EXTRAITS)</b>					
AGG1 - de 20 000	6	1.6	2.5	1.15	1.78
AGG2 - de 20 a 100 000	5	.3	.0	.23	.01
AGG3 - Plus de 100 000	5	-1.5	-1.1	-1.25	-.86
AGG4 - Paris	4	-2.6	-.1	-2.42	-.11

Autrement dit, dans l'hypothèse d'un tirage au hasard, la valeur-test d'une catégorie supplémentaire a 95 chances sur 100 d'être comprise dans l'intervalle  $[-1.96 \text{ et } +1.96]^1$ . Comme on le lit sur le tableau 1.2 - 7, la valeur-test du point "Paris" sur l'axe horizontal est de -2.6. C'est une modalité dont la position est significativement différente de l'origine.

La figure 1.2 - 24, tout comme le tableau 1.2 - 7, montrent que les trois variables nominales permettent surtout d'identifier le premier axe, opposant les jeunes instruits urbains aux personnes plus âgées et moins instruites. Seules les communes rurales (Agglo1) semblent liées au second axe.

Le lecteur de ces graphiques doit garder à l'esprit le fait qu'il s'agit ici d'identification passive par des variables nominales d'une analyse réalisée uniquement à partir des temps d'activité. Il ne s'agit pas d'une étude des liaisons existant entre ces variables nominales, même si certaines proximités peuvent paraître familières.

<sup>1</sup> Naturellement, l'intervalle de confiance précédent est trop étroit, car le test est répété sur plusieurs modalités ; il convient de ne le considérer que comme donnant un ordre de grandeur.

## Analyse des Correspondances

L'analyse des correspondances, présentée sous ce nom et développée par Benzécri (1969), a un certain nombre de précurseurs, parmi lesquels il faut citer Guttman (1941), Hayashi (1956).

Comme l'analyse en composantes principales, l'analyse des correspondances peut être présentée selon divers points de vue. Il est d'ailleurs difficile de faire l'historique précis de cette méthode. Les principes théoriques remontent probablement aux travaux de Fisher (1940) sur les tables de contingences, dans un cadre de statistique inférentielle classique. Depuis les travaux de Benzécri (1973) et de Escofier-Cordier (1965), on utilise surtout les propriétés algébriques et géométriques de l'outil descriptif que constitue l'analyse<sup>1</sup>. Cette méthode n'est pas un cas particulier de l'analyse en composantes principales bien que l'on puisse se ramener à cette technique en faisant des changements de variables appropriés (à condition de traiter chaque espace séparément). On peut la présenter comme un cas particulier de l'analyse canonique (cf. section 3.1) lorsque les données ont un codage "disjonctif" et également comme un cas particulier de l'analyse discriminante (cf. section 3.3). On peut enfin montrer qu'il s'agit de la recherche de la meilleure *représentation simultanée* de deux ensembles constituant les lignes et les colonnes d'un tableau de données (cf. § 1.3.3).

### 1.3.1 Domaine d'application

L'analyse des correspondances a un domaine d'application différent de l'analyse en composantes principales. Alors que l'on réserve cette dernière aux tableaux de mesures éventuellement hétérogènes et au traitement de variables numériques continues, l'analyse des correspondances est une méthode adaptée aux *tableaux de contingence* et permet d'étudier les éventuelles relations existant entre deux variables nominales. Nous verrons à la section suivante (1.4) qu'elle fournit, par extension, des descriptions satisfaisantes de certains tableaux de codages discontinus.

Le tableau de contingence (dit aussi de dépendance, ou tableau croisé) est obtenu en ventilant une population selon deux variables nominales.

---

<sup>1</sup> Les ancêtres les plus lointains de l'analyse des correspondances seraient, de façon tout à fait indépendante, Richardson et Kuder (1933) et Hirschfeld (1935). Les premiers auteurs visaient une meilleure sélection des vendeurs de la société *Procter and Gamble*, alors que le dernier étudiait une propriété de statistique mathématique. Cette variété de contextes est caractéristique de l'analyse des correspondances, méthode aussi utile en pratique que stimulante du point de vue théorique. Cf. les références historiques de Hill (1974), Benzécri (1982 a).

L'ensemble des colonnes du tableau désigne les modalités d'une variable et l'ensemble des lignes correspond à celles de l'autre variable. De ce fait, les lignes et les colonnes, qui désignent deux partitions d'une même population, jouent des rôles symétriques et sont traitées de façon analogue.

### 1.3.2 Démarche et principe : une introduction élémentaire

Nous allons utiliser, pour illustrer notre propos, une table de contingence de faible dimension pour laquelle le recours à l'analyse des correspondances ne se justifie pas vraiment, mais qui va permettre de présenter de façon simple les principes de cette méthode et les propriétés qui en découlent<sup>1</sup>. Bien que les lignes et les colonnes jouent un rôle similaire, nous conservons les mêmes notations que pour l'analyse générale.

Considérons le tableau de contingence  $K$  à  $n$  lignes et  $p$  colonnes obtenu en ventilant une population de 592 femmes suivant leurs couleurs des yeux et des cheveux.

**Tableau 1.3 - 1**  
**Tableau de contingence,**  
**répartition de 592 femmes suivant les couleurs des yeux et des cheveux.**

		couleur des cheveux				Total
		brun	châtain	roux	blond	
couleur des yeux	marron	68	119	26	7	220
	noisette	15	54	14	10	93
	vert	5	29	14	16	64
	bleu	20	84	17	94	215
Total		108	286	71	127	592

Source : Snee (1974)

En lignes est présentée la variable "couleur des yeux" à  $n = 4$  modalités (ou catégories) et en colonnes est donnée la variable "couleur des cheveux" à  $p = 4$  modalités.

A l'intersection d'une ligne et d'une colonne, nous avons le nombre  $k_{ij}$  de femmes ayant simultanément la couleur  $i$  des yeux et la couleur  $j$  de cheveux. Le total marginal  $k_i$  est le nombre de femmes ayant les yeux de couleur  $i$ , alors que le total marginal  $k_j$  est le nombre de femmes ayant les cheveux de couleur  $j$ .

On a les relations suivantes :

$$k_i = \sum_j k_{ij} \qquad k_j = \sum_i k_{ij} \qquad k = \sum_{i,j} k_{ij}$$

<sup>1</sup> Une présentation technique plus détaillée sera l'objet des paragraphes suivants de la même section.

qui, en termes de fréquences relatives, donnent lieu aux relations :

$$f_{ij} = \frac{k_{ij}}{k} \quad f_{i.} = \sum_j f_{ij} \quad f_{.j} = \sum_i f_{ij} \quad \sum_{i,j} f_{ij} = 1$$

Y-a-t-il indépendance entre la couleur des yeux et celle des cheveux ? Sinon quels types d'associations existent entre ces couleurs ?

### a – Transformations du tableau de contingence

Pour analyser un tableau de contingence, ce n'est pas le tableau d'effectifs bruts qui nous intéresse mais les tableaux des profils-lignes et celui des profils-colonnes c'est-à-dire les répartitions en pourcentage à l'intérieur d'une ligne ou d'une colonne.

On note les profils-lignes :  $\frac{f_{ij}}{f_{i.}} = \frac{k_{ij}}{k_{i.}}$

**Tableau 1.3 - 2**  
Profils-lignes (pourcentages-lignes arrondis)

		couleur des cheveux				total
		brun	châtain	roux	blond	
couleur des yeux	marron	31	54	12	3	100
	noisette	16	58	15	11	100
	vert	8	45	22	25	100
	bleu	9	39	8	44	100
profil moyen		18	48	12	22	100

et les profils-colonnes :  $\frac{f_{ij}}{f_{.j}} = \frac{k_{ij}}{k_{.j}}$

Le tableau 1.3 - 2 des profils-lignes (multipliés par 100) indique la répartition de la couleur des cheveux pour chaque modalité de couleur des yeux. Ce sont en somme les probabilités conditionnelles d'avoir les cheveux de la couleur  $j$  sachant que les yeux ont la couleur  $i$ . Cette répartition sur l'ensemble de la population étudiée donne le profil moyen :

$$f_{.j} = \frac{k_{.j}}{k}$$

**Tableau 1.3 - 3**  
Profils-colonnes (pourcentages-colonnes arrondis)

		couleur des cheveux				profil moyen
		brun	châtain	roux	blond	
couleur des yeux	marron	63	42	37	6	37
	noisette	14	19	20	8	16
	vert	5	10	20	13	11
	bleu	19	29	24	74	36
total		100	100	100	100	100

Le tableau 1.3 - 3 des profils-colonnes (multipliés par 100) fournit la répartition de la couleur des yeux suivant chaque modalité de couleur des cheveux et le profil moyen de la couleur des yeux :

$$f_{i.} = \frac{k_i}{k}$$

### b – Hypothèse d'indépendance

On s'intéresse aux liens éventuels entre couleurs des yeux et des cheveux.

On sait qu'il y a indépendance entre deux variables aléatoires  $i$  et  $j$  prenant leurs valeurs sur deux ensembles de tailles  $n$  et  $p$ , dont la loi jointe est  $p_{ij}$  et les lois marginales  $p_{i.}$  et  $p_{.j}$ , si pour tout  $i$  et pour tout  $j$  on a (avec les notations usuelles) :

$$p_{ij} = p_{i.} p_{.j}$$

La traduction de cette relation en termes d'estimations empiriques est la suivante :

$$f_{ij} = f_{i.} f_{.j}$$

Naturellement, même sous l'hypothèse d'indépendance, une telle relation n'est qu'approximativement vraie. Le classique test du  $\chi^2$  de Karl Pearson pour les tables de contingence permet précisément d'apprécier l'écart entre les lois empiriques  $f_{ij}$  et  $f_{i.} f_{.j}$ .

Consultons le tableau 1.3 - 4 des fréquences observées  $f_{ij}$  qui n'est autre que la tableau 1.3 - 1 divisé par sa somme (592) et multiplié par 100 pour plus de lisibilité.

**Tableau 1.3 - 4**  
**Tableau de fréquences observées**

		couleur des cheveux				profil moyen
		brun	châtain	roux	blond	
couleur des yeux	marron	11	20	4	1	37
	noisette	3	9	2	2	16
	vert	1	5	2	3	11
	bleu	3	14	3	16	36
profil moyen		18	48	12	21	100

Parmi les 37% de femmes aux yeux marrons par exemple, on devrait observer, sous l'hypothèse d'indépendance, 18% de femmes brunes (ce qui ferait alors 7% de l'ensemble des femmes, au lieu des 11% réellement observés), 48% aux cheveux châains (ce qui ferait 18% au lieu de 20%), etc.

Construisons le tableau de "fréquences théoriques"  $f_{i.} f_{.j}$  sous l'hypothèse d'indépendance (cf. tableau [1.3 - 5]) :

**Tableau 1.3 - 5**  
**Tableau de fréquences théoriques**

		couleur des cheveux				profil moyen
		brun	châtain	roux	blond	
couleur des yeux	marron	7	18	4	8	37
	noisette	3	8	2	3	16
	vert	2	5	1	2	11
	bleu	7	18	12	8	36
profil moyen		18	48	12	21	100

Cette hypothèse s'exprime aussi sur les profils-lignes. En effet, il en découle que, quelque soit  $j$  :

$$\frac{f_{ij}}{f_{i.}} = f_{.j}$$

Si tous les profils "couleurs des yeux" sont identiques entre eux, et par conséquent identiques au profil moyen correspondant, il y a indépendance entre les couleurs des yeux et celles de cheveux puisque la connaissance d'une couleur des yeux ne change pas la répartition de la couleur des cheveux.

Il en est de même pour les profils-colonnes où, quelque soit  $i$  :

$$\frac{f_{ij}}{f_{.j}} = f_{i.}$$

Ainsi, examiner les proximités entre les profils revient à examiner la proximité entre chaque profil et son profil moyen, ce qui permet d'étudier la liaison entre deux variables nominales, c'est-à-dire l'écart à l'indépendance. Sur un tableau de dimension importante, la lecture directe des profils-lignes et des profils-colonnes est difficile, ainsi que la comparaison de ces profils avec leur profil moyen.

Nous allons voir comment la construction du nuage, le choix du critère d'ajustement et celui de la distance, s'imposent de par la nature même des données analysées.

### c – Construction des nuages

Pour l'analyse d'un tableau de contingence, nous raisonnerons en termes de profils, ce qui permet de rendre comparables les modalités d'une même variable. Les proximités entre les points s'interpréteront en terme de similitude.

#### - Nuage des $n$ lignes

L'ensemble des profils-lignes forme un nuage de  $n$  points dans l'espace des  $p$  colonnes et représente ici le nuage des 4 modalités de couleurs des yeux. Chaque point  $i$  a pour coordonnées dans  $\mathbb{R}^p$  :

$$\left\{ \frac{f_{ij}}{f_{i.}} ; j = 1, 2, \dots, p \right\}$$

Il est affecté d'une masse  $f_i$  qui est sa fréquence relative.

Puisque  $\sum_{j=1}^p \frac{f_{ij}}{f_i} = 1$ , les  $n$  points du nuage sont situés dans un sous-espace à  $p-1$  dimensions.

Le centre de gravité de ce nuage est la moyenne des profils-lignes affectés de leurs masses et correspond au profil moyen, c'est-à-dire au profil de la couleur des cheveux sur l'ensemble de la population. Sa  $j^{\text{ème}}$  composante vaut :

$$\sum_{i=1}^n f_i \frac{f_{ij}}{f_i} = f_{.j}$$

C'est la fréquence marginale des colonnes.

#### - Nuage des $p$ colonnes

De la même façon, l'ensemble des  $p$  profils-colonnes constitue un nuage de  $p$  points dans l'espace des  $n$  lignes et représente ici le nuage des 4 modalités de couleur des cheveux.

Les coordonnées dans  $\mathbb{R}^n$  du point  $j$  sont données par :

$$\left\{ \frac{f_{ij}}{f_{.j}} ; i = 1, 2, \dots, n \right\}$$

Chaque point est affecté d'une masse  $f_{.j}$ .

Les  $p$  points du nuage sont situés dans un sous-espace à  $n-1$  dimensions

puisque  $\sum_{i=1}^n \frac{f_{ij}}{f_{.j}} = 1$ .

Le centre de gravité du nuage des profils-colonnes est le profil moyen de la couleur des yeux. Sa  $i^{\text{ème}}$  composante vaut :

$$\sum_{j=1}^p f_{.j} \frac{f_{ij}}{f_{.j}} = f_i$$

C'est la fréquence marginale des lignes.

### d – Critère d'ajustement

On cherche à représenter géométriquement les similitudes entre les différentes modalités d'une même variable, ce qui nous conduit à représenter les proximités entre les profils et le profil moyen défini sur l'ensemble de la population<sup>1</sup>. Ceci nous amène, comme en analyse en

<sup>1</sup> Un nuage de points concentré autour de son centre de gravité est un nuage dont les points-profil sont proches du profil moyen, et donc traduira une certaine indépendance entre les deux variables nominales.

composantes principales dans le cas des points-individus, à considérer le nuage de points centré sur son centre de gravité.

Dans la construction des nuages de  $\mathbb{R}^p$  et de  $\mathbb{R}^n$  (cf. tableaux 1.3 - 2 et 1.3 - 3), le choix des profils comme coordonnées donne à toutes les modalités de couleur des yeux et celles de cheveux la même importance. L'importance est cependant restituée au travers de la masse affectée à chaque point (proportionnelle à sa fréquence), afin de ne pas privilégier les classes d'effectifs faibles et de respecter la répartition réelle de la population. Cette masse interviendra d'une part lors du calcul des coordonnées du centre de gravité du nuage et d'autre part dans le critère d'ajustement.

Pour le calcul de l'ajustement, la quantité à rendre maximale sera donc la somme pondérée des carrés des distances entre les points et le centre de gravité du nuage (c'est-à-dire l'inertie de la droite d'allongement maximum du nuage) en utilisant une distance entre profils qu'il reste à définir.

### e – Choix des distances

La distance euclidienne usuelle entre deux points-lignes exprimée sur le tableau d'effectifs bruts ne ferait que traduire les différences d'effectifs entre deux modalités de couleurs des yeux. En revanche, la distance euclidienne usuelle entre deux profils-lignes traduit bien la ressemblance ou la différence entre les deux couleurs des yeux sans tenir compte des effectifs totaux de ces modalités :

$$d^2(i, i') = \sum_{j=1}^p \left( \frac{f_{ij}}{f_{i.}} - \frac{f_{i'j}}{f_{i'.}} \right)^2$$

Cependant, cette distance favorise les colonnes qui ont une masse  $f_{.j}$  importante c'est-à-dire les couleurs de cheveux qui sont bien représentées dans la population étudiée.

Pour remédier à cela, et aussi pour d'autres propriétés qui seront développées ci-dessous, on pondère chaque écart par l'inverse de la masse de la colonne et l'on calcule une nouvelle distance appelée<sup>1</sup> la distance du  $\chi^2$  :

$$d^2(i, i') = \sum_{j=1}^p \frac{1}{f_{.j}} \left( \frac{f_{ij}}{f_{i.}} - \frac{f_{i'j}}{f_{i'.}} \right)^2 \quad [1.3 - 1]$$

On définit de la même manière la distance entre les profils-colonnes par :

$$d^2(j, j') = \sum_{i=1}^n \frac{1}{f_{i.}} \left( \frac{f_{ij}}{f_{.j}} - \frac{f_{i'j'}}{f_{.j'}} \right)^2 \quad [1.3 - 2]$$

---

<sup>1</sup> L'inertie totale des nuages de points lignes (ou de points colonnes) calculée avec cette distance est proportionnelle au classique  $\chi^2$  de Karl Pearson utilisé pour éprouver l'indépendance des lignes et des colonnes d'une table de contingence. D'où le nom de distance du  $\chi^2$ .



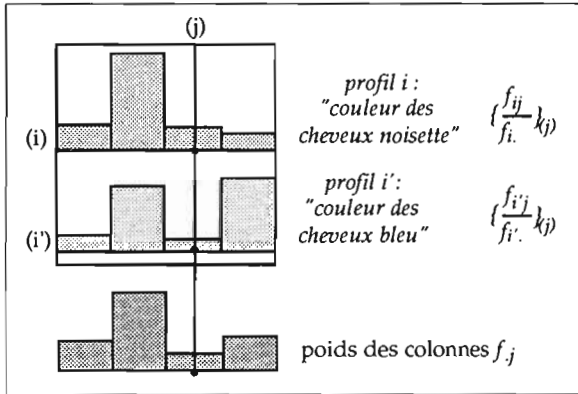


Figure 1.3 - 1  
Distance du  $\chi^2$

C'est cette distance pondérée, ainsi que le rôle symétrique joué par les lignes et les colonnes du tableau de contingence, qui particularisent l'analyse des correspondances et lui assurent des propriétés remarquables que ne possède pas l'analyse en composantes principales : l'équivalence distributionnelle et les relations de transition.

### f - Equivalence distributionnelle

La propriété d'équivalence distributionnelle permet d'agréger deux modalités d'une même variable ayant des profils identiques en une nouvelle modalité affectée de la somme de leurs masses, sans rien changer, ni aux distances entre les modalités de cette variable, ni aux distances entre les modalités de l'autre variable.

Si par exemple les deux profils-lignes  $i'$  et  $i''$  sont identiques dans  $\mathbb{R}^p$ , on les agrège en un profil-ligne  $i$  dont la masse sera la somme des fréquences des deux profils  $i'$  et  $i''$ . Les deux points  $i'$  et  $i''$  étant confondus cela ne modifie pas la configuration du nuage de points dans  $\mathbb{R}^p$ .

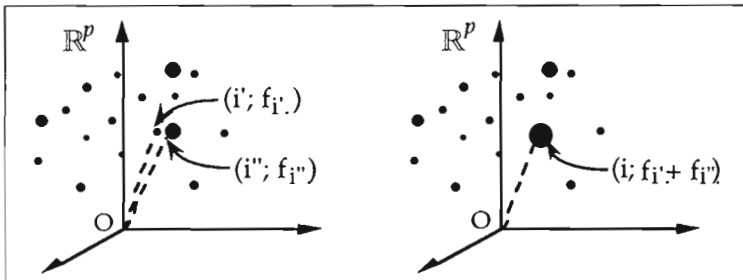


Figure 1.3 - 2  
Equivalence distributionnelle : points-lignes confondus

Mais surtout, les distances entre colonnes restent inchangées. Il en est de même pour des profils-colonnes dans  $\mathbb{R}^n$  ayant les mêmes propriétés.

Cette propriété est fondamentale puisqu'elle garantit une certaine invariance des résultats vis-à-vis de la nomenclature choisie pour la construction des modalités d'une variable, sous condition de regrouper des modalités aux profils semblables.

On ne perd pas d'information en agrégeant certaines classes et l'on n'en gagne pas en subdivisant des classes homogènes.

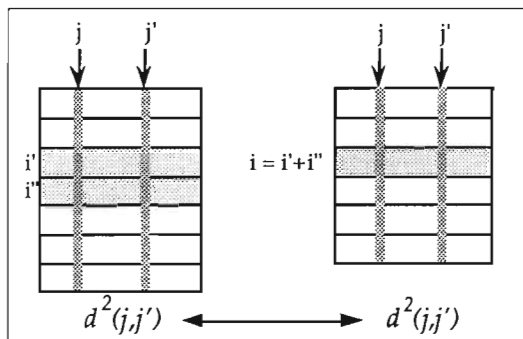


Figure 1.3 - 3  
Equivalence distributionnelle : invariance des distances entre colonnes vis-à-vis de l'agrégation des lignes

Prenons le cas de deux tables de contingences issues du recensement de la population, l'une croisant professions et départements, l'autre professions et régions. Sous l'hypothèse d'homogénéité des départements d'une même région par rapport aux professions, il sera équivalent de réaliser l'analyse des correspondances sur les départements et sur les régions. Les configurations du nuage des professions, pour les deux analyses, seront semblables (voir la démonstration au § 1.3.3.a).

### g – Relations de transition ou quasi-barycentriques

Une des caractéristiques de l'analyse des correspondances est l'existence de relations de type barycentrique qui lient graphiquement les deux variables représentées en ligne et en colonne.

L'idée est simple et revient à représenter les histogrammes des profils-colonnes dans le nuage des profils-lignes et réciproquement.

Supposons fixé le nuage des couleurs des yeux (nuage des profils-lignes) dans un espace à 2 dimensions comme représenté sur la figure 1.3 - 4. Le centre du graphique représente le profil moyen (la distribution marginale) des couleurs des yeux.

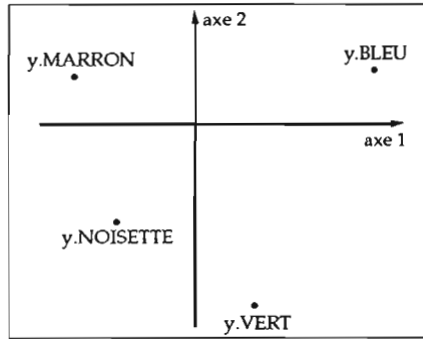


Figure 1.3 - 4  
Nuage des couleurs des yeux

Considérons maintenant l'histogramme décrivant le profil des cheveux bruns suivant la couleur de yeux (cf. tableau 1.3 - 3 des profils-colonnes) représenté figure 1.3 - 5.

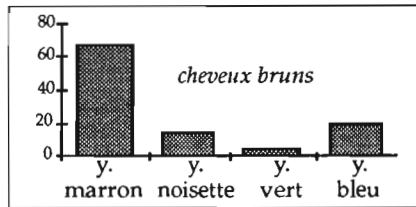


Figure 1.3 - 5  
Histogramme des cheveux bruns

Cet histogramme va permettre de positionner le point-colonne "cheveux bruns" dans le nuage des points-lignes (le nuage des couleurs des yeux) : chaque point  $i$  représentant une couleur des yeux est pondéré par sa fréquence relative telle qu'elle est décrite par l'histogramme.

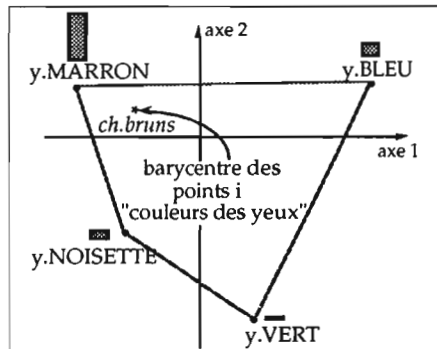


Figure 1.3 - 6  
Position du point "cheveux bruns"  
comme barycentre des points "couleurs des yeux"

On construit ainsi le barycentre de ces points qui correspond au point "cheveux bruns". Il est contenu dans une enveloppe convexe constituée par l'ensemble des points pondérés (cf. figure 1.3 - 6). Cette modalité sera attirée par les yeux marrons, compte tenu de sa masse plus élevée. Elle sera par contre éloignée des yeux verts.

Chaque point  $j$  "couleur des cheveux" est ainsi un barycentre particulier des points  $i$  "couleur des yeux", le point  $i$  étant affecté de la masse "part de la couleur  $i$  des yeux sachant que la couleur des cheveux est  $j$ ", (c'est-à-dire le profil-colonne  $f_{ij}/f_{.j}$ ).

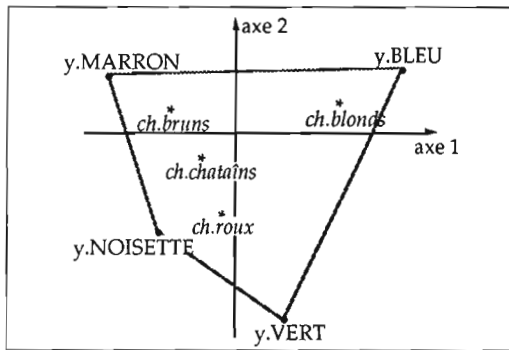


Figure 1.3 - 7

**Représentation des points "couleurs des yeux" et positionnement des points "couleurs des cheveux" en barycentres**

Si l'on considère maintenant le nuage des profils-colonnes, c'est-à-dire le nuage des couleurs des cheveux, il est naturel de procéder de la même façon et de représenter l'histogramme de chaque couleur des yeux dans ce nuage.

On positionne donc chaque point-ligne  $i$  "couleur des yeux" comme barycentre des points  $j$  "couleurs des cheveux" pondérés par la part de la couleur  $j$  des cheveux dans la couleur  $i$  des yeux, donnée par les profils-lignes  $\{f_{ij}/f_{.i}\}$  (cf. figure 1.3 - 8).

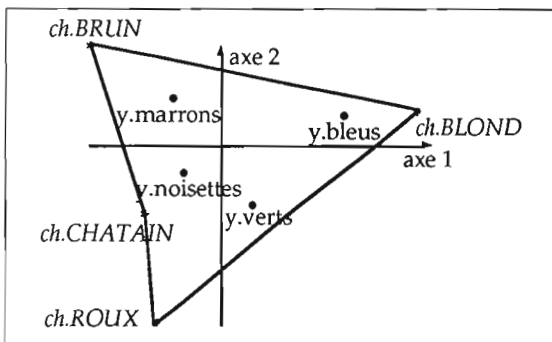


Figure 1.3 - 8

**Représentation des points "couleurs des cheveux" et positionnement des points "couleurs des yeux" en barycentres**

Les relations barycentriques vont justifier et donner un sens à la représentation simultanée des deux nuages définis dans les deux espaces.

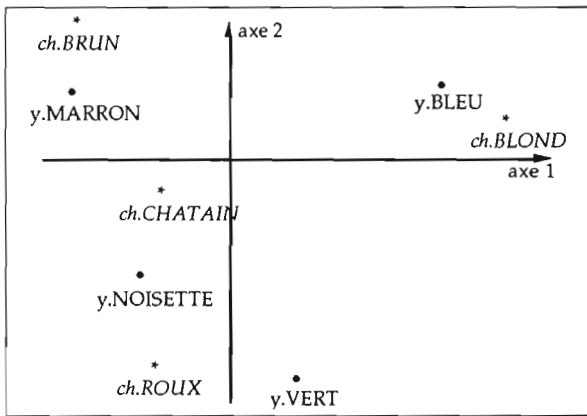
### h – Justification de la représentation simultanée

D'après le schéma de l'analyse générale, on pourrait envisager l'analyse des deux nuages de points de manière indépendante et l'interpréter comme une analyse en composantes principales où toute l'information entre les deux nuages transite par les facteurs de mêmes rangs. Compte tenu des relations barycentriques, il en est autrement en analyse des correspondances.

Ces relations montrent qu'il existe une possibilité de représentation particulière<sup>1</sup> : il est possible de positionner chaque point d'un nuage parmi l'ensemble des points de l'autre nuage.

Ainsi, dans le nuage des profils-lignes, chaque profil-colonne est au barycentre des points du nuage. Projeté sur un plan, nous disposons d'une première représentation simultanée (cf. figure 1.3 - 7). De même, chaque profil-ligne est barycentre de l'ensemble des profils-colonnes et constitue, avec les axes de mêmes rangs, une deuxième représentation simultanée (cf. figure 1.3 - 8).

Mais nous voulons une seule représentation simultanée des deux nuages de points et la situation idéale serait de les superposer.



**Figure 1.3 - 9**  
Représentation simultanée;  
Relations quasi-barycentriques

Ceci est *a priori* impossible par définition même du barycentre puisque chaque ensemble devrait alors être contenu dans l'autre. Il est cependant possible de forcer cette représentation en dilatant (sur chaque axe) les centres de gravité (figure 1.3 - 9). On pourra alors représenter sur de mêmes axes (et

<sup>1</sup> Cette possibilité est due au fait que les coordonnées d'origine (les profils) sont des nombres positifs dont la somme vaut 1.

donc sur un même plan) l'ensemble des lignes et des colonnes afin d'approcher au mieux la situation idéale. Les relations seront quasi-barycentriques (cf. § 1.3.3).

Les yeux bleus s'associent aux cheveux blonds, les yeux marrons aux cheveux bruns. Les cheveux roux sont attirés par les yeux noisettes et verts qui les caractérisent. La catégorie des cheveux châtain est assez proche de l'origine du plan représentant le profil moyen et n'est spécifique d'aucune couleur des yeux<sup>1</sup>.

Nous verrons que le déroulement de l'analyse des correspondances, compte tenu des rôles symétriques des lignes et des colonnes du tableau de contingence et des propriétés de la distance du  $\chi^2$ , aboutit naturellement aux relations barycentriques (à un coefficient près qui est le coefficient de dilatation permettant la représentation simultanée unique).

### 1.3.3 Schéma général de l'analyse des correspondances

L'analyse des correspondances revient à effectuer l'analyse générale d'un nuage de points pondérés dans un espace muni de la métrique du  $\chi^2$ . On fera donc référence à l'analyse générale avec des métriques et des critères quelconques (cf. § 1.1.6.a).

#### a – Géométrie des nuages et éléments de base

Contrairement à l'analyse en composantes principales, le tableau de données subit deux transformations, l'une en profils-lignes, l'autre en profils-colonnes, à partir desquelles vont être construits les nuages de points dans  $\mathbb{R}^p$  et dans  $\mathbb{R}^n$  (figure 1.3 - 10).

Pour faire le lien avec l'analyse générale (cf. section 1.1), nous adopterons des notations matricielles (figure 1.3 - 11).

Les transformations opérées sur le tableau des données peuvent s'écrire à partir des trois matrices  $F$ ,  $D_n$  et  $D_p$  qui définissent les éléments de base de l'analyse.

$F$  d'ordre  $(n,p)$  désigne le tableau des fréquences relatives ;  $D_n$  d'ordre  $(n,n)$  est la matrice diagonale dont les éléments diagonaux sont les marges en lignes  $f_i$  ;  $D_p$  est la matrice diagonale d'ordre  $(p,p)$  des marges en colonnes  $f_j$ .

---

<sup>1</sup> On dispose le plus souvent d'un tableau de données de dimension importante et la représentation du nuage des points non dilaté et des barycentres correspondants, dans un des espaces, fournit un graphique confus puisque les barycentres seront souvent rassemblés près de l'origine du plan. Une seule représentation simultanée, la représentation dite quasi-barycentrique, du fait de la dilatation des nuages de points qu'elle nécessite, offre l'avantage d'une lecture plus facile du graphique.

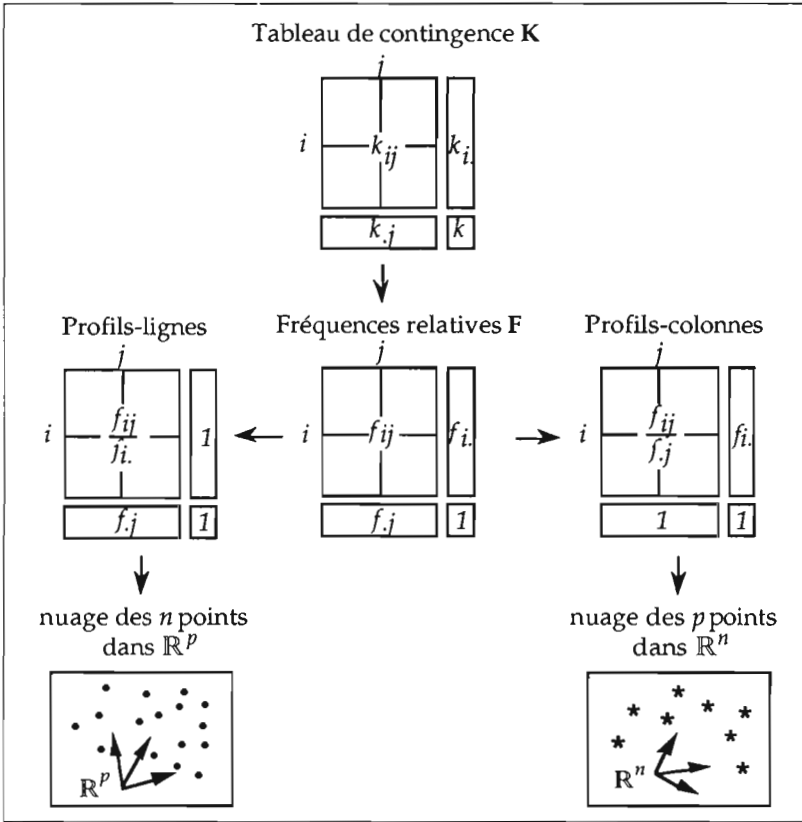


Figure 1.3 - 10 Transformations du tableau de contingence

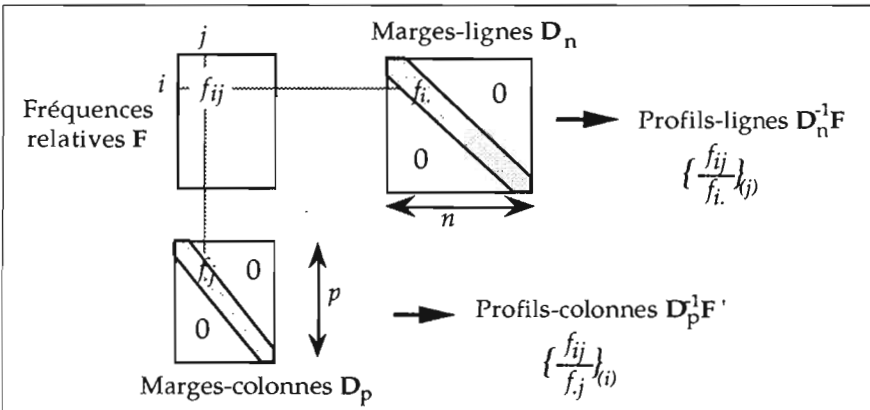


Figure 1.3 - 11 Fréquences, marges, profils

Les deux nuages de points (dans l'espace des colonnes et dans l'espace des lignes) sont construits de manière analogue. Nous récapitulons ici les éléments de base de l'analyse qui vont permettre la construction des facteurs.

Tableau 1.3 - 6  
Les éléments de base de l'analyse : récapitulation

Nuage de $n$ points-lignes dans l'espace $\mathbb{R}^p$	← Eléments → de base	Nuage de $p$ points-colonnes dans l'espace $\mathbb{R}^n$
$X = D_n^{-1}F$ <p><math>p</math> coordonnées (point-ligne <math>i</math>)</p> $\frac{f_{ij}}{f_i}, \text{ pour } j=1, 2, \dots, p.$	Analyse du tableau $X$	$X = D_p^{-1}F'$ <p><math>n</math> coordonnées (point-colonne <math>j</math>)</p> $\frac{f_{ij}}{f_j}, \text{ pour } i=1, 2, \dots, n.$
$M = D_p^{-1}$ $d^2(i, i') = \sum_{j=1}^p \frac{1}{f_j} \left( \frac{f_{ij}}{f_i} - \frac{f_{i'j}}{f_i'} \right)^2$	avec la métrique $M$	$M = D_n^{-1}$ $d^2(j, j') = \sum_{i=1}^n \frac{1}{f_i} \left( \frac{f_{ij}}{f_j} - \frac{f_{ij'}}{f_j'} \right)^2$
$N = D_n$ <p>masse du point <math>i</math>: <math>f_i</math></p>	et le critère $N$	$N = D_p$ <p>masse du point <math>j</math>: <math>f_j</math></p>

### Remarques

- 1) La matrice  $N$  des masses dans un espace est liée à la métrique  $M$  utilisée dans l'autre espace.
- 2) Il existe une différence fondamentale avec l'analyse en composantes principales : les transformations faites sur les données brutes dans les deux espaces sont identiques (car les ensembles mis en correspondance jouent des rôles analogues). Elles correspondent à des transformations analytiques différentes : le tableau des nouvelles coordonnées dans l'espace des colonnes n'est pas le simple transposé de celui des nouvelles coordonnées dans l'espace des lignes. En composantes principales, des transformations très différentes conduisaient à une même formule analytique.

### Démonstration de l'équivalence distributionnelle

La distance du  $\chi^2$  a pour effet d'accorder une même importance, d'une part aux colonnes quelles que soient leurs fréquences relatives dans le calcul de la distance entre deux profils-lignes, et d'autre part aux lignes s'il s'agit du calcul de la distance entre profils-colonnes.

Elle offre l'avantage de vérifier le principe d'équivalence distributionnelle<sup>1</sup> (cf. figure 1.3 - 2). Ce principe assure la robustesse des résultats de l'analyse

<sup>1</sup> La distance euclidienne usuelle entre profils ne possède pas la propriété d'équivalence distributionnelle, mais d'autres distances possèdent cette propriété (cf. Escofier, 1978).



des correspondances vis à vis de l'arbitraire du découpage en modalités des variables nominales. Il s'exprime de la façon suivante dans  $\mathbb{R}^p$ :

si deux points-lignes  $i_1$  et  $i_2$  sont confondus dans  $\mathbb{R}^p$ , on a pour tout  $j$  :

$$\frac{f_{i_1j}}{f_{i_1.}} = \frac{f_{i_2j}}{f_{i_2.}} = \frac{f_{i_0j}}{f_{i_0.}} \quad [1.3 - 3]$$

On a en particulier :

$$\frac{f_{i_1j} + f_{i_2j}}{f_{i_1.} + f_{i_2.}} = \frac{f_{i_0j}}{f_{i_0.}}$$

D'où, puisque les dénominateurs sont égaux, on a pour tout  $j$  :

$$f_{i_1j} + f_{i_2j} = f_{i_0j}$$

Les calculs des quantités  $f_{.j} = \sum_i f_{ij}$  ne sont donc pas affectés et les distances  $d^2(i, i')$  données par la formule [1.3 - 1] sont invariantes.

Montrons maintenant que les distances entre colonnes ne changent pas. La distance  $d^2(j, j')$  donnée par la formule [1.3 - 2] contient entre autres les deux termes  $A(i_1)$  et  $A(i_2)$  correspondant aux indices  $i_1$  et  $i_2$  :

$$A(i_1) + A(i_2) = \frac{1}{f_{i_1.}} \left\{ \frac{f_{i_1j}}{f_{.j}} - \frac{f_{i_1j'}}{f_{.j'}} \right\}^2 + \frac{1}{f_{i_2.}} \left\{ \frac{f_{i_2j}}{f_{.j}} - \frac{f_{i_2j'}}{f_{.j'}} \right\}^2$$

Ces deux termes sont remplacés par un seul terme  $A(i_0)$  tel que :

$$A(i_0) = \frac{1}{f_{i_0.}} \left\{ \frac{f_{i_0j}}{f_{.j}} - \frac{f_{i_0j'}}{f_{.j'}} \right\}^2$$

Remarquons par exemple que :

$$A(i_1) = \frac{1}{f_{i_1.}} \left\{ \frac{f_{i_1j}}{f_{i_1.} f_{.j}} - \frac{f_{i_1j'}}{f_{i_1.} f_{.j'}} \right\}^2$$

$A(i_1)$  et  $A(i_2)$  s'écrivent de la même façon et les quantités entre accolades sont égales, d'après la relation [1.3 - 3], à un même nombre que l'on notera  $B$ . On a donc :

$$A(i_1) + A(i_2) = f_{i_1.} B + f_{i_2.} B = f_{i_0.} B = A(i_0)$$

D'où l'invariance de la distance  $d^2(j, j')$ .

## b – Critère à maximiser et matrice à diagonaliser

Nous voulons représenter graphiquement les proximités entre profils. Nous nous plaçons donc, dans les deux espaces, aux centres de gravité des nuages. Cependant, et c'est là une des particularités de l'analyse des correspondances, il est équivalent de procéder à l'analyse par rapport à

l'origine ou par rapport aux centres de gravité, à condition de négliger dans le premier cas l'axe factoriel qui joint l'origine au centre de gravité<sup>1</sup>.

Nous commencerons par effectuer l'analyse générale par rapport à l'origine, l'expression des formules étant plus simple, puis nous montrerons, au paragraphe 1.3.7, l'équivalence avec l'analyse effectuée par rapport aux centres de gravité.

Plaçons-nous dans l'espace des colonnes<sup>2</sup>  $\mathbb{R}^p$  et cherchons l'axe d'inertie maximum du nuage des points-lignes passant par l'origine  $O$  et engendré par un vecteur-unitaire  $\mathbf{u}$  pour la métrique  $\mathbf{D}_p^{-1}$ . Ceci nous amène à maximiser la somme pondérée des carrés des projections sur l'axe (cf. § 1.1.1) c'est-à-dire :

$$\text{Max}_{\mathbf{u}} \left\{ \sum_i f_i \cdot d^2(i, O) \right\}$$

et à rendre maximale la quantité :

$$\mathbf{u}' \mathbf{D}_p^{-1} \mathbf{F}' \mathbf{D}_n^{-1} \mathbf{F} \mathbf{D}_p^{-1} \mathbf{u}$$

avec la contrainte :

$$\mathbf{u}' \mathbf{D}_p^{-1} \mathbf{u} = 1$$

$\mathbf{u}$  est vecteur propre de la matrice :

$$\mathbf{S} = \mathbf{F}' \mathbf{D}_n^{-1} \mathbf{F} \mathbf{D}_p^{-1}$$

associé à la plus grande valeur propre  $\lambda$  différente de 1.

La matrice à diagonaliser est donc la matrice  $\mathbf{S}$  de terme général :

$$s_{jj'} = \sum_{i=1}^n \frac{f_{ij} f_{ij'}}{f_i \cdot f_{j'}}$$

De la même façon, on doit rendre maximum dans  $\mathbb{R}^n$ , la quantité :

$$\mathbf{v}' \mathbf{D}_n^{-1} \mathbf{F} \mathbf{D}_p^{-1} \mathbf{F}' \mathbf{D}_n^{-1} \mathbf{v}$$

avec la contrainte :

$$\mathbf{v}' \mathbf{D}_n^{-1} \mathbf{v} = 1$$

$\mathbf{v}$  est vecteur propre de la matrice :

$$\mathbf{T} = \mathbf{F} \mathbf{D}_p^{-1} \mathbf{F}' \mathbf{D}_n^{-1}$$

<sup>1</sup> Cet axe est associé à la valeur propre égale à 1, appelée valeur propre triviale.

<sup>2</sup> Compte tenu de la symétrie du tableau de contingence, les démonstrations dans l'autre espace se déduisent par permutation des indices  $i$  et  $j$  (c'est-à-dire transposition de  $\mathbf{F}$  et permutation des matrices  $\mathbf{D}_p$  et  $\mathbf{D}_n$ ).

### c – Axes factoriels et facteurs

Nous supposons ici que  $p$  correspond à la plus petite dimension du tableau de données. Après avoir écarté la valeur propre triviale égale à 1 et le vecteur propre associé, nous retenons, de la diagonalisation de la matrice, les  $p-1$  valeurs propres non nulles et les vecteurs propres associés. Nous obtenons ainsi au plus  $p-1$  axes factoriels.

Tableau 1.3 - 7  
Eléments de construction de l'analyse

Dans $\mathbb{R}^p$	← Eléments de construction →	Dans $\mathbb{R}^n$
$S = F' D_n^{-1} F D_p^{-1}$	Matrice à diagonaliser	$T = F D_p^{-1} F' D_n^{-1}$
$S u_\alpha = \lambda_\alpha u_\alpha$	Axe factoriel	$T v_\alpha = \lambda_\alpha v_\alpha$
$\psi_\alpha = D_n^{-1} F D_p^{-1} u_\alpha$ $\psi_{\alpha i} = \sum_{j=1}^p \frac{f_{ij}}{f_{i.} f_{.j}} u_{\alpha j}$	Coordonnées factorielles	$\varphi_\alpha = D_p^{-1} F' D_n^{-1} v_\alpha$ $\varphi_{\alpha j} = \sum_{i=1}^n \frac{f_{ij}}{f_{i.} f_{.j}} v_{\alpha i}$

Les coordonnées factorielles sont centrées :

$$\sum_{i=1}^n f_{i.} \psi_{\alpha i} = \sum_{j=1}^p f_{.j} \varphi_{\alpha j} = 0 \quad [1.3 - 4]$$

et de variance égale à  $\lambda_\alpha$  :

$$\sum_{i=1}^n f_{i.} \psi_{\alpha i}^2 = \sum_{j=1}^p f_{.j} \varphi_{\alpha j}^2 = \lambda_\alpha \quad [1.3 - 5]$$

### d – Relation entre les deux espaces

L'analyse générale a montré que les matrices  $S$  et  $T$  ont les mêmes valeurs propres non nulles  $\lambda_\alpha$  et qu'entre le vecteur propre unitaire  $u_\alpha$  de  $S$  associé à  $\lambda_\alpha$  et le vecteur propre unitaire  $v_\alpha$  de  $T$  relatif à la même valeur propre, il existe les relations dites de transition :

$$\left\{ \begin{array}{l} v_\alpha = \frac{1}{\sqrt{\lambda_\alpha}} F D_p^{-1} u_\alpha \\ u_\alpha = \frac{1}{\sqrt{\lambda_\alpha}} F' D_n^{-1} v_\alpha \end{array} \right. \quad [1.3 - 6]$$

$$[1.3 - 7]$$

La comparaison de ces relations avec les expressions des coordonnées factorielles :

$$\psi_\alpha = D_n^{-1} F D_p^{-1} u_\alpha \quad [1.3 - 8]$$

et

$$\varphi_\alpha = D_p^{-1} F' D_n^{-1} v_\alpha \quad [1.3 - 9]$$

montre que celles-ci sont liées aux composantes des axes de l'autre espace par les formules :

$$\begin{cases} \psi_{\alpha} = \sqrt{\lambda_{\alpha}} \mathbf{D}_n^{-1} \mathbf{v}_{\alpha} & [1.3 - 10] \\ \varphi_{\alpha} = \sqrt{\lambda_{\alpha}} \mathbf{D}_p^{-1} \mathbf{u}_{\alpha} & [1.3 - 11] \end{cases}$$

C'est-à-dire, explicitement :

$$\begin{cases} \psi_{\alpha i} = \frac{\sqrt{\lambda_{\alpha}}}{f_i} v_{\alpha i} \\ \varphi_{\alpha j} = \frac{\sqrt{\lambda_{\alpha}}}{f_j} u_{\alpha j} \end{cases}$$

### e – Relations de transition (ou quasi-barycentriques)

Les substitutions dans la relation [1.3 - 9] de  $\mathbf{v}_{\alpha}$  par sa valeur tirée de [1.3 - 10] et dans la relation [1.3 - 8] de  $\mathbf{u}_{\alpha}$  par sa valeur tirée de [1.3 - 11] conduisent aux relations fondamentales existant entre les coordonnées des points-lignes et des points-colonnes sur l'axe  $\alpha$ , les relations quasi-barycentriques :

$$\begin{cases} \psi_{\alpha i} = \frac{1}{\sqrt{\lambda_{\alpha}}} \sum_{j=1}^p \frac{f_{ij}}{f_i} \varphi_{\alpha j} & [1.3 - 12] \\ \varphi_{\alpha j} = \frac{1}{\sqrt{\lambda_{\alpha}}} \sum_{i=1}^n \frac{f_{ij}}{f_j} \psi_{\alpha i} & [1.3 - 13] \end{cases}$$

Ainsi, au coefficient de dilatation  $\frac{1}{\sqrt{\lambda_{\alpha}}}$  près, les projections des points représentatifs d'un nuage sont, sur un axe, les *barycentres* des projections des points représentatifs de l'autre nuage.

La matrice de terme général  $\left( \frac{f_{ij}}{f_i} \right)$  permettant de calculer les coordonnées d'un point  $i$  à partir de tous les points  $j$  (relation [1.3 - 12]) n'est autre que le tableau des profils-lignes.

La coordonnée de la modalité  $i$  d'une des variables est la moyenne des modalités  $j$  de l'autre variable pondérées par les fréquences conditionnelles du profil de  $i$ . De même, la relation [1.3 - 13] montre que la coordonnée de la modalité  $j$  est la moyenne de l'ensemble des modalités  $i$  pondérées par les fréquences conditionnelles du profil de  $j$ .

### Remarques

1) Toutes les valeurs propres sont nécessairement inférieures ou égales à 1. En effet puisque :

$$\sqrt{\lambda_{\alpha}} \psi_{\alpha i} = \sum_{j=1}^p \frac{f_{ij}}{f_i} \varphi_{\alpha j}$$

on a :

$$\min_{(j)} \{\varphi_{\alpha j}\} \leq \sqrt{\lambda_{\alpha}} \psi_{\alpha i} \leq \max_{(j)} \{\varphi_{\alpha j}\}$$

d'où :

$$\max_{(i)} \{\sqrt{\lambda_{\alpha}} \psi_{\alpha i}\} \leq \max_{(j)} \{\varphi_{\alpha j}\}$$

De la même manière, on a :

$$\max_{(j)} \{\sqrt{\lambda_{\alpha}} \varphi_{\alpha j}\} \leq \max_{(i)} \{\psi_{\alpha i}\}$$

comme  $\lambda_{\alpha} \geq 0$  :

$$\max_{(j)} \{\varphi_{\alpha j}\} \leq \max_{(j)} \{\varphi_{\alpha j}\}$$

et finalement:

$$\lambda_{\alpha} \leq 1$$

2) Les relations quasi-barycentriques ne sont pas des cas particuliers des relations de transitions établies lors de l'analyse générale car les matrices "de passage" ne sont pas transposées l'une de l'autre.

## f – Représentation simultanée

Les relations quasi-barycentriques justifient la représentation simultanée des lignes et des colonnes. La figure 1.3 - 12 illustre schématiquement le processus de l'analyse des correspondances.

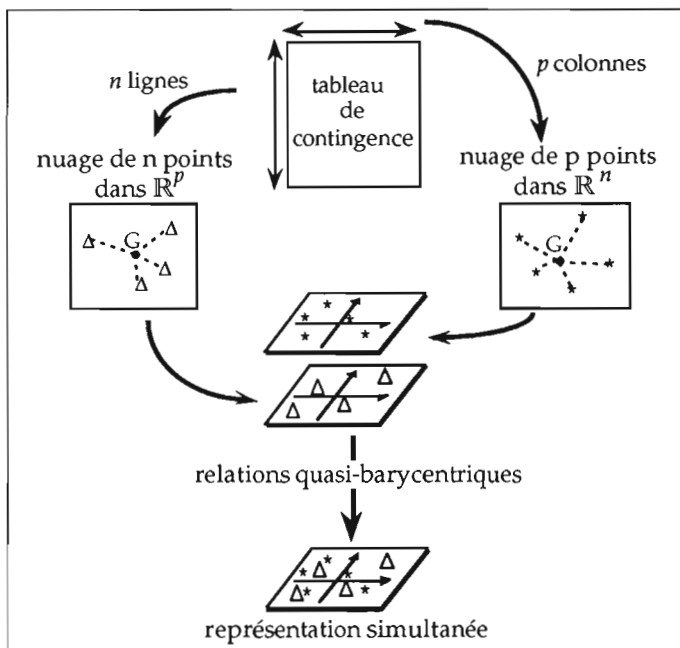


Figure 1.3 - 12  
Schéma de la représentation simultanée

Si les méthodes factorielles sont fondées sur le calcul des distances entre points-lignes et entre points-colonnes, la distance entre un point-ligne et un point-colonne n'a pas de sens puisque ces points sont dans des espaces différents. L'analyse des correspondances offre cependant la possibilité de positionner et d'interpréter un point d'un ensemble relatif à un espace par rapport à l'ensemble des autres points définis dans l'autre espace.

### g – Autre présentation de l'analyse des correspondances

La recherche de la meilleure représentation simultanée des ensembles lignes et colonnes, introduite au paragraphe 1.3.2, est une façon de présenter l'analyse des correspondances qui conduit directement aux formules de calculs analytiques des facteurs. Nous cherchons à représenter sur un même axe l'ensemble des lignes et l'ensemble des colonnes, pour approcher la situation idéale suivante :

- [i] Chaque point-colonne  $j$  est barycentre des points-lignes  $i$ , ceux-ci étant affectés d'une masse  $p_i$  proportionnelle à leur importance dans la modalité  $j$  c'est-à-dire de la masse :  $p_i = \frac{f_{ij}}{f_j}$

Ces masses constituent, pour chaque colonne  $j$ , les profils-colonnes du tableau de données avec  $\sum_{i=1}^n p_i = 1$ .

- [ii] Chaque point-ligne  $i$  est barycentre des points-colonnes  $j$ , chaque point-colonne étant affecté de la masse  $q_j$  représentant la part de la modalité  $j$  dans la modalité  $i$  c'est-à-dire :  $q_j = \frac{f_{ij}}{f_i}$

Ces masses constituent, pour chaque ligne  $i$ , les profils-lignes du tableau de données avec  $\sum_{j=1}^p q_j = 1$ .

Nous définissons ainsi des relations strictement barycentriques entre les deux ensembles. Si  $\varphi_j$  désigne la coordonnée du point-colonne  $j$  sur un axe, et si  $\psi_i$  désigne la coordonnée du point-ligne  $i$  sur ce même axe, les conditions [i] et [ii] s'écrivent respectivement :

$$\begin{cases} \varphi = \mathbf{D}_p^{-1} \mathbf{F}' \psi \\ \psi = \mathbf{D}_n^{-1} \mathbf{F} \varphi \end{cases} \quad \text{soit} \quad \begin{cases} \varphi_j = \sum_{i=1}^n \frac{f_{ij}}{f_j} \psi_i \\ \psi_i = \sum_{j=1}^p \frac{f_{ij}}{f_i} \varphi_j \end{cases}$$

Ces relations sont en général impossibles à réaliser simultanément, car elles impliquent que chaque ensemble soit contenu dans l'autre. (Il existe une solution triviale, pour laquelle tous les points des deux ensembles sont confondus avec le point d'abscisse 1).

Pour approcher cette situation idéale, nous cherchons un coefficient  $\beta$  positif et le plus proche possible de 1, tel que l'on ait les relations :

$$\begin{cases} \varphi = \beta D_p^{-1} F' \psi & [1.3 - 14] \\ \psi = \beta D_n^{-1} F \varphi & [1.3 - 15] \end{cases}$$

Remarquons que  $\beta$  est nécessairement supérieur (ou égal) à 1 sinon les relations [1.3 - 14] et [1.3 - 15] impliqueraient encore que chacun des deux ensembles recouvre un intervalle de l'axe strictement contenu dans l'intervalle recouvert par l'autre. On est donc conduit à chercher le plus petit  $\beta$  positif tel que [1.3 - 14] et [1.3 - 15] soient vérifiées.

Dans [1.3 - 14], par exemple, remplaçons  $\psi$  par sa valeur tirée de [1.3 - 15] :

$$D_p^{-1} F' D_n^{-1} F \varphi = \frac{1}{\beta^2} \varphi$$

Prémultipliant l'équation de l'axe factoriel  $u$  dans  $\mathbb{R}^p$  par  $D_p^{-1}$  :

$$D_p^{-1} F' D_n^{-1} F D_p^{-1} u = \lambda D_p^{-1} u$$

On rappelle que les coordonnées factorielles dans  $\mathbb{R}^n$  valent (cf. formule [1.3 - 11]):

$$\varphi = \sqrt{\lambda} D_p^{-1} u$$

On a donc :

$$D_p^{-1} F' D_n^{-1} F \varphi = \lambda \varphi$$

Et par identification, on obtient :

$$\lambda = \frac{1}{\beta^2} \quad \text{d'où} \quad \beta = \frac{1}{\sqrt{\lambda}}$$

Les relations [1.3 - 14] et [1.3 - 15] ne sont autres que les relations quasi-barycentriques [1.3 - 12] et [1.3 - 13] définies précédemment<sup>1</sup>.

On peut étendre la recherche de la meilleure représentation  $\beta$ -barycentrique sur un axe, à celle de la meilleure représentation  $(\beta_1, \beta_2)$ -barycentrique dans un plan repéré par deux axes orthogonaux, puis généraliser à un sous-espace de dimension quelconque. On trouve alors la représentation simultanée fournie par l'analyse des correspondances<sup>2</sup>.

<sup>1</sup> Puisque le coefficient  $\beta$  doit être supérieur ou égal à 1, on démontre également de cette façon le résultat déjà établi au paragraphe (e) ci-dessus selon lequel, en analyse des correspondances, toutes les valeurs propres sont inférieures ou égales à 1.

<sup>2</sup> Nous verrons également au chapitre 3 d'autres présentations de l'analyse des correspondances (cas particulier des analyses canoniques, discriminantes). D'autres points de vue complémentaires sont développés par Escoufier (1985, 1988).

### h – Formule de reconstitution des données

Les calculs du paragraphe 1.1.5 s'appliquent également au cas de l'analyse des correspondances, en notant toutefois que les vecteurs  $\mathbf{u}_\alpha$  et  $\mathbf{v}_\alpha$  sont maintenant orthonormés pour les métriques  $D_p^{-1}$  et  $D_n^{-1}$ . En partant des relations [1.3 - 6] et [1.3 - 7] (cf. § 1.3.3.d), et en suivant un raisonnement analogue à celui du paragraphe 1.1.5, on obtient la formule :

$$f_{ij} = f_i.f_j \sum_{\alpha=1}^p \sqrt{\lambda_\alpha} \varphi_{\alpha j} \psi_{\alpha i} \quad [1.3 - 16]$$

qui s'écrit aussi, en faisant intervenir la première valeur propre qui vaut 1, et les facteurs correspondants (voir plus bas, paragraphe 1.3 - 7a) :

$$f_{ij} = f_i.f_j \left( 1 + \sum_{\alpha=2}^p \sqrt{\lambda_\alpha} \varphi_{\alpha j} \psi_{\alpha i} \right) \quad [1.3 - 17]$$

### 1.3.4 Règles d'interprétation : inertie, formes de nuages

Les nuages de points-lignes et de points-colonnes vont être représentés dans les plans de projection formés par les premiers axes factoriels pris deux à deux. La lecture des graphiques nécessite cependant des règles d'interprétation, en particulier pour apprécier les proximités, identifier les éléments responsables de la formation des facteurs et ceux qui en sont des caractéristiques. Ces règles s'appuient sur le bilan de l'opération de réduction que constitue la séquence des valeurs propres et des pourcentages d'inertie, ainsi que sur un ensemble de coefficients classiques : les contributions absolues et les cosinus carrés, qui seront étudiés au paragraphe 1.3.5.

La valeur de l'inertie globale n'a pas toujours une interprétation intéressante. En analyse en composantes principales normée (section 1.2) et, nous verrons, en analyse des correspondances multiples (section 1.4), l'inertie totale dépend uniquement du nombre de variables. On interprète, en revanche, les pourcentages d'inertie de chaque axe les uns par rapport aux autres.

Les problèmes de validité et de portée des résultats seront abordés au chapitre 4 dans un cadre général. On se bornera dans cette section à quelques considérations générales.

#### a – Inertie et test d'indépendance

En analyse des correspondances, nous l'avons vu (§ 1.3.2.e), la valeur de l'inertie globale est liée au test classique du  $\chi^2$ .

L'inertie totale  $I$  du nuage de points par rapport au centre de gravité s'écrit par définition :



$$I = \sum_{i=1}^n f_i d^2(i, G) = \sum_{j=1}^p f_j d^2(j, G) = \sum_{j=1}^p \sum_{i=1}^n \left( \frac{f_{ij} - f_i f_j}{f_i f_j} \right)^2$$

L'effectif total étant  $k$ , on reconnaît en  $kI$  la statistique qui est asymptotiquement distribuée suivant la loi du  $\chi^2$  à  $(n-1)(p-1)$  degrés de liberté (sous l'hypothèse d'indépendance) :

$$\chi^2 = kI$$

L'inertie s'exprime également par :

$$I = \sum_{\alpha=1}^{p-1} \lambda_{\alpha}$$

La somme des valeurs propres non triviales d'une analyse des correspondances a donc une interprétation statistique simple. On pourra rejeter l'hypothèse nulle d'indépendance des variables en lignes et en colonnes si la valeur observée  $\chi^2$  dépasse la valeur  $\chi_0^2$  qui a une probabilité d'être dépassée inférieure à un seuil fixé au préalable<sup>1</sup>.

La valeur de l'inertie est un indicateur de la dispersion du nuage et mesure la liaison entre les deux variables.

Cependant, on ne s'intéresse pas seulement à la dispersion du nuage mais surtout à l'existence de directions privilégiées dans ce nuage.

On consulte les inerties de chaque axe (valeurs propres) ainsi que les taux d'inertie correspondants. Cet examen nous renseigne sur la forme du nuage : forme "sphérique" (pas de direction privilégiée) ou forme non sphérique (directions privilégiées).

**Tableau 1.3 - 8**  
**Valeurs propres, pourcentages d'inertie pour la table 1.3 - 1**

N0	VALEUR PROPRE	POUR-CENTAGE	POURCENT. CUMULE	
1	.2088	89.37	89.37	*****
2	.0222	9.51	98.89	***
3	.0026	1.11	100.00	*
Trace	.2336	( = INERTIE TOTALE)		

Le tableau 1.3 - 8 donne les valeurs des trois valeurs propres non nulles de l'analyse de la table 1.3 - 1. L'inertie totale (0.2336), somme des trois valeurs propres, multipliée par l'effectif total de la table (592) donne la valeur 138.29

<sup>1</sup> Cette façon d'opérer un test d'hypothèse correspond à l'usage classique des tables statistiques donnant les valeurs  $\chi_0^2$  pour chaque degré de liberté et pour certains seuils conventionnels (0.05 ou 0.01 en général). Souvent les logiciels donnent directement la probabilité que le  $\chi^2$  calculé soit dépassé. Il suffit alors, sans recours à une table, de comparer cette probabilité aux seuils précédents.

qui doit être une réalisation d'un  $\chi^2$  à 9 degrés de liberté dans l'hypothèse d'indépendance des lignes et des colonnes de la table. Un tel  $\chi^2$  ne dépasse 21.7 que dans 1% des cas (seuil 0.01).

L'hypothèse d'indépendance des couleurs des yeux et des cheveux est donc rejetée. C'est dans une telle circonstance qu'interviendra utilement l'analyse des correspondances, pour décrire cette dépendance entre lignes et colonnes.

D'une façon générale, deux variables sont indépendantes si les profils de leurs modalités sont identiques (aux fluctuations d'échantillonnage près) aux profils moyens (cf. 1.3.3.b) : l'inertie totale est faible et il n'existe pas de direction privilégiée. Géométriquement, cela signifie que tous les points sont concentrés autour du centre de gravité du nuage suivant une forme sphérique. Ceci se traduit par le schéma de la figure 1.3 - 13.

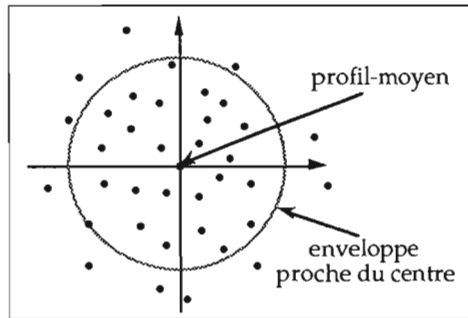


Figure 1.3 - 13  
Situation d'indépendance

Ces indicateurs, portant d'une part sur l'inertie totale et d'autre part sur l'inertie des axes et leurs taux d'inertie, ont donc un intérêt au moment de l'interprétation.

On schématise les principaux cas sur la figure 1.3 - 14. On remarque que, dans les situations 2 et 4, les nuages ont des taux d'inertie identiques mais une inertie totale différente. Par ailleurs, les situations 3 et 4 révèlent deux nuages de même inertie totale et des taux d'inertie différents.

Le test du  $\chi^2$  permet de détecter ces deux dernières situations, mais ne permet pas de mettre en évidence la situation 2 (cf. section 4.1 pour l'étude de cette question).

Enfin, l'inertie d'un facteur mesure la liaison qu'il met en évidence. Elle ne peut être supérieure à 1 (cf. § 1.3.3.f). Une valeur propre qui tend vers 1 indique une dichotomie au niveau des données ; on obtient pour chaque variable deux groupes de modalités séparant le nuage de points en deux sous-nuages. Cela peut signifier également l'existence d'un groupe de points isolés des autres points (constituant alors l'autre groupe).

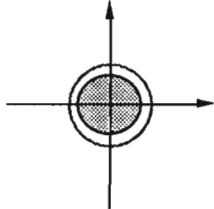
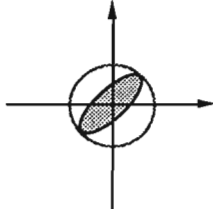
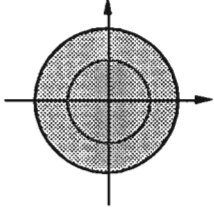
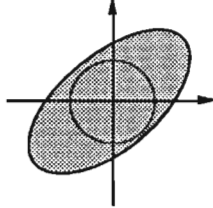
Nuage		Directions Taux d'inerties des axes	
		Forme "sphérique"	Forme "non-sphérique"
Inertie	Faible inertie	 <p>1- INDÉPENDANCE</p> <ul style="list-style-type: none"> <li>• faible inertie totale</li> <li>• pas de direction privilégiée</li> </ul>	 <p>2- DÉPENDANCE</p> <ul style="list-style-type: none"> <li>• faible inertie totale</li> <li>• direction privilégiée</li> </ul>
	Forte inertie	 <p>3- DÉPENDANCE</p> <ul style="list-style-type: none"> <li>• forte inertie totale</li> <li>• pas de direction privilégiée</li> </ul>	 <p>4- DÉPENDANCE</p> <ul style="list-style-type: none"> <li>• forte inertie totale</li> <li>• direction privilégiée</li> </ul>

Figure 1.3 - 14  
Indépendance et dépendances

Lorsque deux valeurs propres sont proches de 1, on obtient trois sous-nuages et les modalités des variables se décomposent en trois groupes. Si toutes les valeurs propres sont proches de 1, chaque modalité d'une variable est en correspondance presque exclusive avec une seule modalité de l'autre variable.

Cependant des valeurs propres faibles (signifiant que les profils sont proches du profil moyen) ne doivent pas empêcher une interprétation des axes d'inertie associés. Ceux-ci peuvent révéler une structure intéressante et plus difficilement perceptible. Ce point sera repris au chapitre 4, § 4.1.3.

**b – Quelques formes caractéristiques de nuages de points**

Envisageons quelques formes classiques de nuages afin de montrer comment la configuration du nuage de points projeté permet de réorganiser le tableau de données, par permutation des lignes et des colonnes et ainsi de mieux l'interpréter.

- *Le nuage de points est scindé en deux sous-nuages*

Le tableau de données peut être réorganisé en ordonnant les coordonnées des lignes et des colonnes sur le premier facteur. On obtient de façon schématique :

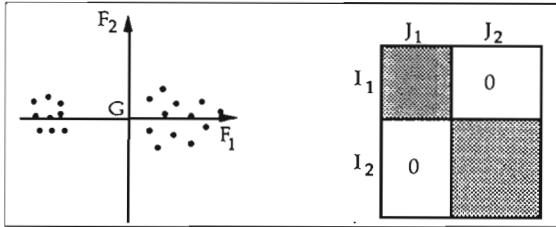


Figure 1.3 - 15

Nuage de points scindé en deux sous-nuages

Il peut être intéressant d'analyser séparément les deux sous-nuages définis par les deux tableaux de correspondances  $(I_1, J_1)$  et  $(I_2, J_2)$ .

- *Le nuage se décompose en trois sous-nuages de points*

On réorganise de la même manière le tableau de données par permutation des lignes et des colonnes. Les trois sous-nuages peuvent également faire l'objet d'analyses séparées.

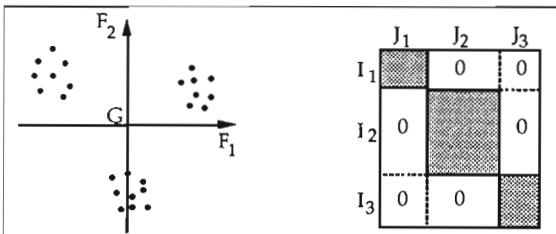


Figure 1.3 - 16

Nuage de points scindé en trois sous-nuages

- *"L'effet Guttman"*

On peut aboutir ainsi à la situation où le nuage de points a une forme parabolique. Le tableau correspondant est réordonné suivant une diagonale chargée :

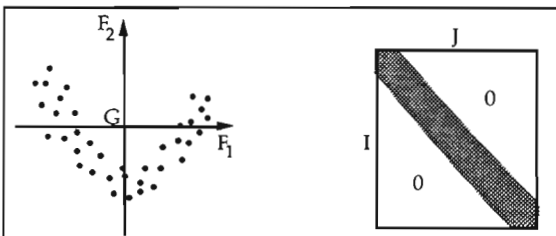


Figure 1.3 - 17

Effet Guttman et structure possible du tableau

Cette situation met en évidence "l'effet Guttman" qui traduit une redondance des deux variables : de la connaissance de la ligne  $i$ , on en déduit la colonne  $j$ . Toute l'information est quasiment donnée par le premier facteur.

Pourtant le tableau n'est pas de rang 1 et l'on disposera de  $p-1$  facteurs. Mais le deuxième facteur est une fonction du second degré du premier facteur, le troisième est une fonction du troisième degré, etc. L'information donnée par les axes de rang ultérieurs traduit le même phénomène. Cependant l'examen du deuxième facteur affine l'interprétation du premier axe<sup>1</sup>.

Généralement l'effet Guttman apparaît lorsque les variables sont ordonnées (variables continues transformées en variables nominales). Un axe (souvent le premier) oppose les valeurs extrêmes et un autre axe oppose les valeurs intermédiaires aux valeurs extrêmes. L'effet Guttman met parfois en évidence une structure triviale qui pourra cependant être intéressante si la forme parabolique n'est pas parfaite. Les points de rupture sont alors intéressants à analyser.

### 1.3.5 Règles d'interprétation : contributions et cosinus

Deux séries de coefficients apportent une information supplémentaire par rapport aux coordonnées factorielles :

- les *contributions*, parfois appelées contributions absolues, qui expriment la part prise par une modalité de la variable dans l'inertie (ou variance) "expliquée" par un facteur;
- les *cosinus carrés*, parfois appelés contributions relatives ou qualité de représentation, qui expriment la part prise par un facteur dans la dispersion d'une modalité de la variable.

C'est après l'examen de ces coefficients que l'on pourra interpréter les graphiques factoriels en tenant compte des relations de transition.

#### a – Contributions

On cherche à connaître les éléments responsables de la construction de l'axe  $\alpha$ . Calculons la variance des coordonnées des  $n$  points-lignes  $i$  sur l'axe  $\alpha$ , chacun d'eux étant muni de la masse  $f_i$ .

L'origine étant prise au centre de gravité, les coordonnées factorielles sont centrées (cf. formule [1.3 - 4]) et la variance vaut  $\lambda_\alpha$  (cf. formule [1.3 - 5]).

Ainsi le quotient :

---

<sup>1</sup> Sur l'effet Guttman en analyse des correspondances, cf. Benzécri (1973, chapitre II.B-7 et II.B-10), Heiser (1986), Van Rijckevorsel (1987) ; Tenenhaus (1994, chapitre 7, §9).

$$Cr_{\alpha}(i) = \frac{f_i \cdot \psi_{\alpha i}^2}{\lambda_{\alpha}}$$

mesure la part de l'élément  $i$  dans la variance prise en compte sur l'axe  $\alpha$ .

Ce quotient est appelé *contribution* de l'élément  $i$  à l'axe  $\alpha$  et permet de savoir dans quelle proportion un point  $i$  contribue à l'inertie  $\lambda_{\alpha}$  du nuage projeté sur l'axe  $\alpha$ .

On notera que pour tout axe  $\alpha$  :

$$\sum_{i=1}^n Cr_{\alpha}(i) = 1$$

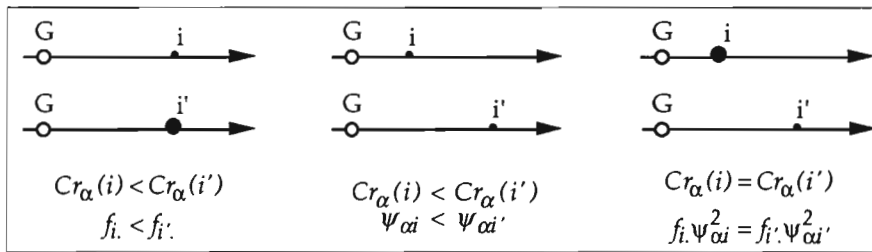


Figure 1.3 - 18  
Contributions à l'axe  $\alpha$  : trois cas de figure.

De la même façon on définit la contribution de l'élément  $j$  à l'axe  $\alpha$  par :

$$Cr_{\alpha}(j) = \frac{f_j \cdot \varphi_{\alpha j}^2}{\lambda_{\alpha}}$$

avec la relation :

$$\sum_{j=1}^p Cr_{\alpha}(j) = 1$$

Pour trouver une éventuelle signification à un axe, on s'intéresse d'abord aux points ayant une forte contribution. Ce sont eux qui fixent la position de l'axe (dans  $\mathbb{R}^p$  pour les points  $i$ , et dans  $\mathbb{R}^n$  pour les points  $j$ ).

### b – Cosinus carrés

On cherche à apprécier si un point est bien représenté sur un sous-espace factoriel.

Les axes factoriels de chaque espace constituent des bases orthonormées. Le carré de la distance d'un point au centre de gravité se décompose en somme de carrés des coordonnées sur ces axes.

Pour un point  $i$  de  $\mathbb{R}^p$ , on a :

$$d^2(i, G) = \sum_{j=1}^p \frac{1}{f_{.j}} \left( \frac{f_{ij}}{f_{.i}} - f_{.j} \right)^2$$

On remarque que la distance s'annule lorsque le profil du point est égal au profil moyen.

Le carré de la projection de la variable  $i$  sur l'axe  $\alpha$  vaut :

$$d_{\alpha}^2(i, G) = \psi_{\alpha i}^2$$

Notons que :

$$\sum_{\alpha} d_{\alpha}^2(i, G) = d^2(i, G)$$

Un point  $i$  dans  $\mathbb{R}^p$  est plus ou moins proche de l'axe  $\alpha$ . La proximité entre deux points projetés sur l'axe  $\alpha$  correspond d'autant mieux à leur distance réelle que les points sont plus proches de l'axe.

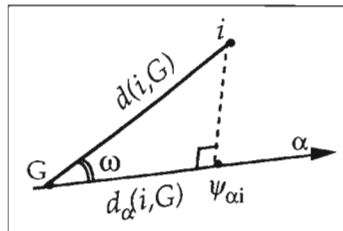


Figure 1.3 - 19  
Projection du point  $i$  sur l'axe  $\alpha$

La "qualité" de la représentation du point  $i$  sur l'axe  $\alpha$  peut être évaluée par le cosinus de l'angle entre l'axe et le vecteur joignant le centre de gravité du nuage au point  $i$  :

$$\text{Cos}_{\alpha}^2(i) = \frac{d_{\alpha}^2(i, G)}{d^2(i, G)} = \frac{\psi_{\alpha i}^2}{d^2(i, G)}$$

Cette quantité, appelée *cosinus carré*, représente la part de la distance au centre prise en compte dans la direction  $\alpha$ . On l'appelle aussi la *contribution relative* du facteur à la position du point  $i$ .

Plus le cosinus carré est proche de 1, plus la position du point observé en projection est proche de la position réelle du point dans l'espace (figure 1.3 - 20).

On apprécie la qualité de la représentation d'un point dans un plan en faisant la somme des cosinus carrés sur les axes étudiés.

Notons que pour tout  $i$  :

$$\sum_{\alpha} \text{Cos}_{\alpha}^2(i) = 1$$

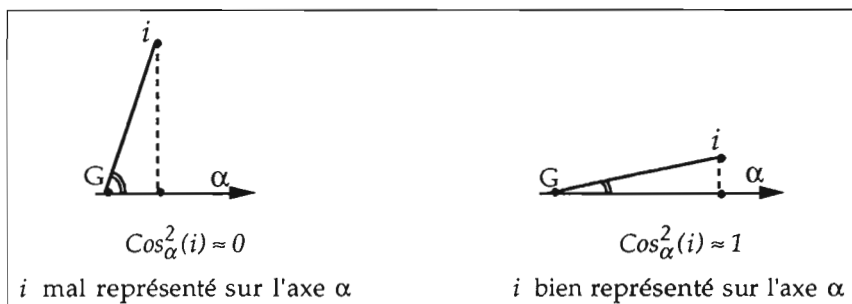


Figure 1.3 - 20  
Qualité de représentation d'un point  $i$  sur l'axe  $\alpha$

Ce qui vient d'être dit des  $n$  points-lignes peut être transposé aux  $p$  éléments de l'autre ensemble. On mesure la contribution relative du facteur  $\alpha$  à la position du point  $j$  par le cosinus carré de  $j$  :

$$\text{Cos}_\alpha^2(j) = \frac{\varphi_{\alpha j}^2}{d^2(j, G)}$$

et l'on a également pour tout  $j$  :

$$\sum_{\alpha} \text{Cos}_\alpha^2(j) = 1$$

Pour analyser les proximités entre points, on s'intéresse surtout aux points ayant un cosinus carré élevé. Les proximités entre ces points, observés dans le sous-espace factoriel, donnent une bonne image de leurs proximités réelles.

#### Remarque

Pour les contributions ainsi que pour les cosinus carrés, il n'y a pas de valeurs "seuils" à partir desquelles on peut dire que telle ou telle valeur est "forte" ou "faible". Les appréciations se font empiriquement, en fonction de l'ensemble des valeurs calculées et varient d'un jeu de données à un autre<sup>1</sup>.

#### c – Exemple numérique

L'exemple concerne toujours l'analyse des correspondances de la table 1.3 - 1. Les coordonnées sur le premier axe (tableau 1.3 - 9) montrent que la couleur des cheveux "blond" s'oppose à toutes les autres sur le premier axe, mais surtout à "brun". Le point "blond" a une contribution de 71.7% au premier axe et un cosinus carré de 0.99 : il est pratiquement sur cet axe et ne pourra donc pas caractériser les axes ultérieurs. Notons que le point "roux" a une contribution très faible sur le premier axe (1.0%).

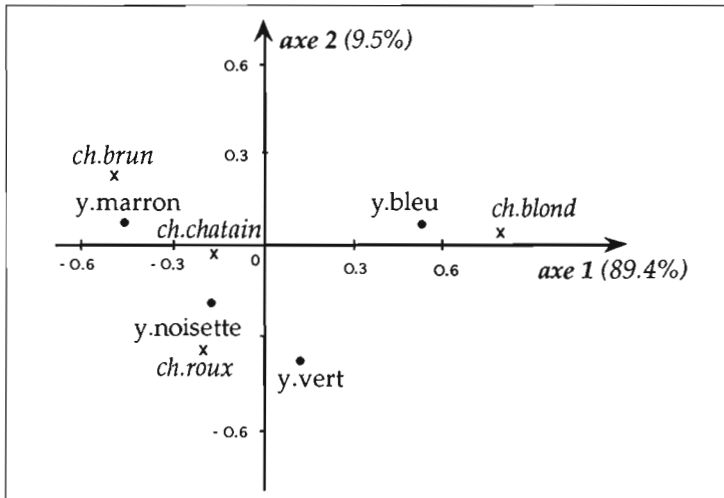
<sup>1</sup> Notons qu'il est usuel de multiplier par 100 les contributions, de façon à exprimer en pourcentage la participation de chaque point.



**Tableau 1.3 - 9**  
**Coordonnées, contributions, cosinus carrés**  
**pour l'analyse des correspondances de la table 1.3 - 1**

COLONNES	COORDONNEES			CONTRIBUTIONS			COSINUS CARRÉS		
	1	2	3	1	2	3	1	2	3
<b>CHEVEUX</b>									
Ch.Brun	-.50	.21	-.06	22.2	37.9	21.6	.84	.15	.01
Ch.chatain	-.15	-.03	.05	5.1	2.3	44.3	.86	.04	.09
Ch.roux	-.13	-.32	-.08	1.0	55.1	31.9	.13	.81	.05
Ch.blond	.84	.07	-.02	71.7	4.7	2.2	.99	.01	.00
<b>LIGNES</b>									
	1	2	3	1	2	3	1	2	3
<b>YEUX</b>									
y.marron	-.49	.09	-.02	43.1	13.0	6.7	.97	.03	.00
y.noisette	-.21	-.17	.10	3.4	19.8	61.1	.54	.34	.12
y.vert	.16	-.34	-.09	1.4	55.9	31.9	.18	.77	.05
y.bleu	.55	.08	.00	52.1	11.2	.3	.98	.02	.00

Le second axe (dont on a vu qu'il correspondait à une valeur propre près de dix fois plus petite que le premier) est essentiellement construit par la couleur "roux" (55.1 %) qui s'oppose simultanément à "brun" et "blond". La couleur "roux" est le seul point bien représenté sur l'axe 2 (cosinus carré de 0.81). Pour les points-lignes, le premier axe est construit presque exclusivement par les yeux "marrons" et "bleus" (contributions de 43.1% et 52.1%), points situés pratiquement sur l'axe (cosinus carrés de 0.97 et 0.98), le second axe étant surtout lié aux yeux "verts".



**Figure 1.3 - 21**  
**Premier plan factoriel pour l'analyse de la table 1.3 - 1**

On note que la consultation des coordonnées pouvait faire penser que les yeux "noisettes" et "verts" jouaient un certain rôle dans la construction du

premier axe. La figure 1.3 - 21 qui utilise les deux premières coordonnées, montre le caractère suggestif de la représentation graphique simultanée des lignes et des colonnes. Elle permet d'interpréter les proximités ou distances entre points d'un même ensemble par leur association avec ceux de l'autre ensemble.

Pourquoi par exemple le point "ch.blond" est-il plus excentré que le point "y.bleu" sur ce premier axe très dominant ? Parce que les cheveux blonds sont beaucoup mieux caractérisés par les yeux bleus que l'inverse : d'après le tableau 1.3 - 3 (profils colonnes), 74% des blonds ont les yeux bleus, alors que d'après le tableau 1.3 - 2 (profils lignes) 44% des personnes ayant les yeux bleus ont des cheveux blonds.

En d'autres termes, dans la relation quasi-barycentrique qui permet de positionner le point "ch.blond", le point "y.bleu" a un poids relatif de 0.74, alors que dans la relation quasi-barycentrique qui permet de positionner le point "y.bleu", le point "blond" n'a qu'un poids relatif de 0.44.

### 1.3.6 Éléments supplémentaires

On dispose par exemple de  $p_s$  colonnes supplémentaires qui concernent des modalités de variables nominales, analogues aux colonnes de la table de contingence.

Il s'agit de situer ces nouveaux points-colonnes par rapport aux  $p$  points analysés. Soit  $k_{ij}^+$  la  $i^{\text{ème}}$  coordonnée de la  $j^{\text{ème}}$  colonne supplémentaire. Son profil est donné par :

$$\left\{ \frac{k_{ij}^+}{k_j^+}; i = 1, 2, \dots, n \right\} \quad \text{avec} \quad k_j^+ = \sum_{i=1}^n k_{ij}^+$$

On projette ce point  $j$  sur l'axe  $\alpha$  en utilisant la même formule de transition [1.3 - 13] que pour les colonnes du tableau de fréquences :

$$\varphi_{\alpha j}^+ = \frac{1}{\sqrt{\lambda_\alpha}} \sum_{i=1}^n \frac{k_{ij}^+}{k_j^+} \psi_{\alpha i}$$

Pour une modalité  $i$  d'une variable portée en ligne supplémentaire, on aura de façon analogue (formule de transition 1.3 - 12) :

$$\psi_{\alpha i}^+ = \frac{1}{\sqrt{\lambda_\alpha}} \sum_{j=1}^p \frac{k_{ij}^+}{k_i^+} \varphi_{\alpha j}$$

A l'instar des éléments analysés, les modalités supplémentaires se calculent et s'interprètent comme des quasi-barycentres.

**Remarques**

1) Les éléments en supplémentaires, n'intervenant pas dans la construction du nuage, sont affectés d'un poids nul et leur contribution est donc nulle. En revanche, les cosinus carrés restent des aides à l'interprétation de ces éléments <sup>1</sup>.

2) La somme des cosinus carrés d'un élément supplémentaire sur l'ensemble des facteurs peut être inférieure à 1 alors que pour les éléments actifs elle est exactement égale à 1.

En effet, supposons  $n > p$  et plaçons-nous dans l'espace des lignes. Un point-colonne actif  $j$  est défini dans  $\mathbb{R}^n$  mais il est situé, par l'analyse, dans l'espace factoriel à  $p - 1$  dimensions. Il suffit de  $p - 1$  coordonnées pour positionner cet élément. Un élément-colonne supplémentaire  $j^*$  sera positionné dans l'espace à  $p - 1$  dimensions construit par l'analyse alors qu'il appartient à  $\mathbb{R}^p$ . Les éléments supplémentaires ne sont donc pas entièrement contenus dans l'espace factoriel<sup>2</sup>.

**1.3.7 Mise en œuvre des calculs**

La distance du  $\chi^2$  ne diffère en fait de la métrique euclidienne usuelle que par l'introduction d'une pondération. On peut se ramener à la métrique euclidienne usuelle par un changement de coordonnées initial. Les calculs en sont simplifiés et, notamment, la matrice à diagonaliser devient symétrique. Par ailleurs, l'analyse par rapport aux centres de gravité est équivalente à l'analyse par rapport à l'origine.

**a – Analyse par rapport à l'origine ou au centre de gravité du nuage**

Nous raisonnerons, pour fixer les idées, dans  $\mathbb{R}^p$ .

Le centre de gravité  $G$  du nuage des profils-lignes a pour  $j^{\text{ième}}$  composante :

$$g_j = \sum_{i=1}^n f_i \cdot \frac{f_{ij}}{f_i} = f_j$$

L'analyse par rapport au centre de gravité revient à remplacer  $\frac{f_{ij}}{f_i}$  par  $\frac{f_{ij}}{f_i} - f_j$

c'est-à-dire par  $\frac{f_{ij} - f_i f_j}{f_i}$ .

Remarquons que le nuage est contenu dans un hyperplan  $\mathcal{H}$  à  $p - 1$  dimensions défini pour tout  $i$  par la relation :

$$\sum_{j=1}^p \frac{f_{ij}}{f_i} = 1$$

<sup>1</sup> Pour une vue d'ensemble sur le rôle et l'utilisation des variables supplémentaires en analyse des correspondances, cf. Cazes (1982).

<sup>2</sup> Cette remarque vaut également pour l'analyse en composantes principales.

Ce sous-espace contient le centre de gravité  $G$  et les axes factoriels de l'analyse par rapport à  $G$ . La somme des composantes de ces facteurs est nulle.

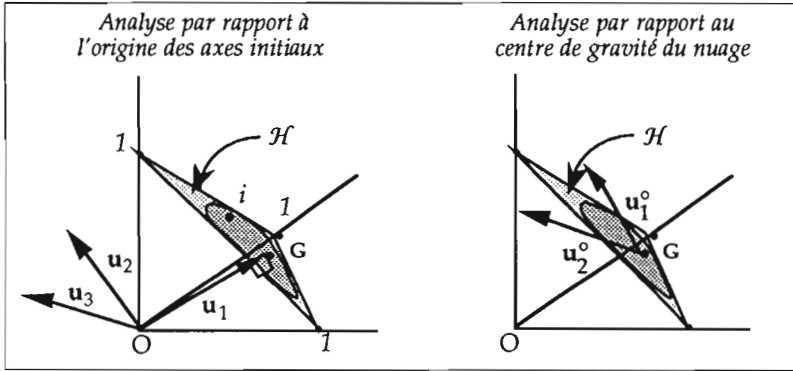


Figure 1.3 - 22  
Analyse dans  $\mathbb{R}^3$

Dans l'analyse par rapport à l'origine, la première direction  $u_1$  est l'axe joignant l'origine au centre de gravité du nuage orthogonalement à  $\mathcal{H}$ . L'inertie projetée sur cet axe vaut 1, égale à la distance entre l'origine et le centre de gravité, puisque la projection des points du nuage sur cet axe est confondue avec le centre de gravité. Les  $p - 1$  axes suivants ( $u_2, \dots, u_\alpha, \dots, u_p$ ) contenus dans  $\mathcal{H}$  constituent une base définissant des directions de droites d'inertie maximum du nuage. Ils coïncident avec les  $p - 1$  premiers axes de l'analyse par rapport au centre de gravité ( $u_1^0, \dots, u_\alpha^0, \dots, u_{p-1}^0$ ).

Le  $p$ <sup>ième</sup> axe correspond à  $u_1$  et n'indique aucune direction dans  $\mathcal{H}$  puisqu'il n'est pas contenu dans  $\mathcal{H}$ . Son inertie (valeur propre) associée, est nulle.

$S$  étant la matrice à diagonaliser du nuage non centré et  $S^0$  celle du nuage centré, on a les relations :

$$s_{jj}^0 = s_{jj} - f_j$$

et pour  $1 < \alpha < p - 1$  :

$$\begin{aligned} u_\alpha^0 &= u_{\alpha+1} & \text{et} & & \lambda_\alpha^0 &= \lambda_{\alpha+1} \\ u_p^0 &= u_1 & \text{et} & & \lambda_p^0 &= 0 \quad \text{et} \quad \lambda_1 = 1 \end{aligned}$$

Ainsi dans  $\mathbb{R}^p$  (et il en est de même dans  $\mathbb{R}^n$ ), il est équivalent de réaliser l'analyse des correspondances sur le tableau de données centrées de terme général :

$$\frac{f_{ij}}{f_i} - f_j$$

ou sur le tableau de données non centrées de terme général :

$$\frac{f_{ij}}{f_i}$$

On peut donc diagonaliser la matrice  $\mathbf{S}$  de l'analyse par rapport à l'origine<sup>1</sup>, en prenant soin d'éliminer le premier vecteur propre reliant l'origine au centre de gravité du nuage et la valeur propre associée égale à 1.

### b – Symétrisation de la matrice à diagonaliser

La matrice à diagonaliser  $\mathbf{S} = \mathbf{F}' \mathbf{D}_n^{-1} \mathbf{F} \mathbf{D}_p^{-1}$ , dans  $\mathbb{R}^p$ , n'est pas en général symétrique. Son terme général s'écrit :

$$s_{jj'} = \sum_{i=1}^n \frac{f_{ij} f_{ij'}}{f_i f_{i,j'}}$$

Considérons la matrice  $\hat{\mathbf{A}} = \mathbf{F}' \mathbf{D}_n^{-1} \mathbf{F}$  symétrique et la matrice  $\mathbf{D}_p^{-1}$  diagonale. On exprime alors  $\mathbf{S}$  de la manière suivante :

$$\mathbf{S} = \hat{\mathbf{A}} \mathbf{D}_p^{-1/2} \mathbf{D}_p^{-1/2}$$

Partant de la relation  $\mathbf{S} \mathbf{u} = \lambda \mathbf{u}$ , il vient :

$$\hat{\mathbf{A}} \mathbf{D}_p^{-1/2} \mathbf{D}_p^{-1/2} \mathbf{u} = \lambda \mathbf{u}$$

Prémultiplions les deux membres par  $\mathbf{D}_p^{-1/2}$  et en posant  $\mathbf{D}_p^{-1/2} \mathbf{u} = \mathbf{w}$ , on obtient :

$$\mathbf{D}_p^{-1/2} \hat{\mathbf{A}} \mathbf{D}_p^{-1/2} \mathbf{w} = \lambda \mathbf{w}$$

La matrice  $\mathbf{A}$  est symétrique :

$$\mathbf{A} = \mathbf{D}_p^{-1/2} \hat{\mathbf{A}} \mathbf{D}_p^{-1/2} = \mathbf{D}_p^{-1/2} \mathbf{F}' \mathbf{D}_n^{-1} \mathbf{F} \mathbf{D}_p^{-1/2}$$

et :

$$\mathbf{A} \mathbf{w} = \lambda \mathbf{w}$$

Les matrices  $\mathbf{S}$  et  $\mathbf{A}$  ont mêmes valeurs propres  $\lambda$ . Leurs vecteurs propres sont liés par la relation :

$$\mathbf{u} = \mathbf{D}_p^{-1/2} \mathbf{w}$$

Il est plus facile de diagonaliser la matrice  $\mathbf{A}$  de terme général :

$$a_{jj'} = \sum_{i=1}^n \frac{f_{ij} f_{ij'}}{f_i \sqrt{f_{i,j} f_{i,j'}}$$

<sup>1</sup> Compte tenu du critère d'ajustement, on considère l'inertie totale du nuage centré, égale à la trace  $tr(\mathbf{S}^\circ)$  de  $\mathbf{S}^\circ$  et l'on a :  $tr(\mathbf{S}^\circ) = tr(\mathbf{S}) - 1$ .

**Remarque :**

C'est la matrice à diagonaliser si l'on choisit de prendre comme coordonnées initiales du point  $i$ , les  $p$  quantités :

$$x_{ij} = \frac{f_{ij}}{f_i \cdot \sqrt{f_{.j}}} \quad (j = 1, \dots, p)$$

Dans ce cas, la distance du  $\chi^2$  entre deux points  $i$  et  $i'$  devient, avec les nouvelles coordonnées, la distance euclidienne usuelle :

$$d^2(i, i') = \sum_{j=1}^p \left( \frac{f_{ij}}{f_i \cdot \sqrt{f_{.j}}} - \frac{f_{i'j}}{f_{i'} \cdot \sqrt{f_{.j}}} \right)^2$$

Cette transformation du tableau des fréquences relatives conduit à la diagonalisation d'une matrice symétrique.

Notons que les coordonnées du centre de gravité  $G$  sont alors :

$$G_j = \sqrt{f_{.j}}$$

et les coordonnées du point  $i$  après recentrage :

$$\frac{f_{ij}}{f_i \cdot \sqrt{f_{.j}}} - \sqrt{f_{.j}} = \frac{f_{ij} - f_i \cdot f_{.j}}{f_i \cdot \sqrt{f_{.j}}}$$

**1.3.8 Exemple d'application**

L'exemple concerne l'analyse d'un tableau de contingence qui croise 8 professions et catégories socioprofessionnelles (PCS) et 6 types de médias pour un échantillon de 12 388 "contacts média" relatifs à 4433 personnes interrogées. L'individu statistique sera pour nous le "contact média" et non la personne interrogée dans l'enquête. Comme ce fut le cas pour l'exemple traité au paragraphe 1.2.11, les données sont extraites de l'*Enquête Budget-temps Multimédia 1991-1992* du CESP.

Afin d'interpréter plus efficacement les représentations obtenues, on projettera en éléments supplémentaires certaines autres caractéristiques de la population enquêtée telles que le sexe, l'âge, le niveau d'instruction.

Nous disposons des tables de contingence suivantes (cf. tableau 1.3 - 10). Pour le premier blocs  $K$  de 8 lignes (lignes actives) on trouve, à l'intersection de la ligne  $i$  et de la colonne  $j$  le nombre  $k_{ij}$  d'individus appartenant à la catégorie  $i$  et ayant eu la veille (un jour de semaine) au moins un contact avec le type de média  $j$ . Les blocs suivants (lignes supplémentaires) s'interprètent de façon analogue. Une personne interrogée pouvant avoir des contacts avec plusieurs médias, les sommes en ligne représentent des "nombres de contacts"<sup>1</sup>.

<sup>1</sup> Il y a 12 388 contacts pour 4433 individus concernés. Les chiffres publiés ici ayant été arrondis après un redressement, les totaux relatifs aux différentes partitions de la population peuvent ne pas coïncider.

**Tableau 1.3 - 10**  
**Tables de contingence croisant les types de contacts-média (colonnes)**  
**avec professions, sexe, âge, niveau d'éducation (lignes).**

	Radio	Tél.	Quot.N.	Quot	R. P.	Mag. P. TV
<b>Professions</b>						
Agriculteur	96	118	2	71	50	17
Petit patron	122	136	11	76	49	41
Prof. Cad. S.	193	184	74	63	103	79
Prof. interm	360	365	63	145	141	184
Employé	511	593	57	217	172	306
Ouvrier qual	385	457	42	174	104	220
Ouvrier n-q	156	185	8	69	42	85
Inactif	1474	1931	181	852	642	782
<b>Sexe</b>						
Homme	1630	1900	285	854	621	776
Femme	1667	2069	152	815	683	938
<b>Age</b>						
15-24 ans	660	713	69	216	234	360
25-34 ans	640	719	84	230	212	380
35-49 ans	888	1000	130	429	345	466
50-64 ans	617	774	84	391	262	263
65 ans ou +	491	761	70	402	251	245
<b>Education</b>						
Primaire	908	1307	73	642	360	435
Secondaire	869	1008	107	408	336	494
Techn. prof.	901	1035	80	140	311	504
Supérieur	619	612	177	209	298	281

On cherche à décrire les éventuelles affinités entre les groupes socioprofessionnels et les différents types de médias.

L'analyse des correspondances de la table **K** conduit aux valeurs propres consignées dans le tableau 1.3 - 11.

**Tableau 1.3 - 11**  
**Valeurs propres, pourcentages d'inertie pour la table K**  
**"Professions-Contacts média" (8 premières lignes de la table 1.3 - 11)**

NUM.	VALEUR PROPRE	POURCENT.	POURCENT. CUMULE	
1	.0139	62.20	62.20	*****
2	.0072	32.37	94.56	*****
3	.0008	3.70	98.26	**
4	.0003	1.36	99.63	*
5	.0001	.37	100.00	*
SOMME	.0223			

Le produit de la trace  $t = 0.0223$  par l'effectif total  $k = 12\ 388$  vaut :

$$kt = 276.25$$

Dans l'hypothèse d'indépendance des lignes et des colonnes de la table, cette quantité serait une réalisation d'un  $\chi^2$  à 35 degrés de liberté (noté  $\chi_{35}^2$ ) [ $35 = (8-1)(6-1)$ ].

Lorsque le nombre de degrés de liberté  $n$  dépasse 30, on considère que la variable  $u = \frac{\chi_n^2 - n}{\sqrt{2n}}$  est une variable normale (de Laplace-Gauss) centrée réduite. Ici,  $u = 28.8$  (28.8 écarts-types de la moyenne). L'hypothèse d'indépendance est évidemment rejetée.

Deux facteurs sont dominants et représentent près de 95% de l'inertie totale. Les coordonnées et les aides à l'interprétation correspondants figurent dans le tableau 1.3 - 12. Celui-ci donne également les coordonnées et les cosinus carrés des lignes supplémentaires.

**Tableau 1.3 - 12**  
Poids relatifs (P.REL), Distances à l'origine (DIS), coordonnées, contributions et cosinus carrés des éléments sur les trois premiers axes

LIBELLES	FREQUENCES		COORDONNEES			CONTRIBUTIONS			COSINUS CARRES		
	P.REL	DIS	1	2	3	1	2	3	1	2	3
<b>COLONNES ACTIVES</b>											
Radio	26.61	.00	-.01	.02	-.05	.4	1.8	70.4	.08	.17	.75
Télévision	32.04	.00	.05	.00	.02	6.6	.0	10.5	.85	.00	.08
Quotidien natio	3.54	.29	-.54	-.01	.02	74.6	.0	1.8	.99	.00	.00
Quotidien regio	13.46	.02	.11	-.11	.01	11.5	22.4	.4	.49	.49	.00
Presse Magazine	10.52	.03	-.09	-.13	.02	6.8	25.6	4.5	.32	.62	.01
Presse Mag. T.V.	13.84	.03	.01	.16	.03	.1	50.1	12.4	.00	.96	.03
<b>LIGNES ACTIVES</b>											
Agriculteur	2.86	.13	.17	-.31	-.07	5.7	38.0	17.9	.21	.74	.04
Petit patron	3.51	.03	.07	-.14	-.06	1.2	10.0	17.7	.15	.67	.14
Prof. Cadre Sup	5.62	.19	-.43	-.06	.00	75.0	2.9	.1	.98	.02	.00
Prof. interm	10.15	.01	-.11	.03	-.03	8.3	1.5	11.8	.80	.08	.07
Employé	14.98	.01	.02	.10	-.01	.3	18.9	.5	.03	.93	.00
Ouvrier qual	11.16	.01	.04	.10	-.02	1.5	15.9	5.1	.14	.74	.03
Ouvrier n-q	4.40	.02	.12	.09	-.04	4.4	5.5	8.4	.56	.36	.06
Inactif	47.32	.00	.03	-.03	.03	3.6	7.3	38.7	.37	.39	.24
<b>LIGNES ILLUSTRATIVES (SUPPLEMENTAIRES)</b>											
Homme	48.97	.01	-.05	-.02	-.01	.0	.0	.0	.48	.11	.02
Femme	51.05	.00	.05	.02	.01	.0	.0	.0	.49	.10	.02
15-24 ans	18.18	.02	-.02	.10	-.04	.0	.0	.0	.02	.56	.08
25-34 ans	18.28	.02	-.03	.12	-.01	.0	.0	.0	.05	.87	.01
35-49 ans	26.30	.00	-.03	.01	-.01	.0	.0	.0	.61	.10	.07
50-64 ans	19.30	.01	.02	.10	.00	.0	.0	.0	.05	.80	.00
65 ans ou +	17.92	.03	.07	-.14	.07	.0	.0	.0	.14	.58	.16
Primaire	30.07	.03	.13	-.08	.02	.0	.0	.0	.63	.24	.02
Secondaire	26.01	.00	.00	.04	.00	.0	.0	.0	.00	.69	.00
Techn. prof.	23.98	.07	-.03	.18	-.04	.0	.0	.0	.01	.46	.02
Supérieur	17.73	.09	-.29	-.02	-.01	.0	.0	.0	.99	.00	.00

On note que l'élément "Quotidien national" dont la fréquence relative (colonne P.REL) est très faible (3.54%) a une distance au point moyen (colonne DIS) très élevée : le profil correspondant est donc atypique. Il contribue pour 74.6% à la construction du premier axe, qui en est très proche (cosinus carré : 0.99). Ce même premier axe est caractérisé par la ligne active



"Prof.Cadre" (profession libérale, cadres supérieurs) et par la ligne supplémentaire "Supérieur" (niveau d'étude supérieur).

Le second axe sépare la "Presse Magazine de Télévision" (associée aux catégories employés et ouvriers, et aux classes d'âges plutôt jeunes) de la presse magazine (Presse TV exclue) et de la presse quotidienne régionale, toutes deux associées aux agriculteurs et aux petits patrons, et à des catégories d'âge plus élevées.

Les figures 1.3 - 23 et 1.3 - 24 résument ce réseau d'associations.

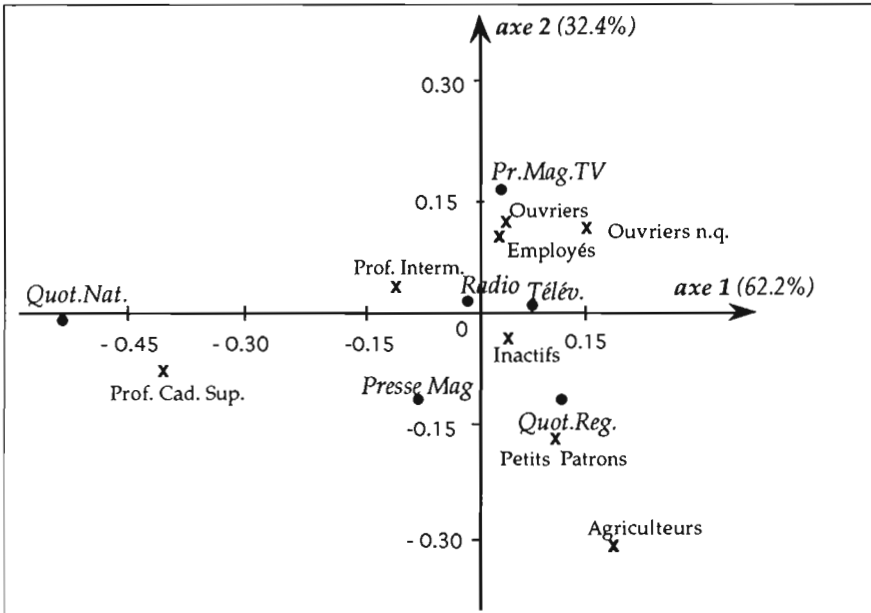


Figure 1.3 - 23  
Variables actives dans le premier plan factoriel

Il est clair dans une analyse de ce type que le premier axe correspond à une interprétation ponctuelle : les contacts média avec la presse quotidienne nationale sont, de façon significative, surtout le fait de cadres supérieurs et/ou de personnes d'un haut niveau d'éducation. Ce résultat n'est cependant pas d'emblée visible sur le tableau 1.3 - 10.

En revanche, les positions des points sur les deux figures donnent une interprétation plus nuancée du second axe : les professions salariées, de niveau d'éducation moyen, composées surtout de jeunes (contact média : Presse magazine TV), s'opposent aux petits patrons et agriculteurs, en moyenne sensiblement plus âgés et moins instruits (contacts : presse magazine autre que TV, et presse quotidienne régionale).

Que se passe-t-il si l'on supprime, au sein des colonnes actives, la colonne "Quot. N." dont le rôle est prédominant, pour la positionner en élément supplémentaire ?

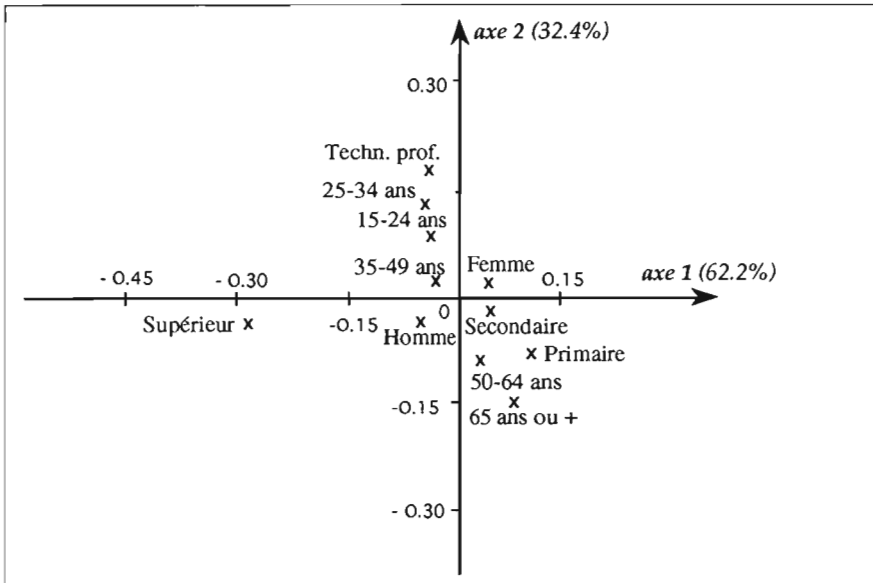


Figure 1.3 - 24

## Variables supplémentaires ou illustratives dans le premier plan factoriel

On a vu que cette colonne est presque située sur l'axe 1 (cosinus carré de 0.99). Sa suppression enlèverait 74.6% de l'inertie dans cette direction (valeur de la contribution), et donc l'inertie dans cette direction serait inférieure à celle du second axe actuel<sup>1</sup> sur lequel la colonne supprimée a d'ailleurs une contribution nulle. Donc le nouveau premier axe d'inertie maximale sera très voisin de l'ancien second axe. Tous calculs faits, on trouve, après suppression de la colonne en question, une première valeur propre de 0.0074 (la seconde valeur propre valait 0.0072) et des coordonnées sur ce nouveau premier axe qui diffèrent d'au plus de 0.01 de celles de l'ancien second axe. Le nouveau second axe (sur lequel la colonne supplémentaire "Presse Quot." a une coordonnée de 0.54 et un cosinus carré de 0.88) est très voisin de l'ancien premier axe.

Cet exemple aura illustré le positionnement de lignes supplémentaires et de colonnes supplémentaires, l'usage simultané des trois types d'aides à l'interprétation (valeurs propres, contributions, cosinus carrés) ainsi que le caractère itératif de l'analyse, qui fait penser à un "épluchage" progressif des nuages de points profils. L'exemple du paragraphe 2.4.4 illustrera aussi cette démarche en montrant la complémentarité de l'analyse factorielle avec la classification automatique.

<sup>1</sup> 25.4 % (complément à 100 de 74.6 %) de 0.0139 (première valeur propre  $\lambda_1$ ) est en effet très inférieur à 0.0072 (seconde valeur propre  $\lambda_2$ ).

## Analyse des Correspondances Multiples

L'analyse des correspondances introduite dans la section précédente peut se généraliser de plusieurs façons au cas où plus de deux ensembles sont mis en correspondance. Une des généralisations la plus simple et la plus utilisée est l'*analyse des correspondances multiples* qui permet de décrire de vastes tableaux binaires, dont les fichiers d'enquêtes socio-économiques constituent un exemple privilégié : les lignes de ces tableaux sont en général des individus ou observations (il peut en exister plusieurs milliers) ; les colonnes sont des modalités de variables nominales, le plus souvent des modalités de réponses à des questions. Il s'agit en fait d'une simple extension du domaine d'application de l'analyse des correspondances, avec cependant des procédures de calcul et des règles d'interprétation spécifiques.

On peut faire remonter les principes de cette méthode à Guttman (1941), mais aussi à Burt (1950) ou à Hayashi (1956). D'autres types d'extension ont été proposés par Benzécri (1973), Escofier-Cordier (1965), et par Masson (1974) qui s'appuie notamment sur les travaux de Carroll (1968), Horst (1961) et Kettenring (1971)<sup>1</sup>.

### 1.4.1 Domaine d'application

L'analyse des correspondances multiples est une analyse des correspondances simple appliquée non plus à une table de contingence, mais à un *tableau disjonctif complet*. Les propriétés d'un tel tableau sont intéressantes, les procédures de calculs et les règles d'interprétation des représentations obtenues sont simples et spécifiques.

L'extension du domaine d'application de l'analyse des correspondances se fonde sur l'équivalence suivante : si pour  $n$  individus, on dispose des valeurs (réponses) prises par deux variables nominales ayant respectivement  $p_1$  et  $p_2$  modalités, il est alors équivalent de soumettre à l'analyse des correspondances le tableau de contingence  $(p_1, p_2)$  croisant les deux variables ou d'analyser le tableau binaire à  $n$  lignes et  $(p_1 + p_2)$

---

<sup>1</sup> L'analyse des correspondances multiples a été développée également sur le nom d'*Homogeneity Analysis* par l'équipe de J. de Leeuw depuis 1973 (cf. Gifi, 1990) et sous le nom de *Dual Scaling* par Nishisato (1980). Une application de l'analyse des correspondances à un tableau disjonctif complet se trouve dans Nakache (1973). L'ensemble des résultats et propriétés présentés dans cette section figurent, avec des programmes et des exemples, dans Lebart et Tabard (1973). Le nom d'*analyse des correspondances multiples* figure pour la première fois dans Lebart (1975 a). Un exposé synthétique de ces diverses approches a été réalisé par Tenenhaus et Young (1985).

colonnes décrivant les réponses. L'analyse de ce dernier tableau est plus coûteuse, mais plus intéressante, car elle se généralise immédiatement au cas de plus deux variables nominales.

### 1.4.2 Notations et définitions

Une partie généralement importante des fichiers d'enquête se compose de réponses à des questions mises sous forme *disjonctive complète* : les diverses modalités de réponses s'excluent mutuellement et une modalité est obligatoirement choisie.

Par exemple à la question :

*Etes-vous ?*

1- célibataire,

2- marié(e) ou vivant maritalement,

3- veuf(ve),

4- divorcé(e),

5- non réponse,

cinq modalités de réponses (dont une non-réponse) sont possibles.

Une variable continue peut être transformée en variable nominale par le découpage en classes des valeurs de la variable. Par exemple, à la question "âge de l'enquêté", on prévoit 8 modalités de réponse :

1- moins de 25 ans;

2- de 25 à 29 ans;

3- de 30 à 34 ans;

4- de 35 à 39 ans;

5- de 40 à 44 ans;

6- de 45 à 49 ans;

7- de 50 ans et plus;

8- non-réponse.

Si l'on désigne par  $s$  le nombre des questions posées à  $n$  individus, on dispose ainsi d'un tableau de données  $\mathbf{R}$  ayant  $n$  lignes et  $s$  colonnes mis sous forme de codage condensé, illustré sur la figure 1.4 - 1 par un tableau pour lequel  $s = 3$  et  $n = 12$ .

Le terme général  $r_{iq}$  désigne la modalité de la question  $q$  choisie par le sujet  $i$ . En notant  $p_q$  le nombre des modalités de réponses à une question  $q$ , on a :  $r_{iq} \leq p_q$ .

Mais un tel tableau n'est pas exploitable : les sommes en ligne et en colonne n'ont pas de sens. Il faut recoder les variables.

	$s=3$		
1	2	2	4
	2	1	3
	3	1	2
	1	2	4
	1	2	3
	2	2	3
	3	1	1
	1	1	1
	2	1	2
	2	2	3
	3	2	2
n	1	1	4

$\mathbf{R} = (n, s)$

Figure 1.4 - 1

Tableau de données sous forme de codage condensé

### a – Hypercube de contingence

Pour disposer de toute l'information, on peut construire l'hypercube de contingence  $H$  croisant les  $s$  questions et dont les éléments constituent l'éventail des réponses possibles des sujets enquêtés. On dispose d'un ensemble-produit des modalités des  $s$  questions dont les éléments sont constitués des suites de  $s$  modalités, chacune étant prise dans une question différente.

Pour  $s=3$  questions ayant respectivement 3, 2 et 4 modalités, il existe 24 combinaisons possibles de réponses selon lesquelles sont réparties les individus. Dans le cas de deux questions, l'hypertable est le tableau de contingence. Pour un nombre important de questions, l'hypertable sera en général presque vide. Si l'on pose à 1000 individus 12 questions ayant chacune 10 modalités de réponse, le nombre de réponses possibles distinctes vaut  $10^{12}$ . Au plus une case sur un milliard de l'hypertable ne sera pas vide.

### b – Tableau disjonctif complet

On désigne par  $I$  l'ensemble des  $n$  sujets ayant répondu au questionnaire et par  $p$  le nombre total des modalités des  $s$  questions. On a :

$$p = \sum_{q=1}^s p_q$$

On construit, à partir du tableau de données  $R$ , le tableau  $Z$  à  $n$  lignes et  $p$  colonnes décrivant les  $s$  réponses des  $n$  individus par un codage binaire. Le tableau  $Z$  est la juxtaposition de  $s$  sous-tableaux :

$$Z = [Z_1, Z_2, \dots, Z_q, \dots, Z_s]$$

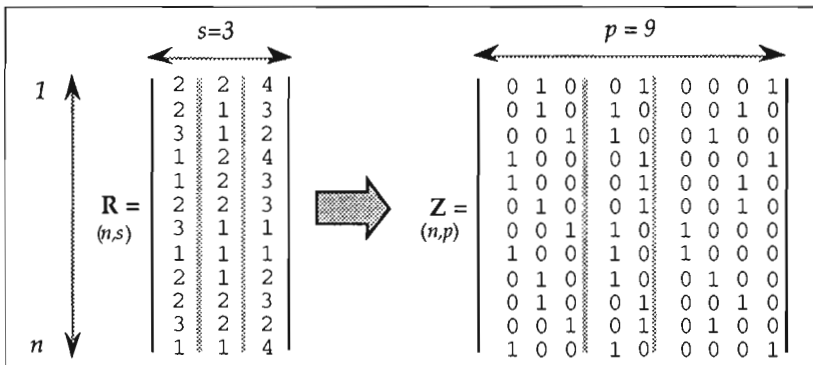


Figure 1.4 - .2  
Construction du tableau disjonctif complet  $Z$

Le sous-tableau  $Z_q$ , à  $n$  lignes et  $p_q$  colonnes, est tel que sa  $i^{\text{ème}}$  ligne contient  $p_q - 1$  fois la valeur 0 et une fois la valeur 1 dans la colonne correspondant à la modalité de la question  $q$  choisie par le sujet  $i$ . Autrement dit le tableau  $Z_q$  décrit la partition des  $n$  individus induite par les réponses à la question  $q$ .

Le tableau  $\mathbf{Z}$  est appelé *tableau disjonctif complet* dont le terme général s'écrit :

$$z_{ij} = 1 \quad \text{ou} \quad z_{ij} = 0$$

selon que le sujet  $i$  a choisi la modalité  $j$  de la question  $q$  ou non.

Les marges en ligne du tableau disjonctif complet sont constantes et égales

au nombre  $s$  de questions :

$$z_{i.} = \sum_{j=1}^p z_{ij} = s$$

Les marges en colonne :

$$z_{.j} = \sum_{i=1}^n z_{ij} \quad \text{correspondent au nombre de}$$

sujets ayant choisi la modalité  $j$  de la question  $q$ .

On vérifie que, pour chaque sous-tableau  $\mathbf{Z}_q$ , l'effectif total est bien :

$$z_q = \sum_{j \in q} z_{.j} = n$$

La somme des marges donne l'effectif total  $z$  du tableau  $\mathbf{Z}$  soit :

$$z = \sum_{i=1}^n \sum_{j=1}^p z_{ij} = ns$$

### c – Tableau des faces de l'hypercube de contingence ou tableau de contingence de Burt

L'ensemble des  $p_q$  modalités de réponse à une question permet de partitionner l'échantillon en au plus  $p_q$  classes. La donnée de deux questions mises sous forme disjonctive complète permet de réaliser deux partitions de l'ensemble des individus enquêtés et l'on obtient un tableau de contingence. L'analyse du tableau croisant les deux partitions peut être généralisée au cas de  $s$  partitions,  $s$  étant un entier supérieur à 2.

On construit, à partir du tableau disjonctif complet  $\mathbf{Z}$ , le tableau symétrique  $\mathbf{B}$  d'ordre  $(p,p)$  qui rassemble les croisements deux à deux de toutes les variables :

$$\mathbf{B} = \mathbf{Z}'\mathbf{Z}$$

$\mathbf{B}$  est appelé *tableau de contingence de Burt*<sup>1</sup> associé au tableau disjonctif complet  $\mathbf{Z}$ .

Le terme général de  $\mathbf{B}$  s'écrit :

$$b_{jj'} = \sum_{i=1}^n z_{ij}z_{ij'}$$

$\mathbf{B}$  est une juxtaposition de tableaux de contingence.

---

<sup>1</sup> Sir Cyril Burt a été un incontestable innovateur au point de vue méthodologique (cf. son article précité de 1950, dans lequel il préconise le calcul de  $\mathbf{B}$ , et sa diagonalisation après une normalisation qui correspond à celle de l'analyse des correspondances multiples). Il est peut-être encore plus célèbre pour les falsifications d'observations et les graves fraudes scientifiques et déontologiques dont il a été l'auteur.

Les marges sont pour tout  $j \leq p$  :

$$b_j = \sum_j^p b_{jj'} = sz_{.j}$$

et l'effectif total  $b$  vaut :

$$b = s^2 n$$

Le tableau  $\mathbf{B}$  est formé de  $s^2$  blocs où l'on distingue :

- le bloc  $\mathbf{Z}'_q \mathbf{Z}_{q'}$  indicé par  $(q, q')$ , d'ordre  $(p_q, p_{q'})$  qui n'est autre que la table de contingence croisant les réponses aux questions  $q$  et  $q'$ .
- le  $q^{\text{ième}}$  bloc carré  $\mathbf{Z}'_q \mathbf{Z}_q$  obtenu par le croisement d'une variable avec elle-même. C'est une matrice d'ordre  $(p_q, p_q)$ , diagonale puisque deux modalités d'une même question ne peuvent être choisies simultanément. Les termes diagonaux sont les effectifs des modalités de la question  $q$ .

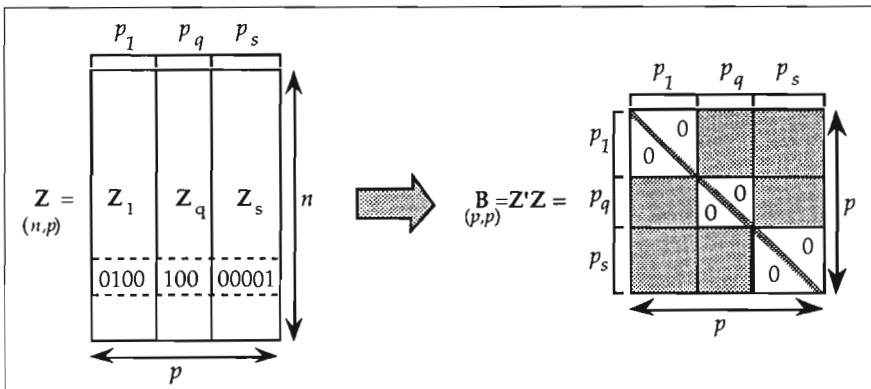


Figure 1.4 - 3  
Construction du tableau des faces de l'hypercube (tableau de Burt)  $\mathbf{B}$   
à partir du tableau disjonctif complet  $\mathbf{Z}$

Nous désignerons par  $\mathbf{D}$  la matrice diagonale, d'ordre  $(p, p)$  ayant les mêmes éléments diagonaux que  $\mathbf{B}$  ; ces éléments sont les effectifs correspondant à chacune des modalités (cf. figure 1.4 - 4) :

$$d_{jj} = b_{jj} = z_{.j}$$

$$d_{jj'} = 0 \quad \text{pour tout } j' \neq j$$

La matrice  $\mathbf{D}$  peut être également considérée comme formée de  $s^2$  blocs.

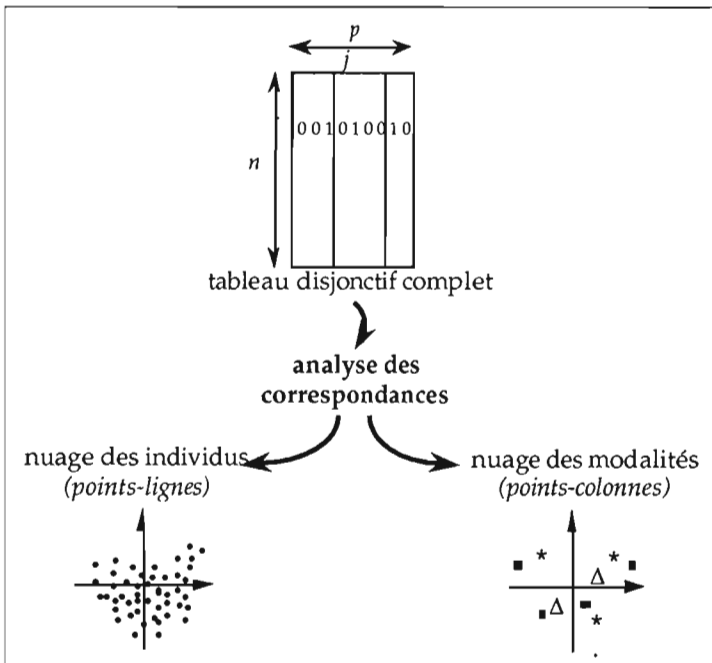
Seules les  $s$  matrices diagonales  $\mathbf{D}_q = \mathbf{Z}'_q \mathbf{Z}_q$  ( $q = 1, \dots, s$ ) constituant les blocs diagonaux de  $\mathbf{B}$  sont des matrices non nulles.

		$p = 9$																		
$\mathbf{B} =$ <small>(p,p)</small>	4	0	0	2	2	1	0	1	2		$\mathbf{D} =$ <small>(p,p)</small>	4	0	0	0	0	0	0	0	0
	0	5	0	2	3	0	1	3	1			0	5	0	0	0	0	0	0	0
	0	0	3	2	1	1	2	0	0			0	0	3	0	0	0	0	0	0
	2	2	2	6	0	2	2	1	1			0	0	0	6	0	0	0	0	0
	2	3	1	0	6	0	1	3	2			0	0	0	0	6	0	0	0	0
	1	0	1	2	0	2	0	0	0			0	0	0	0	0	2	0	0	0
	0	1	2	2	1	0	3	0	0			0	0	0	0	0	0	3	0	0
	1	3	0	1	3	0	0	4	0			0	0	0	0	0	0	0	4	0
	2	1	0	1	2	0	0	0	3			0	0	0	0	0	0	0	0	3

**Figure (1.4 - 4)**  
**Tableau de Burt B et matrice diagonale D associée**  
*(données des figures 1.4 - 1 et 1.4 - 2)*

### 1.4.3 Principes de l'analyse des correspondances multiples

L'analyse des correspondances multiples est l'analyse des correspondances d'un tableau disjonctif complet.



**Figure 1.4 - 5**  
**Analyse des correspondances multiples**

Ses principes sont donc ceux de l'analyse des correspondances à savoir :



- mêmes transformations du tableau de données en profils-lignes et en profils-colonnes;
- même critère d'ajustement avec pondération des points par leurs profils marginaux;
- même distance, celle du  $\chi^2$ .

L'analyse des correspondances multiples présente cependant des propriétés particulières dues à la nature même du tableau disjonctif complet. Nous allons énoncer les principes de cette analyse à partir du tableau disjonctif complet puis nous montrerons l'équivalence avec l'analyse du tableau de Burt.

### a – Critère d'ajustement et distance du $\chi^2$

Les individus sont tous affectés d'une masse identique égale à  $m_i = \frac{1}{n}$  et

chacune des modalités  $j$  est pondérée par sa fréquence  $m_j = \frac{z_{.j}}{ns}$ .

La distance du  $\chi^2$  appliquée à un tableau disjonctif complet conserve un sens. En effet, dans  $\mathbb{R}^n$ , la distance entre modalités s'écrit :

$$d^2(j, j') = \sum_{i=1}^n n \left( \frac{z_{ij}}{z_{.j}} - \frac{z_{ij'}}{z_{.j'}} \right)^2$$

Ainsi deux modalités choisies par les mêmes individus coïncident. Par ailleurs, les modalités de faible effectif sont éloignées des autres modalités.

Dans  $\mathbb{R}^p$ , la distance entre deux individus  $i$  et  $i'$  s'exprime par :

$$d^2(i, i') = \frac{1}{s} \sum_{j=1}^p \frac{n}{z_{.j}} (z_{ij} - z_{i'j})^2$$

Deux individus sont proches s'ils ont choisi les mêmes modalités. Ils sont éloignés s'ils n'ont pas répondu de la même manière<sup>1</sup>.

### b – Axes factoriels et facteurs

En reprenant les résultats de l'analyse des correspondances et les notations adoptées (cf. § 1.3.3.b), on pose<sup>2</sup> :

$$\mathbf{F} = \frac{1}{ns} \mathbf{Z} \quad \text{de terme général} \quad f_{ij} = \frac{z_{ij}}{ns}$$

<sup>1</sup> On note qu'une modalité  $j$  intervient d'autant plus dans le calcul de la distance entre deux individus que sa masse est plus faible.

<sup>2</sup>  $\mathbf{I}_n$  est la matrice identité d'ordre  $(n, n)$  et  $\delta_{ij}$  est tel que :

$$\delta_{ij} = 1 \quad \text{si} \quad i = j \quad \text{et} \quad \delta_{ij} = 0 \quad \text{si} \quad i \neq j$$

$$\mathbf{D}_p = \frac{1}{ns} \mathbf{D} \quad \text{de terme général} \quad f_{.j} = \delta_{ij} \frac{z_{.j}}{ns}$$

$$\mathbf{D}_n = \frac{1}{n} \mathbf{I}_n \quad \text{de terme général} \quad f_{i.} = \frac{\delta_{ij}}{n}$$

Pour trouver les axes factoriels  $\mathbf{u}_\alpha$  on diagonalise la matrice :

$$\mathbf{S} = \mathbf{F}' \mathbf{D}_n^{-1} \mathbf{F} \mathbf{D}_p^{-1} = \frac{1}{s} \mathbf{Z}' \mathbf{Z} \mathbf{D}^{-1}$$

de terme général (attention,  $s$  [sans indice] désigne le nombre de questions dans ce chapitre):

$$s_{jj'} = \frac{1}{s z_{.j'}} \sum_{i=1}^n z_{ij} z_{ij'}$$

Dans  $\mathbb{R}^p$ , l'équation du  $\alpha^{\text{ième}}$  axe factoriel  $\mathbf{u}_\alpha$  est :

$$\frac{1}{s} \mathbf{Z}' \mathbf{Z} \mathbf{D}^{-1} \mathbf{u}_\alpha = \lambda_\alpha \mathbf{u}_\alpha \quad [1.4 - 1]$$

L'équation du  $\alpha^{\text{ième}}$  facteur  $\varphi_\alpha = \mathbf{D}^{-1} \mathbf{u}_\alpha$  s'écrit :

$$\frac{1}{s} \mathbf{D}^{-1} \mathbf{Z}' \mathbf{Z} \varphi_\alpha = \lambda_\alpha \varphi_\alpha \quad [1.4 - 2]$$

De même, l'équation du  $\alpha^{\text{ième}}$  facteur  $\psi_\alpha$  dans  $\mathbb{R}^n$  s'écrit :

$$\frac{1}{s} \mathbf{Z} \mathbf{D}^{-1} \mathbf{Z}' \psi_\alpha = \lambda_\alpha \psi_\alpha$$

Les facteurs  $\varphi_\alpha$  et  $\psi_\alpha$  (de norme  $\lambda_\alpha$ ) représentent les coordonnées des points-lignes et des points-colonnes sur l'axe factoriel  $\alpha$ .

Les relations de transition entre les facteurs  $\varphi_\alpha$  et  $\psi_\alpha$  sont :

$$\begin{cases} \varphi_\alpha = \frac{1}{\sqrt{\lambda_\alpha}} \mathbf{D}^{-1} \mathbf{Z}' \psi_\alpha \\ \psi_\alpha = \frac{1}{s\sqrt{\lambda_\alpha}} \mathbf{Z} \varphi_\alpha \end{cases}$$

### c – Facteurs et relations quasi-barycentriques

La coordonnée factorielle de l'individu  $i$  sur l'axe  $\alpha$  est donnée par :

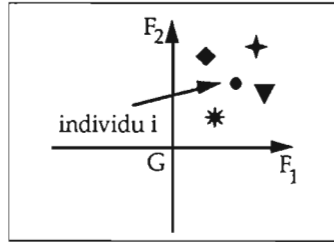
$$\psi_{\alpha i} = \frac{1}{\sqrt{\lambda_\alpha}} \sum_{j=1}^p \frac{z_{ij}}{z_{i.}} \varphi_{\alpha j}$$

c'est-à-dire :

$$\psi_{\alpha i} = \frac{1}{s\sqrt{\lambda_\alpha}} \sum_{j \in p(i)} \varphi_{\alpha j} \quad [1.4 - 3]$$

où  $p(i)$  désigne l'ensemble des modalités choisies par l'individu  $i$ .

Au coefficient  $\frac{1}{\sqrt{\lambda_\alpha}}$  près, l'individu  $i$  se trouve au point moyen du nuage des modalités qu'il a choisies.



**Figure 1.4 - 6**  
Projection d'un individu  
au point moyen des modalités choisies

De même, la coordonnée de la modalité  $j$  sur l'axe  $\alpha$  est donnée par :

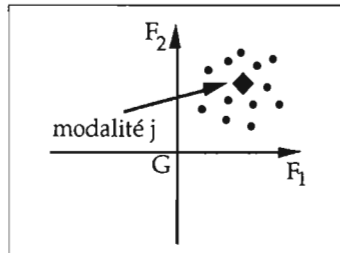
$$\varphi_{\alpha j} = \frac{1}{\sqrt{\lambda_\alpha}} \sum_{i=1}^n \frac{z_{ij}}{z_{.j}} \psi_{\alpha i}$$

c'est-à-dire :

$$\varphi_{\alpha j} = \frac{1}{z_{.j} \sqrt{\lambda_\alpha}} \sum_{i \in I(j)} \psi_{\alpha i} \quad [1.4 - 4]$$

où  $I(j)$  désigne l'ensemble des individus ayant choisi la modalité  $j$ .

Avant la dilatation sur l'axe  $\alpha$ , la modalité  $j$  se trouve au point moyen du nuage des individus qui l'ont choisie comme réponse.



**Figure 1.4 - 7**  
Projection d'une modalité  
au point moyen des individus concernés

Dans le calcul des relations quasi-barycentriques [1.4 - 4], les individus ne sont pas pondérés. Il s'agit de simples calculs de moyennes arithmétiques de coordonnées.

### d – Sous-nuage des modalités d'une même variable

Le nuage des modalités dans  $\mathbb{R}^n$  peut être décomposé en  $s$  sous-nuages, le  $q^{\text{ème}}$  correspondant à l'ensemble des  $p_q$  modalités de la variable  $q$ . Ces sous-nuages ont même centre de gravité  $G$  qui est celui du nuage global.

En effet, les coordonnées des points du sous-nuage relatif à la variable  $q$  sont les colonnes de  $Z_q D_q^{-1}$  et les éléments diagonaux de  $\frac{1}{n} D_q$  sont les masses relatives des  $p_q$  points de ce sous-nuage. Puisque :

$$\sum_{j \in p_q} z_{ij} = 1$$

alors la  $i^{\text{ème}}$  composante du centre de gravité du sous-nuage vaut :

$$G_{qi} = \sum_{j \in p_q} \frac{d_{jj}}{n} \frac{z_{ij}}{d_{jj}} = \frac{1}{n} = G_i$$

où il apparaît que  $G_{qi}$  ne dépend pas de  $q$ .

Les composantes  $\varphi_q$  des modalités d'une variable  $q$  (relatives aux facteurs non-triviaux  $\varphi$ ) sont centrées puisque ces facteurs correspondent à une analyse du nuage après translation de l'origine en  $G$ . Les facteurs opposent les modalités d'une même variable.

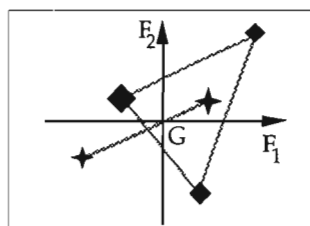


Figure 1.4 - 8  
Composantes centrées

#### Remarques

- 1) Si le tableau disjonctif n'est pas complet (c'est-à-dire si, pour au moins un individu, aucune modalité de réponse à une question n'a été choisie), les modalités d'une même variable ne sont plus centrées sur le centre de gravité du nuage global.
- 2) Le codage disjonctif complet permet de transformer une variable continue en une variable nominale dont les modalités sont des classes ordonnées. Il est alors utile de tracer la trajectoire qui relie les classes, trajectoire qui peut suggérer des liaisons non linéaires entre cette variable et les axes.

### e – Support du nuage des modalités

Les coordonnées des modalités dans  $\mathbb{R}^n$  sont les colonnes de  $Z D^{-1}$ . Elles engendrent un sous-espace dont la dimension est le rang de  $Z D^{-1}$ , donc le rang de  $Z = [Z_1, Z_2, \dots, Z_q, \dots, Z_s]$ .

Tous les sous-espaces engendrés par les  $Z_q$  ont en commun la première bissectrice notée  $\Delta$ . Le rang maximum de  $Z$  est donc :

$$p_1 + (p_2 - 1) + \dots + (p_s - 1) = p - s + 1$$

Le rang maximum de la matrice à diagonaliser  $D^{-1}Z'Z$  sera donc  $p - s + 1$ . Mais dans l'analyse du nuage par rapport à l'origine  $O$ , la première bissectrice est vecteur propre correspondant à la valeur propre 1 (le nuage est contenu dans le sous-espace  $D^{-1}$ -orthogonal à  $\Delta$ ).

Dans l'analyse par rapport au centre de gravité  $G$ , on trouvera donc  $p - s$  valeurs propres non nulles. En choisissant une base dans le support du nuage, on pourra se ramener à la recherche d'éléments propres d'une matrice d'ordre  $p - s$ .

## f – Meilleure représentation simultanée

La présentation de l'analyse des correspondances peut être formulée ici de façon particulière en raison du codage spécifique au tableau disjonctif complet.

Nous cherchons sur un même axe les coordonnées des  $n$  individus et des  $p$  modalités de façon que :

- la coordonnée d'un individu  $i$  soit la moyenne arithmétique des coordonnées des modalités qu'il a choisies (à une dilatation  $\beta$  près, que l'on s'efforcera de rendre minimale).
- la coordonnée d'une modalité  $j$  soit la moyenne arithmétique des coordonnées des individus qui l'ont choisie (à une même dilatation  $\beta$  près).

Bien entendu, on obtient les relations dite quasi-barycentriques issues de l'analyse du tableau disjonctif complet  $Z$  avec, pour le coefficient de

dilatation  $\beta$ , la valeur minimale  $\beta = \frac{1}{\sqrt{\lambda}}$  :

$$\begin{cases} \varphi = \frac{1}{\sqrt{\lambda}} D^{-1} Z' \psi \\ \psi = \frac{1}{s\sqrt{\lambda}} Z \varphi \end{cases}$$

La représentation simultanée des individus et des modalités est importante pour l'interprétation des résultats. Cependant elle n'est pratiquement pas utilisée, d'une part pour des raisons d'encombrement graphique (on dispose souvent de plusieurs centaines voire de plusieurs milliers d'individus) et d'autre part parce que les individus sont, dans la plupart des applications, anonymes. Ils ne présentent de l'intérêt que par l'intermédiaire de leurs caractéristiques. On peut cependant vouloir projeter les individus sur un plan factoriel afin d'apprécier leur répartition et les zones de densité.

### g – Inertie du nuage des modalités et conséquences pratiques

On rappelle que la distance du  $\chi^2$  dans  $\mathbb{R}^n$  est la métrique  $D_n^{-1}$ .

La distance entre la modalité  $j$  et le centre de gravité du nuage  $G$ , dont toutes les  $n$  coordonnées valent  $\frac{1}{n}$ , s'écrit :

$$d^2(j, G) = n \sum_{i=1}^n \left( \frac{z_{ij}}{z_{.j}} - \frac{1}{n} \right)^2 = \frac{n}{z_{.j}} - 1$$

La distance d'une modalité au centre de gravité est d'autant plus grande que l'effectif est plus faible.

#### - Inertie d'une modalité

L'inertie  $I(j)$  de la modalité  $j$  vaut :

$$I(j) = m_j d^2(j, G)$$

avec :

$$m_j = \frac{z_{.j}}{ns}$$

d'où :

$$I(j) = \frac{1}{s} \left( 1 - \frac{z_{.j}}{n} \right)$$

La part d'inertie due à une modalité de réponse est d'autant plus grande que l'effectif dans cette modalité est plus faible.

Le maximum  $\frac{1}{s}$  serait atteint par une modalité d'effectif nul. En conséquence, on évite, au moment du codage, les modalités à faibles effectifs susceptibles de perturber les directions des premiers axes factoriels.

#### - Inertie d'une question

L'inertie de la question  $q$ , notée  $I(q)$ , vaut :

$$I(q) = \sum_{j=1}^{p_q} I(j) = \frac{1}{s} (p_q - 1)$$

Ainsi la part d'inertie due à une question est fonction croissante du nombre de modalités de réponse.

La part minimale  $\frac{1}{s}$  correspond aux questions à 2 modalités. D'où l'intérêt d'équilibrer le système des questions, c'est-à-dire le découpage des variables en modalités, si on veut faire jouer le même rôle à toutes les questions.

**- Inertie totale**

On en déduit que l'inertie totale  $I$  vaut :

$$I = \sum_q I(q) = \sum_{j=1}^p \frac{z_j}{ns} d^2(j, G)$$

d'où :

$$I = \frac{p}{s} - 1$$

En particulier, elle vaut 1 dans le cas où toutes les questions ont deux modalités de réponse (cas où  $p=2s$ ). On verra au paragraphe 1.4.7.a que dans ce cas, analyse des correspondances multiples et analyse en composantes principales donnent des résultats équivalents.

L'inertie totale dépend uniquement du nombre de variables et de modalités et non des liaisons entre les variables. C'est une quantité qui, dans le cadre de l'analyse des correspondances multiples (comme dans celui de l'analyse en composantes principales normée), n'a pas de signification statistique.

**h – Règles d'interprétation**

Dire qu'il existe des affinités entre réponses, c'est dire aussi qu'il existe des individus qui ont choisi simultanément toutes ou presque toutes ces réponses.

L'analyse des correspondances multiples met alors en évidence des types d'individus ayant des profils semblables quant aux attributs choisis pour les décrire. Compte tenu des distances entre les éléments du tableau disjonctif complet et des relations barycentriques particulières, on exprime :

- *la proximité entre individus en terme de ressemblances :*  
deux individus se ressemblent s'ils ont choisi globalement les mêmes modalités.
- *la proximité entre modalités de variables différentes en terme d'association :*  
ces modalités correspondent aux points moyens des individus qui les ont choisies et sont proches parce qu'elles concernent globalement les mêmes individus ou des individus semblables.
- *la proximité entre deux modalités d'une même variable en terme de ressemblance :*  
par construction, les modalités d'une même variable s'excluent. Si elles sont proches, cette proximité s'interprète en terme de ressemblance entre les groupes d'individus qui les ont choisies (vis-à-vis d'autres variables actives de l'analyse).

Les règles d'interprétation des résultats (coordonnées, contributions, cosinus carrés) concernant les éléments actifs d'une analyse des correspondances multiples sont sensiblement les mêmes que celles d'une analyse des

correspondances simple (cf. § 1.3.5). On calcule la contribution et la qualité de représentation de chaque modalité et de chaque individu, si ceux-ci ne sont pas anonymes pour l'analyse.

Cependant, la notion de variable doit être prise en compte au moment de l'interprétation, ceci au travers de ses modalités. Compte tenu de la décomposition de l'inertie du nuage des modalités, on calcule la contribution d'une variable au facteur  $\alpha$  en sommant les contributions de ses modalités sur ce facteur :

$$Cr_{\alpha}(q) = \sum_{j \in q} Cr_{\alpha}(j)$$

On repère ainsi, en plus des modalités responsables des axes factoriels, les variables qui ont participé à la définition du facteur. On obtient un indicateur de liaison entre la variable et le facteur [cf. Escofier, 1979 c].

En revanche, les règles d'interprétation des valeurs propres et des taux d'inertie sont différentes (on a vu que la trace n'avait plus d'interprétation statistique). On se reportera au chapitre 4 sur la validité et portée des résultats pour plus de détails.

### **i – Principes du découpage en classes**

Les variables continues, pour être actives dans une analyse des correspondances multiples, doivent être soit rendues nominales (découpées en classes), soit recodées selon deux colonnes numériques<sup>1</sup>.

Lorsque l'on cherche ainsi à découper une variable en classes, on est confronté à plusieurs problèmes : combien de classes choisir et comment les choisir ? Où placer les bornes des classes d'une variable continue ? La consultation de la distribution de chaque variable (tris-à-plat et histogrammes) est indispensable pour effectuer ces choix.

Certains principes, déduits des propriétés de l'analyse des correspondances multiples (cf. § 1.4.3.g), peuvent être utilisés pour guider la phase de recodage : constituer des modalités d'effectifs semblables, découper les variables de manière à avoir un nombre comparable de modalités. Pour donner un ordre de grandeur, un découpage entre 4 à 8 modalités convient dans la plupart des applications.

Il s'agit par conséquent de trouver un compromis entre un découpage techniquement acceptable selon ces principes et un découpage qui exhibe au mieux l'information à retenir. On ne peut généralement pas avoir recours à des algorithmes aveugles pour élaborer un découpage satisfaisant<sup>2</sup>. On

<sup>1</sup> Cf. le recodage préconisé par Escofier (1979 b) présenté au § 3.8.5.c.

<sup>2</sup> L'algorithme de Fisher (1958) fournit une partition optimale exacte (critère variance inter/variance totale maximal), mais ce critère rend très mal compte des mélanges de distributions ayant des variances très inégales et ne sépare donc pas des classes qu'une inspection visuelle d'histogramme distinguerait sans hésiter.



retiendra par exemple une modalité de faible effectif si celle-ci est importante pour l'étude. De même pour sélectionner les bornes des classes d'une variable continue, on respectera un ou plusieurs seuils naturels dans le contexte de l'étude, ou significatifs après examen de l'histogramme (le découpage en classes d'amplitudes égales est parfois inapproprié).

Ces principes sont moins rigoureux pour une variable supplémentaire. N'intervenant pas dans la formation des facteurs ou des classes, on a parfois intérêt à effectuer un découpage fin pour les variables supplémentaires.

La transformation de variables continues en variables nominales occasionne une perte de l'information brute mais présente certains avantages : exploiter simultanément des variables nominales et continues en correspondances multiples ; valider a posteriori les données en permettant d'observer l'éventuelle contiguïté des classes voisines ; et mettre en évidence les éventuelles liaisons non linéaires entre variables continues.

Pour un exposé de synthèse sur les méthodes de codage, on consultera Cazes (1990), Grelet (1993). L'article précité de Cazes et les travaux de Gallego (1982), van Rijckevorsel (1987) portent en particulier sur l'utilisation du codage flou en analyse des correspondances.

#### 1.4.4 Éléments supplémentaires

L'utilisation des éléments supplémentaires en analyse des correspondances multiples permet de prendre en compte toute l'information susceptible d'aider à comprendre ou à interpréter la typologie induite par les éléments actifs.

Ceci est particulièrement intéressant lorsque l'ensemble des variables se décompose en thème, c'est-à-dire en groupes de variables homogènes quant à leur contenu.

Dans l'analyse du tableau disjonctif complet, on fera intervenir des éléments supplémentaires pour :

- Enrichir l'interprétation des axes par des variables n'ayant pas participé à leur construction. On projettera alors dans l'espace des variables les centres de groupes d'individus définis par les modalités des variables supplémentaires.
- Adopter une optique de prévision en projetant les variables supplémentaires dans l'espace des individus. Celles-ci seront "expliquées" par les variables actives. On peut projeter des individus supplémentaires dans l'espace des variables, pour les situer par rapport aux individus actifs ou par rapport à des groupes d'individus actifs dans une optique de discrimination (cf. section 3.3).

Suivant la nature des variables supplémentaires, nominales ou continues, on interprète différemment leur position sur les axes factoriels.

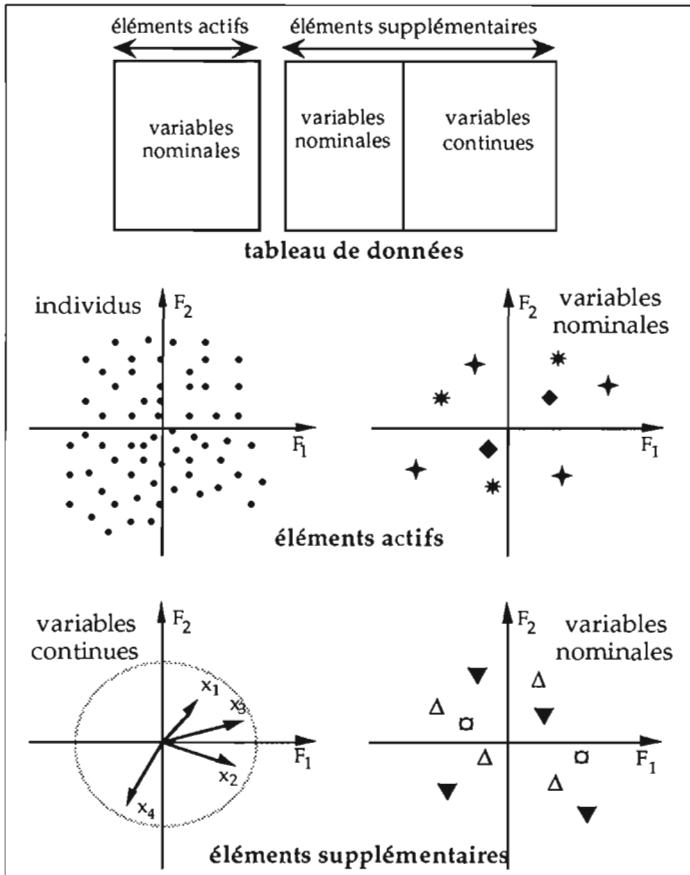


Figure 1.4 - 9  
Représentation des variables supplémentaires  
en analyse des correspondances multiples

**a – Valeurs-test pour les modalités supplémentaires**

Tout comme pour l'analyse des correspondances simples, il n'est pas nécessaire de projeter en supplémentaire toutes les modalités d'une variable nominale.

La coordonnée factorielle  $\varphi_{\alpha j}$  d'une modalité  $j$  sur un axe  $\alpha$  (que cette modalité figure parmi les variables actives ou qu'elle soit supplémentaire) est le produit par le coefficient  $\frac{1}{\sqrt{\lambda_\alpha}}$  de la moyenne arithmétique des coordonnées  $\psi_{\alpha i}$  des individus ayant choisi cette modalité  $j$  de réponse :

$$\varphi_{\alpha j} = \frac{1}{z_j \sqrt{\lambda_\alpha}} \sum_{i \in I(j)} \psi_{\alpha i}$$

où  $I(j)$  est l'ensemble des individus ayant choisi la modalité  $j$ . Ceci suggère alors le test d'hypothèse suivant.

Supposons qu'une modalité supplémentaire  $j$  concerne  $n_j$  individus ( $n_j = z.j$ ). Si ces  $n_j$  individus sont tirés au hasard (hypothèse nulle  $H_0$ ) parmi les  $n$  individus analysés (tirage supposé sans remise), la moyenne de  $n_j$  coordonnées tirées au hasard dans l'ensemble fini des  $n$  valeurs  $\psi_{\alpha i}$  est une variable aléatoire  $X_{\alpha j}$ :

$$X_{\alpha j} = \frac{1}{n_j} \sum_{i \in I(j)} \psi_{\alpha i}$$

avec pour espérance :

$$E(X_{\alpha j}) = 0$$

et pour variance<sup>1</sup> :

$$\text{Var}_{H_0}(X_{\alpha j}) = \frac{n - n_j}{n - 1} \frac{\lambda_{\alpha}}{n_j}$$

La coordonnée  $\varphi_{\alpha j}$  de la modalité supplémentaire est liée à la variable aléatoire  $X_{\alpha j}$  par la relation :

$$\varphi_{\alpha j} = \frac{1}{\sqrt{\lambda_{\alpha}}} X_{\alpha j}$$

On a donc :

$$E(\varphi_{\alpha j}) = 0$$

et :

$$\text{Var}(\varphi_{\alpha j}) = \frac{n - n_j}{n - 1} \frac{1}{n_j}$$

La quantité  $t_{\alpha j}$  :

$$t_{\alpha j} = \sqrt{n_j \frac{n - n_j}{n - 1}} \varphi_{\alpha j}$$

mesure en nombre d'écart-types la distance entre la modalité  $j$ , c'est-à-dire le quasi-barycentre des  $n_j$  individus, et l'origine sur l'axe factoriel  $\alpha$ . On appelle cette quantité "valeur-test". D'après le théorème de la limite centrale, sa distribution tend vers une loi de Laplace-Gauss centrée réduite.

Ainsi, la position d'une modalité est intéressante dans une direction  $\alpha$  donnée si le sous-nuage qu'elle constitue occupe une zone étroite dans cette direction et si cette zone est éloignée du centre de gravité du nuage.

La valeur-test est un critère qui permet d'apprécier rapidement si une modalité a une position "significative" sur un axe. On considère généralement comme occupant une "position significative" les modalités dont les valeurs-test sont supérieures à 2 en valeur absolue, correspondant approximativement au seuil 5%.

<sup>1</sup> Il s'agit de la formule classique donnant la variance d'une moyenne lors d'un tirage sans remise de  $n_j$  objets parmi  $n$ , en fonction de la variance totale  $\lambda_{\alpha}$ .

Le calcul simultané de plusieurs valeurs-test ou de plusieurs seuils de probabilités se heurte à l'écueil des *comparaisons multiples*, bien connu des statisticiens.

Supposons que l'on projette 100 modalités supplémentaires qui soient vraiment tirées au hasard. Les valeurs-test attachées à ces modalités sont alors toutes des réalisations de variables aléatoires normales centrées réduites indépendantes.

Dans ces conditions, en moyenne, sur 100 valeurs-test calculées, 5 seront en dehors de l'intervalle  $[-1.96, +1.96]$ , et 5 dépasseront la valeur 1.65 (test unilatéral). Le seuil de 5% n'a de sens en fait que pour un seul test, et non pour des tests multiples. On résout de façon pragmatique cette difficulté en choisissant un seuil plus sévère<sup>1</sup>.

On note que les valeurs-test n'ont de sens que pour les modalités supplémentaires ou encore pour les modalités actives ayant des contributions absolues faibles, c'est-à-dire se comportant comme des modalités supplémentaires<sup>2</sup>.

Lorsque l'on dispose d'un nombre important de modalités supplémentaires, les valeurs-test permettent de repérer rapidement les modalités utiles à l'interprétation d'un axe ou d'un plan factoriel.

## **b – Variables continues supplémentaires**

Il est possible de positionner des variables continues en élément supplémentaire (sans transformation au préalable en variable nominale par découpage en classes).

On calcule, comme dans l'analyse en composantes principales normée, le coefficient de corrélation de ces variables avec le facteur. Celui-ci fournit la coordonnée de la variable continue sur l'axe factoriel (cf. la schématisation de la figure 1.4 - 9). Les carrés des coefficients obtenus sont l'équivalent des cosinus carrés.

La position d'une variable sur un plan définit donc la direction où se situent les fortes valeurs de la variable. Ceci est d'autant plus vrai que la variable est proche du cercle de corrélations (de rayon 1) : il existe dans ce cas une liaison forte et linéaire entre la variable et les facteurs<sup>3</sup>.

---

<sup>1</sup> Les valeurs-tests permettent surtout de classer les modalités supplémentaires par ordre d'intérêt décroissant, ce qui constitue une aide précieuse à l'interprétation des facteurs.

<sup>2</sup> Les coordonnées sur un axe des individus correspondant à une modalité active ne peuvent être considérées comme tirées au hasard, puisque cette modalité aura contribué à construire l'axe.

<sup>3</sup> La lecture de la trajectoire des classes d'une variable continue transformée en variable nominale apporte souvent plus de précision que la seule position de la variable considérée comme continue (détection éventuelle de liaisons non linéaires).

### 1.4.5 Analyse du tableau de contingence de Burt : Equivalence avec l'analyse du tableau disjonctif complet

Le tableau  $\mathbf{B}$  de correspondance multiple, obtenu à partir d'un tableau disjonctif complet, est un assemblage particulier des tableaux de contingence qui sont les faces de l'hypercube de contingence.

L'analyse des correspondances appliquée à un tableau disjonctif complet  $\mathbf{Z}$  est équivalente à l'analyse du tableau de Burt  $\mathbf{B}$  et produit les mêmes facteurs.

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

Compte tenu de l'équation [1.4 - 2] donnant le  $\alpha^{\text{ième}}$  facteur  $\varphi_\alpha$  de l'analyse du tableau disjonctif complet  $\mathbf{Z}$ , la matrice à diagonaliser est :

$$\mathbf{S} = \frac{1}{s} \mathbf{D}^{-1} \mathbf{Z}' \mathbf{Z} = \frac{1}{s} \mathbf{D}^{-1} \mathbf{B}$$

Pour l'analyse du tableau de  $\mathbf{B}$  associé à  $\mathbf{Z}$ , le tableau des fréquences relatives  $\mathbf{F}$  s'écrit :

$$\mathbf{F} = \frac{1}{ns^2} \mathbf{B}$$

et

$$\mathbf{D}_p = \mathbf{D}_n = \frac{1}{ns} \mathbf{D}$$

On diagonalise la matrice :

$$\mathbf{S}^* = \frac{1}{s^2} \mathbf{D}^{-1} \mathbf{B} \mathbf{D}^{-1} \mathbf{B}$$

ce qui donne :

$$\mathbf{S}^* = \mathbf{S}^2$$

En prémultipliant les deux membres de [1.4 - 2] par  $\frac{1}{s} \mathbf{D}^{-1} \mathbf{B}$ , on obtient :

$$\frac{1}{s^2} \mathbf{D}^{-1} \mathbf{B} \mathbf{D}^{-1} \mathbf{B} \varphi_\alpha = \lambda_\alpha^2 \varphi_\alpha$$

Les facteurs des deux analyses sont donc colinéaires dans  $\mathbb{R}^p$  mais les valeurs propres associées diffèrent. Celles issues de l'analyse de  $\mathbf{B}$ , notées  $\lambda_B$ , sont le carré de celles issues de l'analyse de  $\mathbf{Z}$  :

$$\lambda_B = \lambda^2 \quad [1.4 - 5]$$

Les facteurs  $\varphi_\alpha$  issus de l'analyse de  $\mathbf{Z}$ , représentant les coordonnées factorielles des modalités, ont pour norme  $\lambda$ , alors que le facteur correspondant de l'analyse de  $\mathbf{B}$ , noté  $\varphi_{B\alpha}$ , aura pour norme  $\lambda^2$ .

D'où la relation liant les deux systèmes de coordonnées factorielles :

$$\varphi_{B\alpha} = \varphi_{\alpha} \sqrt{\lambda_{\alpha}} \quad [1.4 - 6]$$

### 1.4.6 Cas de deux questions

Dans le cas de deux questions  $q_1$  et  $q_2$ , le tableau disjonctif complet s'écrit :

$$Z = [Z_1, Z_2]$$

et nous ramène directement à l'analyse du tableau de contingence.

Il est alors équivalent, au point de vue de la description des associations entre modalités, d'effectuer :

- [1] l'analyse des correspondances du tableau  $Z$  d'ordre  $(n, p)$ ;
- [2] l'analyse des correspondances du tableau  $B$  d'ordre  $(p, p)$ ;
- [3] l'analyse des correspondances du tableau  $K = Z_1'Z_2$  d'ordre  $(p_1, p_2)$ .

L'équivalence entre l'analyse des correspondances du tableau disjonctif complet  $Z$  et celle du tableau des correspondances multiples  $B$  a été donnée dans le cas général de plusieurs questions.

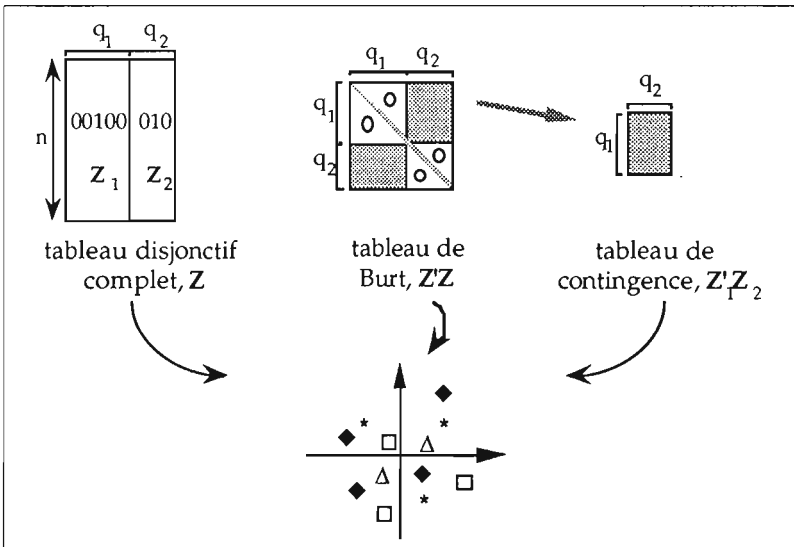


Figure 1.4 - 10  
Equivalence des trois analyses des correspondances

Intéressons-nous maintenant à l'équivalence entre l'analyse des correspondances du tableau disjonctif complet  $Z = [Z_1, Z_2]$  d'ordre  $(n, p)$  et celle du tableau de contingence  $K = Z_1'Z_2$  d'ordre  $(p_1, p_2)$  avec  $p = p_1 + p_2$ .

Montrons que, pour tout couple de facteurs  $(\psi_\alpha, \varphi_\alpha)$  relatifs à une même valeur propre  $\mu_\alpha$  issus de l'analyse du tableau de contingence  $Z_1'Z_2$ , il correspond un facteur  $\Phi_\alpha$  de l'analyse de  $Z$  (ou celle de  $B$ ), avec :

$$\Phi_\alpha = \begin{bmatrix} \psi_\alpha \\ \varphi_\alpha \end{bmatrix}$$

Rappelons que l'on note  $D_1 = Z_1'Z_1$  et  $D_2 = Z_2'Z_2$  et que :

$$D = \begin{bmatrix} D_1 & 0 \\ 0 & D_2 \end{bmatrix}$$

Les éléments diagonaux de  $D_1$  et  $D_2$  sont respectivement les marges en ligne et en colonne du tableau  $Z_1'Z_2$ .

L'analyse de ce tableau nous conduit aux relations de transition :

$$\begin{cases} \psi_\alpha = \frac{1}{\sqrt{\mu_\alpha}} D_1^{-1} Z_1'Z_2 \varphi_\alpha & [1.4 - 7] \\ \varphi_\alpha = \frac{1}{\sqrt{\mu_\alpha}} D_2^{-1} Z_2'Z_1 \psi_\alpha & [1.4 - 8] \end{cases}$$

On peut écrire ces relations sous la forme du système :

$$\begin{cases} D_1^{-1} (D_1 \psi_\alpha + Z_1'Z_2 \varphi_\alpha) = (1 + \sqrt{\mu_\alpha}) \psi_\alpha \\ D_2^{-1} (D_2 \varphi_\alpha + Z_2'Z_1 \psi_\alpha) = (1 + \sqrt{\mu_\alpha}) \varphi_\alpha \end{cases}$$

soit encore :

$$\begin{bmatrix} D_1 & 0 \\ 0 & D_2 \end{bmatrix}^{-1} \begin{bmatrix} D_1 & Z_1'Z_2 \\ Z_2'Z_1 & D_2 \end{bmatrix} \begin{bmatrix} \psi_\alpha \\ \varphi_\alpha \end{bmatrix} = (1 + \sqrt{\mu_\alpha}) \begin{bmatrix} \psi_\alpha \\ \varphi_\alpha \end{bmatrix}$$

Cette équation s'écrit de façon plus condensée :

$$D^{-1}Z'Z \Phi_\alpha = (1 + \sqrt{\mu_\alpha}) \Phi_\alpha \quad [1.4 - 9]$$

Après multiplication des deux membres par  $\frac{1}{s}$ , soit ici  $\frac{1}{2}$ , il vient :

$$\frac{1}{s} D^{-1}Z'Z \Phi_\alpha = \left(\frac{1 + \sqrt{\mu_\alpha}}{2}\right) \Phi_\alpha$$

On y reconnaît la relation [1.4 - 2] avec :

$$\lambda_\alpha = \frac{1 + \sqrt{\mu_\alpha}}{2}$$

Si  $\mu_\alpha$  est la  $\alpha^{\text{ième}}$  plus grande valeur propre issue de l'analyse du tableau de contingence  $Z_1'Z_2$ , alors  $\lambda_\alpha$  est la  $\alpha^{\text{ième}}$  plus grande valeur propre issue de l'analyse de  $Z$ .

Si par exemple  $p_1 \leq p_2$ , l'analyse de  $\mathbf{Z}$  conduit à :

- $p_1$  facteurs du type  $\begin{bmatrix} \psi_\alpha \\ \varphi_\alpha \end{bmatrix}$ , correspondant à la valeur propre  $\frac{1 + \sqrt{\mu_\alpha}}{2}$  ;
- $p_1$  facteurs du type  $\begin{bmatrix} \psi_\alpha \\ -\varphi_\alpha \end{bmatrix}$ , correspondant à la valeur propre  $\frac{1 - \sqrt{\mu_\alpha}}{2}$  ;
- $p_2 - p_1$  facteurs du type<sup>1</sup>  $\begin{bmatrix} 0 \\ \xi_\alpha \end{bmatrix}$ , correspondant à la valeur propre  $\frac{1}{2}$  .

Les résultats relatifs aux trois analyses équivalentes sont rassemblés dans le tableau 1.4 - 1.

**Tableau 1.4 - 1**  
Equivalence des analyses des trois tableaux  
dans le cas de deux questions

Tableau analysé	Dimension	Facteur	Valeur propre
$\mathbf{Z}'_1 \mathbf{Z}'_2$ tableau de contingence	$(p_1, p_2)$	$\psi$ dans $\mathbb{R}^{p_1}$ $\varphi$ dans $\mathbb{R}^{p_2}$	$\mu$
$\mathbf{Z} = [\mathbf{Z}_1, \mathbf{Z}_2]$ tableau disjonctif complet	$(p, n)$ où $p = p_1 + p_2$ .	$\Phi = \begin{bmatrix} \Psi \\ \varphi \end{bmatrix}$	$\lambda = \frac{1 + \sqrt{\mu}}{2}$
$\mathbf{B} = \mathbf{Z}' \mathbf{Z}$ Tableau de Burt	$(p, p)$	$\Phi_B = \Phi \sqrt{\lambda}$	$\lambda^2$

**Remarques :**

1) Les analyses de correspondances appliquées à ces trois types de tableaux, reposant sur la même information brute, donnent les mêmes axes factoriels, mais avec des valeurs propres différentes, donc des taux d'inertie différents. Les relations existant entre les taux d'inertie nous montrent que ceux-ci seront toujours plus élevés pour l'analyse du tableau de contingence  $\mathbf{Z}'_1 \mathbf{Z}'_2$  que pour l'analyse du tableau disjonctif complet  $\mathbf{Z}$ .

Ainsi, la somme des valeurs propres non triviales issues de l'analyse de  $\mathbf{Z}$  vaut :

$$\frac{p_1 + p_2}{2} - 1$$

Comme les valeurs propres sont inférieures ou égales à 1, aucun facteur ne peut avoir un taux d'inertie supérieur en pourcentage à :

$$\frac{2 \times 100}{p_1 + p_2 - 2}$$

Prenons l'exemple du tableau de contingence croisant les 8 professions et les 6 médias (cf. § 1.3.8). Le premier facteur prend en compte 50% de l'inertie totale. La remarque ci-dessus montre que l'analyse du tableau disjonctif correspondant ne

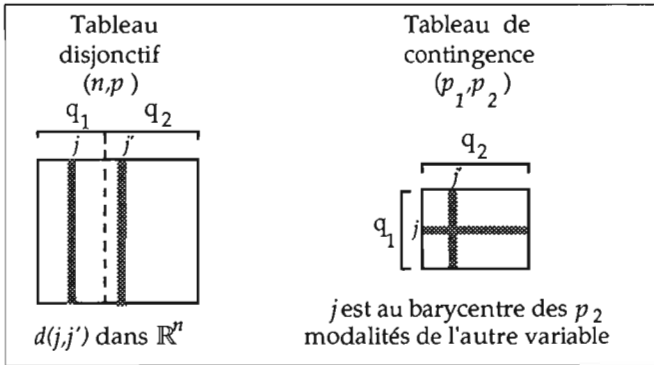
<sup>1</sup> Les axes  $\xi_\alpha$  complètent la base des  $\psi_\alpha$  dans  $\mathbb{R}^p$



peut pas donner un premier facteur expliquant plus de  $\frac{200}{8+6-2} = 16,6\%$ . Les taux d'inertie sont donc dépendants du codage préliminaire de l'information brute. Il faut donc éviter de les interpréter en termes "d'information". On reviendra sur ce point à la section 4.1.

2) Dans l'analyse du tableau disjonctif complet  $Z$ , les points représentant les diverses modalités de réponses aux deux questions sont les éléments d'un même ensemble, l'ensemble des colonnes de  $Z$ .

Au contraire dans l'analyse du tableau de contingence  $Z_1'Z_2$ , ils se scindent en points-lignes et en points-colonnes (cf. figure 1.4 - 11).



**Figure 1.4 - 11**  
Proximité entre deux modalités de variables différentes

Le fait que les représentations obtenues dans l'espace des premiers facteurs soient identiques (à une dilatation près, due au fait que les valeurs propres ne sont pas les mêmes) montre que la représentation simultanée des points-lignes et des points-colonnes en analyse des correspondances n'est pas un simple artifice graphique.

L'interprétation de la position de deux modalités relatives à deux variables différentes dépend du tableau d'analyse. Dans le tableau disjonctif complet, cette position s'interprète en terme de distance. Dans le tableau de contingence, la distance entre une ligne et une colonne n'a pas de sens et une modalité est au "quasi-barycentre" des modalités de l'autre variable. L'analyse de ces deux tableaux fournit des représentations similaires.

### 1.4.7 Cas particuliers

Dans le cas où toutes les variables ont deux modalités, l'analyse des correspondances multiples se ramène à l'analyse en composantes principales des variables caractérisées par une seule de leurs modalités. Dans le cas où l'ensemble des questions peut être partitionné en deux groupes à l'intérieur desquels les questions sont indépendantes, l'analyse des correspondances multiples se ramène à l'analyse de la correspondance entre

les deux groupes : juxtaposition de tables de contingences constituant un sous-tableau du tableau de Burt.

### a – Toutes les questions ont deux modalités

Les variables n'étant représentées que par une seule de leurs modalités  $p - s = \frac{p}{2}$ , on obtient directement la matrice à diagonaliser qui n'est autre que la matrice des corrélations entre variables (Nakhlé, 1976). Rappelons que d'après [1.4 - 2] :

$$\frac{1}{s} \mathbf{D}^{-1} \mathbf{B} \Phi = \lambda \Phi \quad [1.4 - 10]$$

Explicitons cette relation où  $\mathbf{D}$  désigne la matrice diagonale ayant les mêmes éléments diagonaux que  $\mathbf{B}$  et où  $l$  et  $j$  désignent deux modalités :

$$\frac{1}{s} \sum_{j \in p} \frac{b_{lj}}{b_{ll}} \Phi_j = \lambda \Phi_l \quad [1.4 - 11]$$

L'ensemble  $p$  des  $p$  modalités est partitionné en deux sous-ensembles  $p^1$  et  $p^2$  formés respectivement des premières et des deuxièmes modalités de chacune des  $s$  questions :

$$p = p^1 \cup p^2$$

Pour tout  $q \in s$  :

$$p_q = \{j_q^1, j_q^2\}$$

avec  $j_q^1 \in p^1$  et  $j_q^2 \in p^2$ . Notons les relations, pour tout  $q \in s$  :

$$b_{ij_q^1} + b_{ij_q^2} = b_{il} \quad \text{pour tout } l \in p$$

Cette relation exprime que ceux qui ont choisi la réponse  $l$  et l'une ou l'autre des deux modalités de la question  $J_q$  sont simplement ceux qui ont choisi la réponse  $l$ .

$$b_{j_q^1 j_q^1} + b_{j_q^2 j_q^2} = n \quad \text{et} \quad b_{j_q^1 j_q^1} \Phi_{j_q^1} = -b_{j_q^2 j_q^2} \Phi_{j_q^2}$$

La première relation exprime que tous les individus doivent choisir au moins une modalité de réponse pour chaque question, et la seconde traduit le fait que les coordonnées sont centrées pour chaque question.

Il suffit donc de restreindre la sommation de la relation [1.4 - 11] au seul ensemble  $p^1$ , dont l'élément courant sera désormais noté  $j$  :

$$\frac{1}{s b_{ll}} \sum_{j \in p^1} \left( b_{lj} - \frac{(b_{ll} - b_{jj}) b_{jj}}{n - b_{jj}} \right) \Phi_j = \lambda \Phi_l$$

Ce qui peut s'écrire :

$$\sum_{j \in p^1} \frac{n b_{lj} - b_{ll} b_{jj}}{s (n - b_{ll}) b_{ll}} \Phi_j = \lambda \Phi_l \quad [1.4 - 12]$$

Calculons les moments empiriques centrés du second ordre des  $s$  variables caractérisées par leurs premières modalités :

$$\text{Cov}(l, j) = \frac{1}{n}(b_{lj} - \frac{b_{ll}b_{jj}}{n})$$

$$\text{Var}(j) = \frac{1}{n}(b_{jj} - \frac{b_{jj}^2}{n})$$

Le terme général de la matrice des corrélations des  $s$  variables s'écrit :

$$\text{Cor}(l, j) = \frac{n b_{lj} - b_{ll} b_{jj}}{\sqrt{(n - b_{jj}) b_{jj} (n - b_{ll}) b_{ll}}}$$

Il est clair que si  $(\Phi, \lambda)$  est la solution de l'équation [1.4 - 12] alors  $(\Phi^*, \lambda^*)$  est la solution de :

$$\sum_{j \in P^1} \text{Cor}(l, j) \Phi_j^* = \lambda^* \Phi_j^*$$

avec :

$$\Phi_j = \Phi_j^* \frac{\sqrt{n - b_{jj}}}{\sqrt{b_{jj}}}$$

et :

$$\lambda^* = \lambda s$$

Les facteurs et les valeurs propres d'une analyse des correspondances multiples de  $s$  variables à deux modalités ( $p = 2s$ ) sont bien reliés par une relation simple à ceux d'une analyse en composantes principales normées effectuées sur les premières (ou les secondes) modalités de chacune des  $s$  questions (sélection de  $s$  colonnes du tableau disjonctif complet).

## b – Sous-tableau d'un tableau de correspondances multiples

Lorsque l'ensemble des  $s$  questions est partitionné en au moins deux sous-ensembles  $s_1$  et  $s_2$  totalisant respectivement  $p_1$  et  $p_2$  modalités (avec  $p_1 + p_2 = p$ ), on peut vouloir analyser le sous-tableau  $\mathbf{B}_{12}$  croisant ces deux sous-ensembles obtenu à partir du tableau de correspondances multiples.

### - Analyse du sous-tableau

L'analyse du tableau des correspondances multiples  $\mathbf{B}$  permet d'étudier les liaisons entre toutes les questions.

L'analyse du sous-tableau  $\mathbf{B}_{12}$  permet d'étudier les relations existant entre les éléments de  $s_1$  et ceux de  $s_2$  sans tenir compte des dépendances internes à  $s_1$ , ni des dépendances internes à  $s_2$ . Le groupe de questions  $s_1$  est caractérisé par ses associations avec les questions de  $s_2$  et réciproquement (cf. Leclerc, 1975).

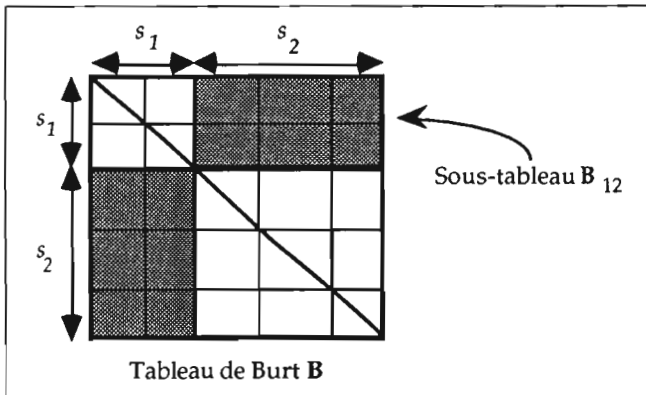


Figure 1.4 - 12  
Sous-tableau  $B_{12}$  du tableau de contingence de Burt  $B$

Lorsqu'un des groupes est réduit à une seule question  $q_0$ , le tableau de données est une bande du tableau des correspondances multiples croisant la variable  $q_0$  avec un groupe de variables ne contenant pas  $q_0$ .

C'est aussi le tableau des barycentres des groupes d'individus définis par les modalités de  $q_0$ .

Nous verrons (§ 3.3.8.b) que l'analyse d'une bande d'un tableau de correspondances multiples constitue une méthode de discrimination appelée analyse discriminante barycentrique.

Les résultats obtenus par l'analyse des correspondances du tableau de Burt  $B$  et celle de la tranche  $B_{12}$  sont en général différents (les nuages relatifs à ces tableaux ne sont pas dans le même espace). Ce sont les objectifs de l'étude qui doivent guider le choix du tableau à analyser.

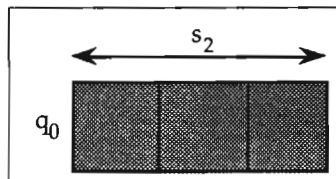


Figure 1.4 - 13  
Bande du tableau  
des correspondances multiples

Cependant, si les variables de chaque sous-ensemble sont indépendantes entre elles, les analyses réalisées à partir des tableau  $B$  et  $B_{12}$  sont équivalentes et celles de chaque sous-ensemble  $s_1$  et  $s_2$  ne présentent pas d'intérêt.

- Cas où l'analyse multiple se ramène à une correspondance binaire

Le cas d'une correspondance binaire s'est révélé particulièrement intéressant du point de vue des calculs à mettre en œuvre. En effet, l'analyse du tableau des correspondances multiples d'ordre  $(p, p)$  est équivalente à l'analyse des correspondances du tableau de contingence croisant les modalités des deux questions, ce qui conduit à diagonaliser une matrice dont l'ordre est déterminé par le plus petit des nombres  $p_1$  et  $p_2$ .

Nous retiendrons la propriété suivante. Si à l'intérieur des deux sous-ensembles  $s_1$  et  $s_2$  les questions sont indépendantes, l'analyse des  $s$  questions se ramène à celle d'une correspondance binaire, et donc à la diagonalisation d'une matrice d'ordre  $\text{Inf}(p_1, p_2)$ .

Nous dirons ici que deux questions  $q$  et  $q'$  sont indépendantes si la table de contingence correspondante vérifie la relation<sup>1</sup> :

$$Z'_q Z_{q'} = \frac{1}{n} d_q d_{q'}$$

où les vecteurs  $d_q$  et  $d_{q'}$  ont respectivement pour composantes les éléments diagonaux de  $Z'_q Z_q$  et  $Z'_{q'} Z_{q'}$  (c'est-à-dire les éléments diagonaux de  $D_q$  et  $D_{q'}$  par définition de ces matrices).

Ecrivons de nouveau la relation [1.4 - 10] en partitionnant  $\Phi$  en deux blocs  $\Phi_{s_1}$  et  $\Phi_{s_2}$ ; on découpe également les matrices  $B$  et  $D$  en quatre blocs, de façon à faire apparaître la partition  $s = s_1 \cup s_2$  :

$$B = \begin{bmatrix} B_{11} & B_{12} \\ B_{21} & B_{22} \end{bmatrix} \quad D = \begin{bmatrix} D_1 & 0 \\ 0 & D_2 \end{bmatrix}$$

On obtient les deux relations :

$$\begin{cases} \frac{1}{s} (D_1^{-1} B_{11} \Phi_{s_1} + D_1^{-1} B_{12} \Phi_{s_2}) = \lambda \Phi_{s_1} \\ \frac{1}{s} (D_2^{-1} B_{21} \Phi_{s_1} + D_2^{-1} B_{22} \Phi_{s_2}) = \lambda \Phi_{s_2} \end{cases}$$

Remarquons que les  $s_1$  (respectivement  $s_2$ ) blocs diagonaux de  $D_1^{-1} B_{11}$  (respectivement  $D_2^{-1} B_{22}$ ) sont des matrices unité dont les ordres correspondent aux cardinaux de chacune des questions.

On a d'autre part, pour  $k \in \{1, 2\}$  :

$$q \in s_k; q' \in s_k; q \neq q' \Rightarrow D_q^{-1} Z'_q Z_{q'} = \frac{1}{n} D_q^{-1} d_q d_{q'}$$

En désignant par  $e_q$  un vecteur dont les  $q$  composantes valent 1 :

<sup>1</sup> Bien entendu, l'indépendance théorique entre les deux questions n'implique pas que cette relation soit exactement vérifiée sur l'échantillon.

$$\mathbf{D}_q^{-1} \mathbf{Z}'_q \mathbf{Z}_q = \frac{1}{n} \mathbf{e}_q \mathbf{d}'_q$$

Les relations  $\mathbf{d}'_q \Phi_q = 0$  (centrage des modalités relatives à chaque question) impliquent finalement :

$$\mathbf{D}_1^{-1} \mathbf{B}_{11} \Phi_{s_1} = \Phi_{s_1} \quad \text{et} \quad \mathbf{D}_2^{-1} \mathbf{B}_{22} \Phi_{s_2} = \Phi_{s_2}$$

Le système ci-dessus s'écrit alors :

$$\begin{cases} \mathbf{D}_1^{-1} \mathbf{B}_{12} \Phi_{s_2} = (\lambda s - 1) \Phi_{s_1} \\ \mathbf{D}_2^{-1} \mathbf{B}_{21} \Phi_{s_1} = (\lambda s - 1) \Phi_{s_2} \end{cases}$$

D'où par substitution :

$$\mathbf{D}_2^{-1} \mathbf{B}_{21} \mathbf{D}_1^{-1} \mathbf{B}_{12} \Phi_{s_2} = (\lambda s - 1)^2 \Phi_{s_2}$$

Ainsi  $\Phi_{s_2}$  est obtenu par diagonalisation d'une matrice d'ordre  $(s_1, s_2)$ . On en déduit facilement  $\Phi_{s_1}$ .

Remarquons que  $\mathbf{B}_{12}$  est obtenu par juxtaposition des tableaux de contingence croisant l'ensemble des modalités des questions du premier groupe avec celles relatives au second groupe. Les marges du tableau  $\mathbf{B}_{12}$  sont les éléments diagonaux de  $s_2 \mathbf{B}_1$  et  $s_1 \mathbf{B}_2$ .

Les facteurs issus de l'analyse des correspondances directe du tableau  $\mathbf{B}_{12}$  considéré comme un tableau de contingence vérifient la relation :

$$\frac{1}{s_1 s_2} \mathbf{D}_2^{-1} \mathbf{B}_{21} \mathbf{D}_1^{-1} \mathbf{B}_{12} \Psi = \lambda \Psi$$

Ils sont donc proportionnels aux facteurs trouvés précédemment<sup>1</sup>.

### 1.4.8 Exemple d'application numérique

L'exemple qui va suivre concerne un petit sous-échantillon (105 individus, 9 questions) de l'enquête "Conditions de vie et aspirations des Français"<sup>2</sup>.

Le tableau 1.4 - 2 est le tableau de données proprement dit, en codage condensé (cf. section 1.4.2 ci-dessus), à l'exception de la variable V2 (âge) qui est numérique.

<sup>1</sup> Ces propriétés concernant les sous-tableaux de tableaux de Burt ont été étudiées par A. Leclerc (1975), puis généralisées par P. Cazes (cf. Cazes, 1977, 1980, 1981).

<sup>2</sup> Pour une présentation générale, des références relatives à cette enquête et des exemples d'application en vraie grandeur, cf. § 2.4.4.

**Tableau 1.4 - 2**  
**Tableau de données R en codage condensé**

n°	V1	V2	V3	V4	V5	V6	V7	V8	V9	n°	V1	V2	V3	V4	V5	V6	V7	V8	V9
1	2	27	1	2	2	1	1	2	1	54	2	54	1	2	2	2	1	1	1
2	2	42	1	3	2	2	1	1	3	55	2	48	1	3	2	2	1	1	1
3	1	71	1	2	2	2	2	1	1	56	2	30	1	3	2	1	1	1	1
4	1	52	1	2	2	1	1	2	1	57	2	50	1	2	2	2	2	1	1
5	2	36	1	2	2	2	2	1	3	58	1	21	1	2	1	2	2	2	2
6	1	22	2	2	2	2	1	2	3	59	2	47	1	2	2	2	2	1	2
7	1	26	2	2	2	2	2	2	2	60	1	51	2	2	2	2	1	1	2
8	2	43	1	2	1	1	2	2	1	61	1	27	2	3	2	1	2	1	2
9	2	33	1	2	2	2	2	1	1	62	2	37	2	3	2	1	2	1	2
10	2	54	2	4	2	2	1	1	3	63	1	67	1	4	2	2	1	1	1
11	1	57	1	3	2	1	1	2	3	64	2	30	2	3	2	2	1	1	3
12	1	33	1	2	2	1	1	1	1	65	1	40	2	2	2	2	2	1	2
13	1	65	1	2	2	2	1	2	1	66	1	67	1	4	2	2	2	1	1
14	2	58	1	2	2	2	2	2	2	67	2	51	1	2	2	2	2	2	1
15	2	33	1	3	2	1	2	1	1	68	1	35	2	2	2	1	1	1	2
16	1	37	1	4	2	2	2	1	1	69	1	24	1	2	2	2	1	1	3
17	1	46	1	3	2	2	1	1	2	70	1	34	2	4	2	2	2	1	4
18	2	30	2	3	2	1	2	1	1	71	1	55	1	4	2	1	1	2	2
19	1	64	1	4	2	2	1	2	1	72	1	41	1	2	2	1	1	1	1
20	2	64	1	2	2	2	2	2	1	73	1	32	1	2	2	1	2	1	2
21	1	41	1	3	2	2	2	2	1	74	1	35	1	2	1	2	1	1	1
22	1	56	1	3	2	2	1	1	1	75	2	27	2	2	2	2	2	1	2
23	2	21	1	3	1	1	1	2	1	76	2	22	2	2	2	2	1	1	2
24	2	49	1	2	2	1	1	1	1	77	2	31	2	2	2	2	1	1	1
25	1	60	2	2	1	2	2	2	2	78	1	35	2	3	2	2	2	1	3
26	1	63	1	1	2	1	1	2	1	79	2	33	2	2	1	1	1	2	1
27	2	46	1	3	1	1	2	1	1	80	1	39	1	2	2	1	2	1	4
28	2	53	2	2	2	2	2	1	3	81	2	21	1	2	2	2	2	2	3
29	2	29	2	3	2	2	1	1	2	82	2	51	1	3	2	2	2	1	3
30	1	59	1	3	2	2	2	1	1	83	2	35	1	3	2	2	1	1	1
31	2	48	1	3	2	2	2	1	3	84	1	58	2	2	2	2	2	1	3
32	2	19	2	2	2	2	2	1	3	85	1	54	1	1	2	1	1	2	3
33	1	56	1	2	2	2	2	2	1	86	2	21	2	3	2	1	2	2	2
34	1	30	1	4	1	2	2	1	3	87	1	29	2	2	2	2	2	1	1
35	2	66	2	3	2	1	1	1	1	88	2	32	1	1	2	2	1	1	3
36	2	30	1	3	2	1	1	1	1	89	2	40	1	2	1	1	2	2	1
37	2	39	1	3	1	1	1	1	1	90	1	34	2	2	1	2	2	2	2
38	1	52	1	2	2	2	2	1	3	91	1	33	2	3	2	2	2	1	2
39	1	23	1	1	2	1	1	1	1	92	2	82	1	1	2	2	1	2	1
40	1	52	1	1	1	2	2	2	3	93	1	69	1	3	2	2	1	2	1
41	1	47	1	1	1	1	2	1	1	94	2	38	2	2	2	2	1	1	3
42	1	47	1	3	2	2	1	1	1	95	1	80	1	3	2	1	1	1	1
43	2	71	2	2	1	2	2	2	1	96	2	39	1	2	1	1	1	1	4
44	2	64	1	2	2	1	1	2	1	97	2	61	1	1	1	2	1	2	1
45	1	37	2	2	1	2	2	2	3	98	1	67	1	2	2	2	1	2	3
46	2	62	1	2	2	2	2	2	1	99	1	24	1	2	1	2	2	2	2
47	1	45	2	1	2	2	2	1	2	100	2	43	1	3	2	2	2	1	1
48	1	26	2	2	2	2	1	2	2	101	1	54	2	1	2	2	2	2	2
49	2	40	1	2	1	1	1	1	1	102	2	76	1	2	2	2	2	2	1
50	1	23	1	3	2	2	2	1	2	103	2	45	1	1	1	1	1	2	2
51	2	28	1	2	1	2	2	1	2	104	2	24	1	2	1	2	2	2	1
52	1	40	2	2	1	2	2	2	2	105	2	80	1	2	2	2	1	2	1
53	1	40	1	2	2	2	1	1	1										

Les libellés des questions figurent dans le tableau 1.4 - 3, les libellés des modalités correspondantes se retrouveront dans les listages de résultats plus bas. Les libellés abrégés en 4 caractères seront utilisés pour les représentations graphiques. Les 4 variables actives servent à calculer les

distances et les axes, les 4 variables illustratives et la variable continue illustrative servent à interpréter *a posteriori* les axes et les proximités.

**Tableau 1.4 - 3**  
**Description des libellés des 9 questions**

<b>4 questions actives</b>	<b>13 modalités associées</b>
-V3- La famille est le seul endroit où l'on se sent bien (2 modalités)	FA01 = oui, FA02 = non.
-V4- Les dépenses de logement sont pour vous une charge (4 modalités)	DL01 = négligeable, DL02 = sans gros problème, DL03 = une lourde charge, DL04 = Une très lourde charge.
-V7- Avez-vous souffert récemment de mal au dos (2 modalités)	MA01 = oui, MA02 = non.
-V8- Vous imposez-vous régulièrement des restrictions (2 modalités)	RE01 = oui, RE02 = non.
<b>4 questions illustratives</b>	<b>10 modalités associées</b>
-V1- Sexe de l'enquêté(e) (2 modalités)	MASC = masculin, FEMI = féminin.
-V5 Disposez-vous d'un magnétophone (2 modalités)	MAG1 = oui, MAG2 = non.
-V6- Avez-vous souffert récemment de maux de tête (2 modalités)	MT01 = oui, MT02 = non.
-V9- Regardez-vous la télévision ? (4 modalités)	TV01 = tous les jours, TV02 = assez souvent, TV03 = pas très souvent, TV04 = jamais.
<b>1 variable continue illustratives</b>	
-V2- Age de l'enquêté(e) (continue)	

Les tableaux disjonctifs complets correspondant aux variables nominales ne sont pas présentés et ne sont jamais développés tels quels dans les calculs. Le tableau de Burt (tableau 1.4 - 4) est calculé directement à partir du codage condensé<sup>1</sup>. Le tableau 1.4 - 4 ne représente que la moitié inférieure du tableau de Burt relatif aux 4 questions actives. On trouve dans ce tableau les 6 tableaux de contingence croisant les 4 questions actives deux à deux. Sur la diagonale se trouvent les questions croisées avec elles-mêmes, et donc les effectifs correspondant à chaque modalité.

On vérifie ensuite (tableau 1.4 - 5) qu'il y a 6 valeurs propres non nulles ( $6 = p - s$ ), et on peut constater que les taux d'inertie correspondant à chaque valeur propre sont modestes, malgré la petite taille de cet exemple pédagogique. Il s'agit là d'une propriété propre à cette méthode : les taux d'inertie sont toujours des mesures très pessimistes de l'information extraite, car le codage disjonctif induit une orthogonalité artificielle des colonnes du tableau. Plusieurs indicateurs de remplacement ont été proposés.

<sup>1</sup> Cette procédure divise le nombre d'opérations par le coefficient  $(s/p)^2$ ,  $s$  étant le nombre de questions actives et  $p$  le nombre total de modalités correspondantes. Dans le cas d'applications courantes ( $p > 100$ ,  $n > 1000$ ,  $n$  étant le nombre d'individus) ce gain est très appréciable.



**Tableau 1.4 - 4**  
**Tableau de Burt des s = 4 questions actives**

	FA01	FA02	DL01	DL02	DL03	DL04	MA01	MA02	RE01	RE02
FA01	72	0								
FA02	0	33								
DL01	9	2	11	0	0	0				
DL02	37	20	0	57	0	0				
DL03	21	9	0	0	30	0				
DL04	5	2	0	0	0	7				
MA01	38	12	7	24	16	3	50	0		
MA02	34	21	4	33	14	4	0	55		
RE01	42	22	4	29	25	6	31	33	64	0
RE02	30	11	7	28	5	1	19	22	0	41

On peut considérer les carrés des valeurs propres, qui sont les valeurs propres de l'analyse des correspondances du tableau de Burt considéré comme tableau de données (cf. § 1.4.5) et qui fournissent des taux d'inertie un peu moins pessimistes. On peut également prendre en compte des fonctions particulières des valeurs propres comme mesures de l'inertie (Benzécri, 1979)<sup>1</sup>.

**Tableau 1.4 - 5**  
**Valeurs propres et taux d'inertie**

NUMERO	VALEUR PROPRE	POURCENT.	POURCENT. CUMULE	
1	.3416	22.77	22.77	*****
2	.3175	21.17	43.94	*****
3	.2520	16.80	60.74	*****
4	.2232	14.88	75.62	*****
5	.2075	13.84	89.46	*****
6	.1582	10.54	100.00	*****
Total	1.5000	100.00		

Le tableau 1.4 - 6 fournit les indicateurs nécessaires pour interpréter les positions des modalités actives.

Les règles de lecture sont semblables à celles du tableau 1.3 - 13 relatif à l'analyse des correspondances simple. Seuls les calculs de contributions cumulées pour les modalités de chaque question ont été ajoutés. Leur interprétation est immédiate. Il est clair, par exemple, que les deux questions relatives aux dépenses de logement et aux restrictions définissent entièrement le premier axe.

<sup>1</sup> Benzécri a proposé la quantité  $\rho(\lambda) = \left(\frac{s}{s-1}\right)^2 \left(\lambda - \frac{1}{s}\right)^2$  qui est voisine de  $\lambda^2$  si le nombre de questions  $s$  est grand, et qui correspond, dans le cas  $s = 2$ , à la valeur propre  $\mu$  de l'analyse des correspondances de la table de contingence croisant les deux questions [dans ce cas, en effet,  $\rho(\lambda) = \mu = (2\lambda - 1)^2$ ]. (voir aussi § 4.1.5.a).

Tableau 1.4.6

Coordonnées, contributions et cosinus carrés des modalités actives sur les axes 1 à 3

MODALITES			COORDONNEES			CONTRIBUTIONS			COSINUS CARRES		
IDEN - LIBELLE	P.REL	DISTO	1	2	3	1	2	3	1	2	3
<i>- la famille est le seul endroit ou l'on se sent bien</i>											
FA01 - - oui -	17.14	.46	.14	-.42	.12	1.0	9.3	.9	.05	.38	.03
FA02 - - non -	7.86	2.18	-.31	.91	-.26	2.3	20.4	2.1	.05	.38	.03
-----+-----						CUMUL =	3.3	29.7	3.0	+-----	
<i>- les dépenses de logement sont pour vous une charge</i>											
DL01 - négligeable	2.62	8.55	1.32	-1.32	.33	13.4	14.4	1.2	.20	.20	.01
DL02 - sans gros problème	13.57	.84	.41	.52	-.11	6.7	11.8	.6	.20	.33	.01
DL03 - une lourde charge	7.14	2.50	-1.00	-.50	-.72	21.1	5.7	14.8	.40	.10	.21
DL04 - très lourde charge	1.67	14.00	-1.11	-.05	3.45	6.0	.0	78.7	.09	.00	.85
-----+-----						CUMUL =	47.2	31.9	95.2	+-----	
<i>- avez-vous souffert récemment de mal au dos</i>											
MA01 - - oui -	11.90	1.10	.03	-.73	-.14	.0	19.8	.9	.00	.48	.02
MA02 - - non -	13.10	.91	-.02	.66	.13	.0	18.0	.8	.00	.48	.02
-----+-----						CUMUL =	.0	37.9	1.86	+-----	
<i>- vous imposez-vous régulièrement des restrictions</i>											
RE01 - - oui -	15.24	.64	-.66	-.06	.01	19.3	.2	.0	.68	.01	.00
RE02 - - non -	9.76	1.56	1.03	.10	-.01	30.2	.3	.0	.68	.01	.00
-----+-----						CUMUL =	49.5	.5	.0	+-----	

Le tableau 1.4 - 7 donne les valeurs-test (cf. section 1.4.4.a ci-dessus) et les coordonnées des modalités supplémentaires sur les trois premiers axes. On note que les seules coordonnées significatives sur le premier axe sont relatives à la possession d'un magnéto-scope (valeurs-test de 2.8). Les mentions de maux de têtes et l'écoute de la télévision - toutes deux liées à l'âge - sont caractéristiques du deuxième axe.

Le tableau 1.4 - 8 est relatif à la variable continue "âge". On y lit sa moyenne, son écart-type, et ses coefficients de corrélation avec les trois premiers axes.

La structure du nuage des modalités actives est décrite par le plan factoriel de la figure 1.4 - 5, qui résume donc les 6 tables de contingence.

Le petit nombre de questions et le faible nombre d'individus limitent l'intérêt des résultats, mais permettent en revanche de comprendre le mécanisme de la méthode. Les deux questions les plus liées (*dépenses de logements* et *restrictions*) emportent le premier axe, la question relative aux dépenses de logement intervenant avec un poids double compte tenu du nombre de ses modalités (cf. § 1.4 .3-g). Les deux questions restantes, plus faiblement liées, caractérisant le deuxième axe.

La représentation simultanée des lignes et des colonnes liée à l'analyse des correspondances n'est pas utilisée sur la figure 1.4 - 5. Les 105 points-lignes correspondent à des individus anonymes ; seules leurs caractéristiques présentent de l'intérêt. Les individus n'interviennent donc que par le truchement des variables supplémentaires.

Les positions des modalités supplémentaires doivent être tempérées par leurs valeur-tests. Dans les études en vraie grandeur où ces modalités peuvent être très nombreuses, seules celles ayant des valeurs-test significatives sont portées sur les graphiques. Ainsi, la variable *sexe* (valeurs-test 0.5 et 0.4 sur les axes 1 et 2) pourrait ne pas figurer dans ce plan factoriel. De même, la modalité *TVO4*, (*ne regarde jamais la télévision*) malgré sa position relativement excentrée à gauche, n'est pas non plus significative (valeur-test = -1.0) car elle ne concerne que 3 individus.

Remarquons que la seule phase du processus permettant de procéder à une inférence statistique est précisément le calcul des valeurs-test relatives aux modalités supplémentaires. Malgré la taille modeste de l'échantillon et le petit nombre de variables, on peut rejeter l'hypothèse d'indépendance entre la possession d'un magnéto-scope (point *MAC2*) et l'aisance financière telle qu'elle est décrite par les modalités (*DL01*, *DL02*, *RE02*).

La variable continue *AGE* est représentée comme un axe, en pointillé. Cette direction a une certaine cohérence, malgré la faible taille de l'échantillon (les individus plus âgés ont des idées plus traditionalistes sur la famille, sont plus souvent propriétaires de leur logements, plus fréquemment téléspectateurs).

Tableau 1.4.7

Coordonnées et valeurs-test des modalités illustratives sur les axes 1 à 3.

MODALITES			VALEURS-TEST			COORDONNEES			
IDEN - LIBELLE	EFF.	P.ABS	1	2	3	1	2	3	DISTO.
<i>- sexe de l'enquêté(e)</i>									
MASC - masculin	53	53.00	.5	.4	2.1	.05	.04	.21	.98
FEMI - feminin	52	52.00	-.5	-.4	-2.1	-.05	-.04	-.21	1.02
<i>- disposez-vous d'un magnétoscope</i>									
MAG1 - oui -	22	22.00	2.8	.7	.5	.54	.13	.09	3.77
MAG2 - non -	83	83.00	-2.8	-.7	-.5	-.14	-.03	-.02	.27
<i>- avez-vous souffert récemment de maux de tete</i>									
MT01 - oui -	33	33.00	.0	-3.1	-1.3	.01	-.45	-.19	2.18
MT02 - non -	72	72.00	.0	3.1	1.3	.00	.21	.09	.46
<i>- regardez-vous la télévision ?</i>									
TV01 - tous les jours	53	53.00	.7	-3.4	-.2	.07	-.33	-.02	.98
TV02 - assez souvent	27	27.00	.1	3.3	-.9	.02	.56	-.16	2.89
TV03 - pas très souvent	22	22.00	-.6	.3	.4	-.11	.07	.08	3.77
TV04 - jamais	3	3.00	-1.0	.7	1.9	-.56	.39	1.11	34.00

Tableau 1.4.8

Coordonnées (corrélations) de la variable continue illustrative sur les axes 1 à 3.

VARIABLE CONTINUE		CARACTERISTIQUES			CORRELATIONS		
(IDEN)	LIBELLE COURT	EFFECTIF	MOYENNE	EC.TYPE	1	2	3
- (age )	age de l'enquete(e)	105	43.89	15.50	.23	-.23	.15

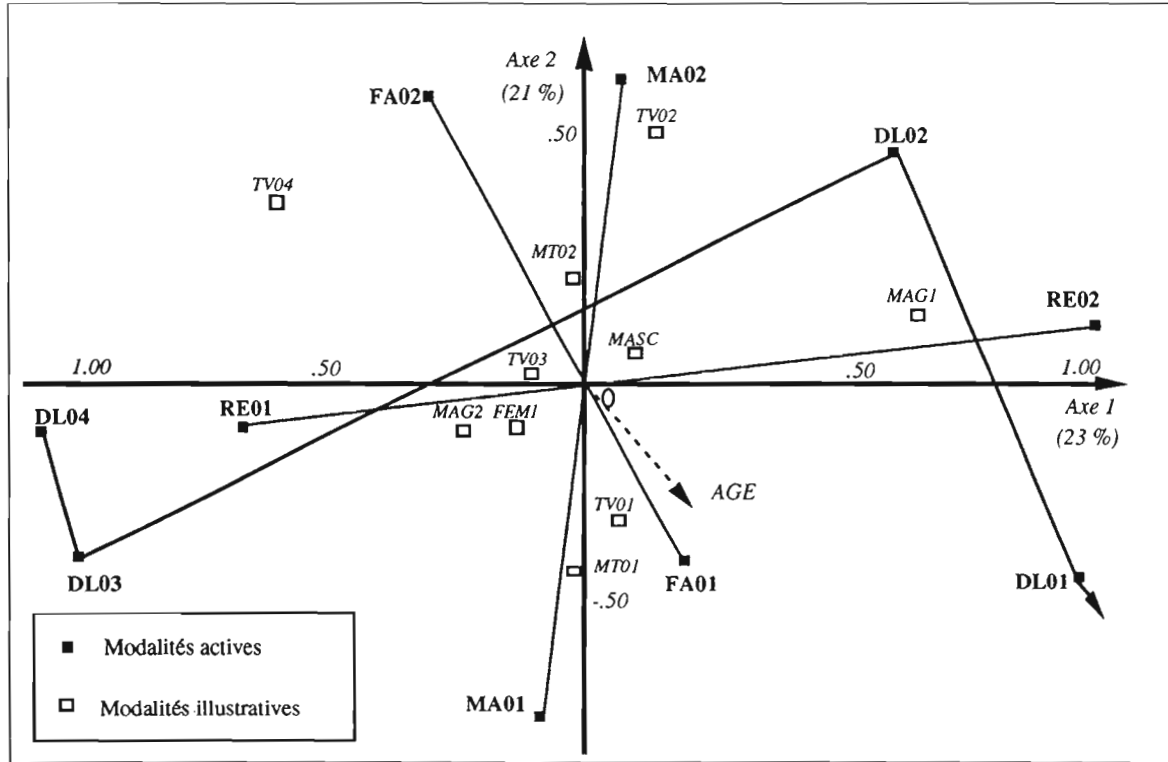


Figure 1.4.5  
Position des modalités actives et illustratives sur le premier plan factoriel.

Les modalités "consécutives" des questions actives sont jointes par des lignes polygonales. On vérifie que l'origine est bien un centre de gravité pour les modalités de chaque question, ce qui implique un alignement avec l'origine pour les questions à 2 modalités. Les variables "restrictions" (RE02 = ne s'impose pas de restriction) et "dépense de logement" (DL01 = négligeables, DL02 = sans gros problème) déterminent le premier axe, illustré a posteriori par la position du point MAG2 (possession d'un magnétoscope). La variable continue AGE est repérée par ses coefficients de corrélation avec les axes (flèche en pointillé).

Chapitre 2

---

**QUELQUES MÉTHODES  
DE CLASSIFICATION**



## Introduction

Les techniques de classification automatique sont destinées à produire des groupements de lignes ou de colonnes d'un tableau. Il s'agit le plus souvent d'objets ou d'individus décrits par un certain nombre de variables ou de caractères. La classification est une branche de l'analyse des données qui a donné lieu à des publications nombreuses et diversifiées. Les ouvrages spécialisés (notamment, en langue française, le tome 1 du traité d'analyse des données de Benzécri, 1973) contiennent en général d'importantes considérations historiques et de rigoureux développements formels sur la notion de classification. L'ouvrage de base, historique, est celui de Sokal et Sneath (1963). Les premiers manuels publiés furent ceux de Lerman (1970), Anderberg (1973), Benzécri (1973), Hartigan (1975), Lerman (1981) et Gordon (1981) auxquels nous ne pouvons que renvoyer le lecteur pour des préalables fondamentaux<sup>1</sup>. Nous nous bornerons ici aux principes de base des méthodes les plus largement utilisées.

Les circonstances d'utilisation sont sensiblement les mêmes que celles des méthodes d'analyse factorielle descriptive présentées au chapitre 1 : l'utilisateur se trouve face à un tableau rectangulaire de valeurs numériques. Ce tableau peut être un tableau de valeurs numériques continues (valeur de la variable  $j$  pour l'individu  $i$ , à l'intersection de la ligne  $i$  et de la colonne  $j$  du tableau), un tableau de contingence (croisant deux partitions d'une même population), ou encore un tableau de présence-absence (valeurs 0 ou 1 selon que tel individu ou objet possède tel caractère ou attribut). Dans certaines applications, l'utilisateur peut disposer d'un tableau carré symétrique de similarités ou de distances.

Le recours aux techniques de classification automatique est sous-tendu par quelques idées générales concernant le champ d'observation. On suppose que certains regroupements doivent exister, ou au contraire on exige que certains regroupements soient effectués. Autrement dit, on ne se satisfait pas d'une visualisation plane et continue des associations statistiques et l'on manifeste, implicitement ou explicitement, un intérêt pour la mise en évidence de *classes* d'individus ou de caractères. Les représentations synthétiques se manifestent soit sous la forme de *partitions* des ensembles étudiés (lignes ou colonnes du tableau analysé), soit sous la forme de *hiérarchie de partitions* que nous définirons de façon plus précise ultérieurement. Quelquefois, il s'agira d'*arbres* au sens de la théorie des

---

<sup>1</sup> Une des premières synthèses historique sur le sujet est celle de Cormack (1971). Une synthèse de travaux plus récents en classification hiérarchique a été faite par Gordon (1987). Cf. également les manuels généraux de Chandon et Pinson (1981), Jambu et Lebeaux (1978), Murtagh (1985), Roux (1985), Kaufman et Rousseeuw (1990).



graphes, arbres dont les sommets sont les objets à classer. Enfin on pourra rechercher des *classes empiétantes* ou simplement mettre en évidence des *zones à forte densité*, laissant de nombreux individus ou caractères non classés.

A une même famille de résultats correspond parfois des démarches et des interprétations différentes. Il peut s'agir de découvrir une partition ayant une existence réelle (cette existence étant conjecturée avant l'analyse statistique ou étant révélée à l'issue des calculs) ou l'on veut au contraire utiliser les partitions produites comme des outils ou des intermédiaires de calculs permettant une exploration des données<sup>1</sup>.

Pour l'essentiel, les techniques de classification font appel à une démarche algorithmique et non aux calculs formalisés usuels. Alors que les valeurs des composantes des axes factoriels, par exemple, sont la solution d'une équation pouvant s'écrire sous une forme très condensée (même si sa résolution est complexe), la définition des classes ne se fera qu'à partir d'une formulation algorithmique: une série d'opérations est définie de façon récursive et répétitive. Il en découle que la mise en œuvre de la plupart des techniques de classification ne nécessite que des notions mathématiques relativement élémentaires.

Il existe plusieurs familles d'algorithmes de classification : les algorithmes conduisant directement à des *partitions* comme les méthodes d'agrégation autour de centres mobiles; les *algorithmes ascendants* (ou encore agglomératifs) qui procèdent à la construction des classes par agglomération successive des objets deux à deux, et qui fournissent une hiérarchie de partitions des objets; enfin les *algorithmes descendants* (ou encore divisifs) qui procèdent par dichotomies successives de l'ensemble des objets, et qui peuvent encore fournir une hiérarchie de partitions. On se limitera ici aux deux premières techniques de classification :

- les groupements peuvent se faire par recherche directe d'une partition, en affectant les éléments à des centres provisoires de classes, puis en recentrant ces classes, et en affectant de façon itérative ces éléments. Il s'agit des techniques *d'agrégation autour de centres mobiles*, apparentées à la méthode des "nuées dynamiques", ou méthode "k-means", qui sont particulièrement intéressantes dans le cas des grands tableaux (section 2.1)
- les groupements peuvent se faire par agglomération progressive des éléments deux à deux. C'est le cas de la classification ascendante hiérarchique qui est présentée ici suivant plusieurs critères d'agrégations. Nous envisagerons d'une part la technique "du saut minimal" équivalente, d'un certain point de vue, à la recherche de l'arbre de longueur minimale, et d'autre part la technique d'agrégation

<sup>1</sup> Cette dernière démarche généralise en quelque sorte la construction d'histogrammes de la statistique unidimensionnelle : en vue d'une étude plus aisée, les observations sont regroupées par paquets homogènes, même si la construction de ces paquets implique un découpage quelque peu arbitraire d'un ensemble continu.

"selon la variance", intéressante par la compatibilité de ses résultats avec certaines analyses factorielles (section 2.2).

Ces techniques présentent des avantages différents et peuvent être utilisées conjointement. Il est ainsi possible d'envisager une stratégie de classification basée sur un *algorithme mixte*, particulièrement adapté au partitionnement d'ensembles de données comprenant des milliers d'individus à classer (section 2.3).

Un des avantages des méthodes de classification est de donner lieu à des éléments (les classes) souvent plus faciles à décrire automatiquement que les axes factoriels. Les outils de description seront évoqués à la section 2.3.

Enfin, la pratique montre que l'utilisateur a intérêt à utiliser de façon conjointe les méthodes factorielles et les méthodes de classification. Les aspects théoriques et pratiques de la complémentarité entre ces deux familles de méthodes exploratoires seront abordés à la section 2.4

# Agrégation autour des centres mobiles

Bien qu'elle ne fasse appel qu'à un formalisme limité et que son efficacité soit dans une large mesure attestée par les seuls résultats expérimentaux, la méthode de *classification autour de centres mobiles* est probablement la technique de partitionnement la mieux adaptée actuellement aux vastes recueils de données ainsi que la plus utilisée pour ce type d'application. Produisant des partitions des ensembles étudiés, elle est utilisée aussi bien comme technique de description et d'analyse que comme technique de réduction, généralement en association avec des analyses factorielles et d'autres méthodes de classification.

L'algorithme peut être imputé principalement à Forgy (1965), bien que de nombreux travaux (parfois antérieurs : Thorndike, 1953), le plus souvent postérieurs (MacQueen, 1967; Ball and Hall, 1967) aient été menés parallèlement et indépendamment pour introduire des variantes ou des généralisations. Cette méthode peut être considérée comme un cas particulier de techniques connues sous le nom de *nuées dynamiques* étudiées dans un cadre formel par Diday (1971).

Elle est particulièrement intéressante pour les gros fichiers numériques car les données sont traitées en *lecture directe* : le tableau des données, conservé sur une mémoire auxiliaire (disque, CD-ROM), est lu plusieurs fois de façon séquentielle, sans jamais encombrer de zones importantes dans la mémoire centrale de l'ordinateur. La lecture directe permet également d'utiliser au mieux les particularités du codage des données, ce qui réduit le temps de calcul dans le cas des codages disjonctifs.

### 2.1.1 Bases théoriques de l'algorithme

Soit un ensemble  $I$  de  $n$  individus à partitionner, caractérisés par  $p$  caractères ou variables. On suppose que l'espace  $\mathbb{R}^p$  supportant les  $n$  points-individus est muni d'une distance appropriée notée  $d$  (souvent distance euclidienne usuelle ou distance du  $\chi^2$ ). On désire constituer au maximum  $q$  classes. Les étapes de l'algorithme sont illustrées par la figure 2.1 - 1.

Étape 0 : On détermine  $q$  centres provisoires de classes (par exemple, par tirage pseudo-aléatoire sans remise de  $q$  individus dans la population à classifier, selon une préconisation de MacQueen). Les  $q$  centres :

$$\{C_1^0, \dots, C_k^0, \dots, C_q^0\}$$

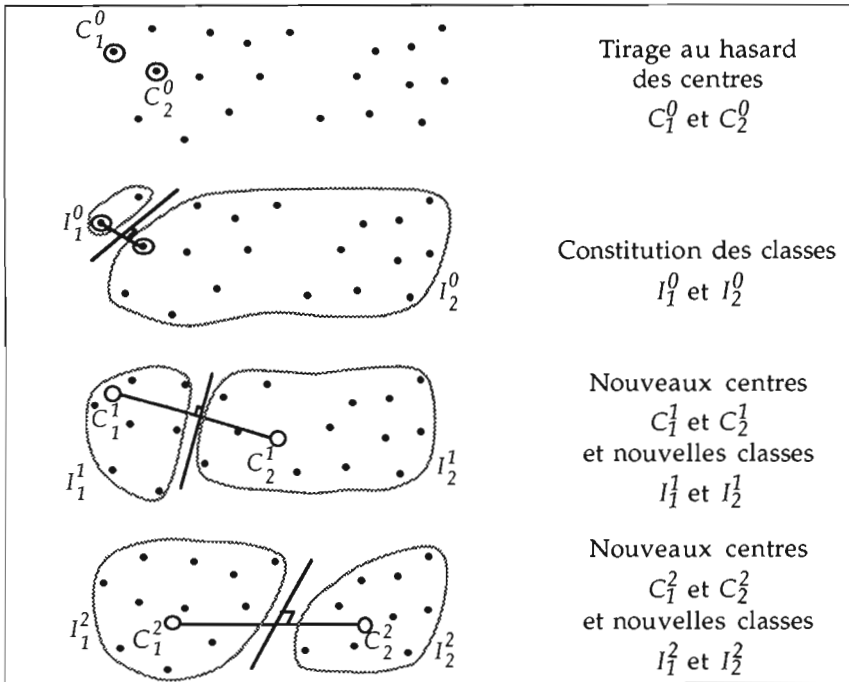


Figure 2.1 - 1  
Etapes de l'algorithme

induisent une première partition  $P^0$  de l'ensemble des individus  $I$  en  $q$  classes :

$$\{I_1^0, \dots, I_k^0, \dots, I_q^0\}$$

Ainsi l'individu  $i$  appartient à la classe  $I_k^0$  s'il est plus proche de  $C_k^0$  que de tous les autres centres<sup>1</sup>.

Étape 1: On détermine  $q$  nouveaux centres de classes :

$$\{C_1^1, \dots, C_k^1, \dots, C_q^1\}$$

en prenant les centres de gravité des classes qui viennent d'être obtenues :

$$\{I_1^0, \dots, I_k^0, \dots, I_q^0\}$$

Ces nouveaux centres induisent une nouvelle partition  $P^1$  de  $I$  construite selon la même règle que pour  $P^0$ .

<sup>1</sup> Les classes sont alors délimitées dans l'espace par les cloisons polyédrales convexes formées par les plans médiateurs des segments joignant tous les couples de centres.

La partition  $P^1$  est formée des classes notées :

$$\{I_1^1, \dots, I_k^1, \dots, I_q^1\}$$

Étape  $m$ : On détermine  $q$  nouveaux centres de classes :

$$\{C_1^m, \dots, C_k^m, \dots, C_q^m\}$$

en prenant les centres de gravité des classes qui ont été obtenues lors de l'étape précédente,

$$\{I_1^{m-1}, \dots, I_k^{m-1}, \dots, I_q^{m-1}\}$$

Ces nouveaux centres induisent une nouvelle partition  $P^m$  de l'ensemble  $I$  formée des classes :

$$\{I_1^m, \dots, I_k^m, \dots, I_q^m\}$$

Le processus se stabilise nécessairement (voir paragraphe suivant) et l'algorithme s'arrête soit lorsque deux itérations successives conduisent à la même partition, soit lorsqu'un critère convenablement choisi (par exemple, la mesure de la variance intra-classes) cesse de décroître de façon sensible, soit encore parce qu'un nombre maximal d'itérations a été fixé *a priori*.

Généralement, la partition obtenue finalement dépend du choix initial des centres.

### 2.1.2 Justification élémentaire de l'algorithme

On va montrer que la variance intra-classes ne peut que décroître (ou rester stationnaire) entre l'étape  $m$  et l'étape  $m + 1$ . Des règles d'affectation<sup>1</sup> permettent de faire en sorte que cette décroissance soit stricte et donc de conclure à la convergence de l'algorithme puisque l'ensemble de départ  $I$  est fini<sup>2</sup>.

Supposons que les  $n$  individus de l'ensemble à classer  $I$  soient munis de masses relatives  $p_i$  (leur somme vaut 1) et soit  $d^2(i, C_k^m)$  le carré de la distance entre l'individu  $i$  et le centre de la classe  $k$  à l'étape  $m$ . Nous nous intéressons à la quantité *critère* :

$$v(m) = \sum_{k=1}^q \left\{ \sum_{i \in I_k^m} p_i d^2(i, C_k^m) \right\}$$

<sup>1</sup> Ces règles sont des conventions de programmation propres à chaque variante ou spécification de l'algorithme.

<sup>2</sup> Bien entendu ce n'est pas la convergence, mais la vitesse de convergence qui justifierait en pratique l'utilisation de la méthode.

Rappelons qu'à l'étape  $m$ , la classe  $I_k^m$  est formée des individus plus proches de  $C_k^m$  que de tous les autres centres (ces centres étant des centres de gravité des classes  $I_k^{m-1}$  de l'étape précédente).

La variance intra-classes à l'étape  $m$  est la quantité :

$$V(m) = \sum_{k=1}^q \left\{ \sum_{i \in I_k^m} p_i d^2(i, C_k^{m+1}) \right\}$$

où  $C_k^{m+1}$  est le centre de gravité de la classe  $I_k^m$ . A l'étape  $m + 1$ , la quantité critère s'écrit :

$$v(m+1) = \sum_{k=1}^q \left\{ \sum_{i \in I_k^{m+1}} p_i d^2(i, C_k^{m+1}) \right\}$$

On va montrer que :

$$v(m) \geq V(m) \geq v(m+1)$$

ce qui établira la décroissance simultanée du critère et de la variance intra-classes. En notant  $p_k$  la somme des  $p_i$  pour  $i \in I_k^m$ , remarquons tout d'abord d'après le théorème de Huygens :

$$v(m) = V(m) + \sum_{k=1}^q p_k d^2(C_k^{m+1}, C_k^m)$$

ce qui établit la première partie de l'inégalité.

La seconde partie découle du fait qu'entre les accolades qui apparaissent dans les définitions de  $V(m)$  et  $v(m)$ , seules changent les affectations des points aux centres. Puisque  $I_k^{m+1}$  est l'ensemble des points plus proches de  $C_k^{m+1}$  que de tous les autres centres, les distances n'ont pu que décroître (ou rester inchangées) au cours de cette réaffectation.

### 2.1.3 Techniques connexes

Il existe de nombreux algorithmes dont le principe général est voisin de l'algorithme d'agrégation autour de centres mobiles mais qui en diffèrent cependant sur certains points<sup>1</sup>.

Ainsi, dans la technique des *nuées dynamiques* (Diday, 1972, 1974), les classes ne sont pas caractérisées par un centre de gravité, mais par un certain nombre d'individus à classer, dénommés "étalons", qui constituent alors un "noyau" ayant pour certaines utilisations un meilleur pouvoir descriptif

<sup>1</sup> Pour des informations plus détaillées sur les techniques d'agrégation autour des centres mobiles, on pourra consulter les ouvrages de Benzécri (1973) et Anderberg (1973).

que des centres ponctuels. Ce formalisme a permis plusieurs généralisations de la méthode.

La méthode dite des *k-means* (*k-moyennes*) introduite par MacQueen (1967) commence effectivement par un tirage pseudo-aléatoire de centres ponctuels. Cependant la règle de calcul des nouveaux centres n'est pas la même. On n'attend pas d'avoir procédé à la réaffectation de tous les individus pour modifier la position des centres : chaque réaffectation d'individus entraîne une modification de la position du centre correspondant. En une seule itération, cette procédure peut ainsi donner une partition de bonne qualité. Mais celle-ci dépendra de l'ordre des individus sur le fichier, ce qui n'est pas le cas pour la technique exposée précédemment<sup>1</sup>.

### 2.1.4 Formes fortes et groupements stables

Les algorithmes d'agrégation autour de centres mobiles convergent vers des *optima locaux*. Le problème de la recherche d'une partition *optimale* en  $q$  classes (en prenant comme critère la variance intra-classes, qu'il faut alors rendre minimale sur l'ensemble des partitions possibles en  $q$  classes) n'a pas jusqu'à présent donné lieu à un algorithme satisfaisant<sup>2</sup>. Les partitions obtenues dépendent en général des premiers centres choisis.

La procédure de recherche de *groupements stables* (ou encore *formes fortes*), suggérée pour l'essentiel par E. Diday (1972), permet de remédier au moins partiellement à cet inconvénient. Elle a surtout l'avantage de nuancer les résultats souvent trop frustes que l'on obtient dans le cadre rigide d'une seule partition, en mettant en évidence les zones à forte densité du nuage des points-individus. Cette technique consiste à effectuer plusieurs partitions à partir de plusieurs ensembles différents de centres, et à retenir comme *groupements stables* les ensembles d'individus qui ont toujours été affectés à une même classe dans chacune des partitions (cf. figure 2.1 - 2).

Supposons que l'on effectue  $s$  partitions  $\{P_1, P_2, \dots, P_s\}$  en  $q$  classes chacune. Dans la *partition-produit*, la classe indexée par  $\{k_1, k_2, \dots, k_s\}$  contient les individus ayant appartenu à la classe  $k_1$  de  $P_1$ , puis à la classe  $k_2$  de  $P_2$ , etc., enfin à la classe  $k_s$  de  $P_s$ . Les classes contenant plus d'un individu de la partition-produit constitueront les groupements stables.

<sup>1</sup> D'autres méthodes diffèrent par le choix initial des *centres* (individus équidistants pour Thorndike (1953), par l'introduction de *seuils* ou de *protections* destinés à modifier éventuellement le nombre des classes. Ainsi la technique proposée sous le nom *Isodata* par Ball et Hall (1965) met en jeu plusieurs paramètres destinés à piloter l'élaboration de la partition.

<sup>2</sup> Dans le cas où les individus ne sont décrits que par un seul paramètre, le calcul d'une partition optimale exacte est possible car il existe une relation d'ordre entre les individus, ce qui limite considérablement l'éventail des partitions à examiner (cf. W.D. Fisher, 1958).





**Remarque**

La recherche des groupements stables constitue une exploration des zones de fortes densité dans l'espace, mais ne fournit pas une partition utilisable en pratique, car le nombre de classes est en général trop élevé, et corrélativement les effectifs de certaines classes sont trop faibles (cf. les 50 groupements du tableau 2.1 - 2). De façon pragmatique, on peut utiliser les premiers groupements stables pour définir une partition de la façon suivante : le nombre de classes pourra être suggéré par le nombre de groupements d'effectifs notables : ainsi, les 7 premiers groupements du tableau 2.1 - 2 ont des effectifs importants (il y a de plus un écart important entre 78 et 26). Les classes seront obtenues par réaffectation des individus restants aux groupements retenus les plus proches (affectation des individus des groupements 8 à 50 autour des centres des 7 premiers groupements pour notre exemple). Mais nous verrons que les méthodes mixtes de la section 2.3 permettent de perfectionner cette démarche.

## Classification hiérarchique

Les principes généraux communs aux diverses techniques de classification ascendante hiérarchique sont également extrêmement simples. Il est difficile de leur trouver une paternité car ces principes relèvent plus du bon sens que d'une théorie formalisée. Les exposés les plus systématiques et les plus anciens sont peut-être ceux de Sokal et Sneath (1963), puis de Lance et Williams (1967). Pour une revue synthétique, cf. Gordon (1987).

### 2.2.1 Principe

Le principe de l'algorithme consiste à créer, à chaque étape, une partition obtenue en agrégeant deux à deux les éléments les plus proches. On désignera alors par *élément* à la fois les individus ou objets à classer eux-mêmes et les regroupements d'individus générés par l'algorithme. Il y a différentes manières de considérer le nouveau couple d'éléments agrégés, d'où un nombre important de variantes de cette technique.

L'algorithme ne fournit pas une partition en  $q$  classes d'un ensemble de  $n$  objets mais une *hiérarchie de partitions*, se présentant sous la forme d'*arbres* appelés également *dendrogrammes* et contenant  $n - 1$  partitions. L'intérêt de ces arbres est qu'ils peuvent donner une idée du nombre de classes existant effectivement dans la population.

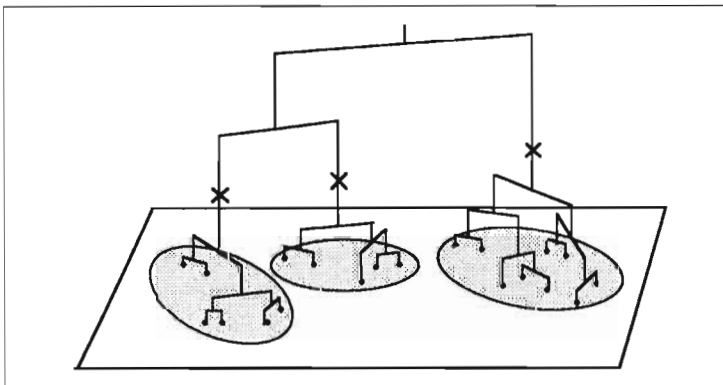


Figure 2.2 - 1  
Dendrogramme ou arbre hiérarchique

Chaque coupure d'un arbre fournit une partition, ayant d'autant moins de classes et des classes d'autant moins homogènes que l'on coupe plus haut.

### a – Distances entre éléments et entre groupes

On suppose au départ que l'ensemble des individus à classer est muni d'une *distance*<sup>1</sup>. Ceci ne suppose donc pas que les distances soient toutes calculées au départ : il faut pouvoir les calculer ou les recalculer à partir des coordonnées des points-individus, celles-ci devant être accessibles rapidement. On construit alors une première matrice de distances entre tous les individus.

Une fois constitué un groupe d'individus, il convient de se demander ensuite sur quelle base on peut calculer une distance entre un individu et un groupe et par la suite une distance entre deux groupes. Ceci revient à définir une stratégie de regroupements des éléments, c'est-à-dire se fixer des *règles de calcul des distances entre groupements* disjoints d'individus, appelées *critères d'agrégation*. Cette distance entre groupements pourra en général se calculer directement à partir des distances des différents éléments impliqués dans le regroupement.

Par exemple, si  $x, y, z$  sont trois objets, et si les objets  $x$  et  $y$  sont regroupés en un seul élément noté  $h$ , on peut définir la distance de ce groupement à  $z$  par la plus petite distance des divers éléments de  $h$  à  $z$  :

$$d(h,z) = \text{Min} \{d(x,z), d(y,z)\}$$

Cette distance s'appelle le *saut minimal (single linkage)* (Sneath,1957 ; Johnson,1967) et constitue un critère d'agrégation.

On peut également définir la distance du *saut maximal* (ou diamètre) en prenant la plus grande distance des divers éléments de  $h$  à  $z$  :

$$d(h,z) = \text{Max} \{d(x,z), d(y,z)\}$$

Une autre règle simple et fréquemment employée est celle de la *distance moyenne* ; pour deux objets  $x$  et  $y$  regroupés en  $h$  :

$$d(h,z) = \frac{d(x,z) + d(y,z)}{2}$$

Plus généralement, si  $x$  et  $y$  désignent des sous-ensembles disjoints de l'ensemble des objets, ayant respectivement  $n_x$  et  $n_y$  éléments,  $h$  est alors un sous-ensemble formé de  $n_x + n_y$  éléments et on définit :

$$d(h,z) = \frac{\{n_x d(x,z) + n_y d(y,z)\}}{n_x + n_y}$$

### b – Algorithme de classification

L'algorithme fondamental de classification ascendante hiérarchique se déroule de la façon suivante :

<sup>1</sup> Il s'agira parfois simplement d'une mesure de dissimilarité. Dans ce cas, l'inégalité triangulaire  $d(x,y) \leq d(x,z) + d(y,z)$  n'est pas exigée).

- Étape 1 : il y a  $n$  éléments à classer (qui sont les  $n$  individus);
- Étape 2 : on construit la matrice de distances entre les  $n$  éléments et l'on cherche les deux plus proches, que l'on agrège en un nouvel élément. On obtient une première partition à  $n-1$  classes;
- Étape 3 : on construit une nouvelle matrice des distances qui résultent de l'agrégation, en calculant les distances entre le nouvel élément et les éléments restants (les autres distances sont inchangées). On se trouve dans les mêmes conditions qu'à l'étape 1, mais avec seulement  $(n-1)$  éléments à classer et en ayant choisi un critère d'agrégation. On cherche de nouveau les deux éléments les plus proches, que l'on agrège. On obtient une deuxième partition avec  $n-2$  classes et qui englobe la première;
- Étape  $m$  : on calcule les nouvelles distances, et l'on réitère le processus jusqu'à n'avoir plus qu'un seul élément regroupant tous les objets et qui constitue la dernière partition.

Nous illustrons cette procédure en prenant comme objets à classer cinq points (figure 2.2 - 2).

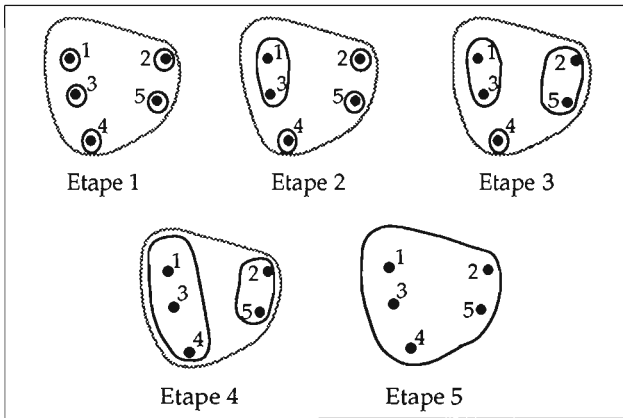


Figure 2.2 - 2  
Agglomération progressive de 5 points

Les regroupements successifs peuvent être représentés par un arbre ou dendrogramme, comme le montre la figure 2.2 - 3 où l'on a porté en ordonnée les valeurs des indices ou encore distances correspondant aux différents niveaux d'agrégation.

### c – Éléments de vocabulaire

Quelques remarques vont nous permettre d'introduire les notions et la terminologie habituellement utilisées en classification ascendante hiérarchique. Le fonctionnement de l'algorithme nous montre que les distances (avec ces règles de calcul) n'interviennent que par les *inégalités*

qui existent entre elles. Le même arbre (à une dilatation près des ordonnées) aurait été obtenu à partir d'un simple classement des couples d'objets dans l'ordre des distances croissantes. Un tel classement s'appelle une *ordonnance* (une *préordonnance* s'il y a des distances égales). Dans ce cas on tracera conventionnellement l'arbre avec des niveaux équidistants.

La famille  $H$  des parties de l'ensemble  $I$  des objets construite à partir d'algorithmes ascendants forme ce que l'on appelle une *hiérarchie*. Cette famille a pour propriété de contenir l'ensemble tout entier ( $I \in H$ ) ainsi que chacun des objets pris isolément ( $i \in I \Rightarrow \{i\} \in H$ ). Les autres couples de parties  $h, h'$  de  $H$  sont alors soit disjointes ( $h \cap h' = \emptyset$ ), soit incluses l'une dans l'autre ( $h \subset h'$ ). En effet lors du fonctionnement de l'algorithme, chaque fois qu'une classe se forme à partir d'éléments disjoints, elle est elle-même considérée comme un nouvel élément, donc strictement incluse dans une classe ultérieure (cf. figure 2.2 - 2).

Les objets ou individus (1, 2, 3, 4, 5) sont les *éléments terminaux* de l'arbre (ou de la hiérarchie). Les classes 6, 7, 8, 9 sont les nœuds de l'arbre : ce sont des classes issues de regroupements de deux éléments (terminaux ou non) numérotés à la suite des éléments terminaux et dont chacune détermine une nouvelle partition. On appelle arbitrairement *aîné* et *benjamin*, les deux éléments groupés constituant un nœud (cf. figure 2.2 - 3).

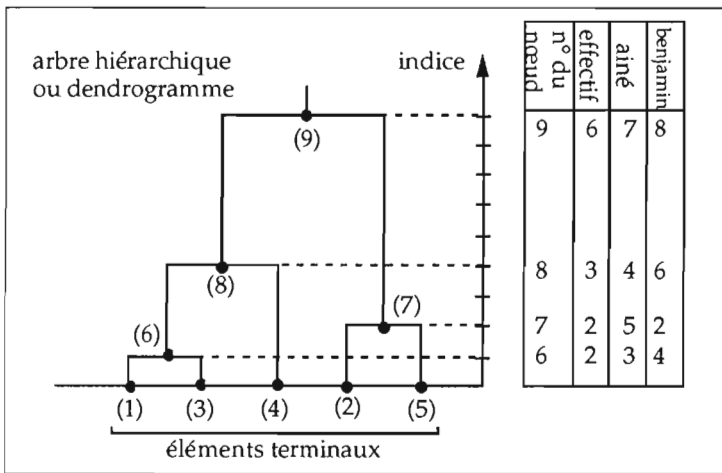


Figure 2.2 - 3  
Arbre hiérarchique et éléments de vocabulaire

On a une *hiérarchie indiquée* si à toute partie  $h$  de la hiérarchie est associée une valeur numérique  $v(h) \geq 0$  compatible avec la relation d'inclusion au sens suivant:

$$\text{si } h \subset h' \text{ alors } v(h) < v(h')$$

La hiérarchie de la figure 2.2 - 3 est indiquée de façon naturelle par les valeurs des distances correspondant à chaque étape d'agrégation (ces distances sont

portées en ordonnées). L'*indice* est la distance déterminant le regroupement.

En "coupant" l'arbre de la figure 2.2 - 3 par une droite horizontale, on obtient une *partition*, d'autant plus fine que la section est proche des éléments terminaux. Si par exemple l'indice est supérieur à 4 et inférieur à 9, on obtient une partition en deux classes {1,3,4} et {2,5}. Si il vaut 3, on obtient trois classes {1,3}, {4} et {2,5}. Une hiérarchie permet donc de fournir une chaîne de  $n$  partitions ayant de 1 à  $n$  classes.

## 2.2.2 Classification ascendante selon le saut minimal et arbre de longueur minimale

Ce mode de classification hiérarchique, présenté lors de l'illustration du paragraphe précédent, est particulièrement simple à mettre en œuvre et possède des propriétés intéressantes que nous allons énoncer et étudier.

### a – Définition d'une ultramétrie

Nous allons montrer que la notion de hiérarchie est étroitement liée à une certaine classe de distances entre individus, que l'on appelle les *distances ultramétriques*. Pour la hiérarchie produite par l'algorithme du saut minimal, on montrera que la distance ultramétrique correspondante est, dans un certain sens, la plus proche de la distance initiale. Ce sera l'*ultramétrie inférieure maximale*, appelée encore *sous-dominante*. On montrera ensuite que l'application de cette méthode est pratiquement équivalente à la résolution d'un problème classique de recherche opérationnelle : la mise en évidence de l'*arbre de longueur minimale* sur un graphe. Rappelons qu'un ensemble  $E$  est muni d'une *métrie* ou *distance*  $d$ , si  $d$  est une application à valeurs positives ou nulles obéissant aux conditions suivantes :

1.  $d(x,y) = 0$  si et seulement si  $x = y$ .
2.  $d(x,y) = d(y,x)$  (symétrie)
3.  $d(x,y) \leq d(x,z) + d(y,z)$  (inégalité triangulaire)

Cette distance sera dite *ultramétrique* si elle vérifie la condition suivante, plus forte que l'inégalité triangulaire :

4.  $d(x,y) \leq \text{Max} \{ d(x,z), d(y,z) \}$

### b – Équivalence entre ultramétrie et hiérarchie indicée

Il est équivalent de munir un ensemble fini  $E$  d'une ultramétrie ou de définir une hiérarchie indicée de parties de cet ensemble. Montrons tout d'abord que toute hiérarchie indicée permet de définir une distance entre éléments ayant les propriétés requises. On prendra comme distance  $d(x,y)$  la

valeur de l'indice correspondant à la plus petite partie contenant à la fois  $x$  et  $y$ .

En remplissant ainsi le tableau des valeurs de  $d$  correspondant à la hiérarchie de la figure 2.2 - 3, on obtient la matrice des distances du tableau 2.2 - 1. On peut noter que l'inégalité 4 ci-dessus est vérifiée par toutes les distances de ce tableau. Ainsi par exemple :

$$d(1,2) \leq \text{Max} \{ d(1,5), d(2,5) \}$$

Tableau 2.2 1  
Matrice des distances

	(1)	(2)	(3)	(4)	(5)
(1)	0	9	1	4	9
(2)	9	0	9	9	2
(3)	1	9	0	4	9
(4)	4	9	4	0	9
(5)	9	2	9	9	0

Montrons plus généralement que l'on a toujours :

$$d(x,y) \leq \text{Max} \{ d(x,z) + d(y,z) \}$$

Rappelons que deux parties de la hiérarchie  $H$  sont soit disjointes, soit liées par une relation d'inclusion. Appelons  $h(x, z)$  la plus petite partie de  $H$  contenant  $x$  et  $z$  (dont l'indice est par conséquent  $d(x, z)$ ). Puisque  $h(x, z)$  et  $h(y, z)$  ne sont pas disjointes, on a par exemple  $h(x, z) \subset h(y, z)$ . Et  $x, y, z$  étant tous trois contenus dans  $h(y, z)$ , on a obligatoirement :

$$h(x, y) \subset h(y, z) \quad \text{d'où} \quad d(x,y) \leq d(y,z)$$

ce qui établit l'inégalité.

Réciproquement, à toute ultramétrique  $d$  on peut faire correspondre une hiérarchie indicée dont  $d$  soit l'indice associé. Il suffit d'appliquer l'algorithme du saut minimal au tableau des distances correspondant. On s'aperçoit alors qu'il est inutile de procéder au calcul des distances à chaque étape : il suffira de rayer l'un des deux éléments agrégés.

En effet, si  $x$  et  $y$  sont agrégés en  $t$ , il faut en principe calculer les distances au nouvel élément  $t$  (cf. figure 2.2 - 4). Or on a obligatoirement, pour tout élément  $z$  non encore agrégé,  $d(z,x) \geq d(x,y)$  et  $d(z,y) \geq d(x,y)$ , sinon  $(z,x)$  ou  $(z,y)$  auraient été agrégés à la place de  $(x,y)$ .

Pour une ultramétrique, cela implique à la fois  $d(z,x) \geq d(z,y)$ , et  $d(z,y) \geq d(z,x)$  c'est-à-dire  $d(z,x) = d(z,y)$ , ce que l'on exprime de façon imagée

en disant que, pour une ultramétrie, tous les triangles sont isocèles, avec le plus petit coté pour base (figure 2.2 - 4).

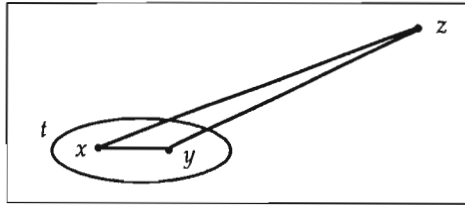


Figure 2.2 - 4  
Agrégation de  $x$  et  $y$  en un nouvel élément  $t$

Il est en effet facile de montrer que si une distance est ultramétrique, tous les triangles sont isocèles.

On a les inégalités :

$$d(z, x) \leq \text{Max} \{ d(x, y), d(y, z) \} \quad \text{donc} \quad d(z, x) \leq d(y, z)$$

De la même façon :

$$d(y, z) \leq \text{Max} \{ d(x, y), d(z, x) \} \quad \text{donc} \quad d(y, z) \leq d(z, x)$$

Il s'ensuit que :

$$d(z, x) = d(y, z)$$

Le calcul des distances de  $z$  à  $t$  est finalement inutile puisque les deux distances mises en cause sont égales. Ceci nous montre comment l'algorithme du saut minimal a opéré sur la matrice des distances : il a transformé la métrique initiale en ultramétrie en diminuant certaines distances à chaque étape.

### c – L'ultramétrie sous dominante

Le passage d'une métrique à une ultramétrie (ou, de façon équivalente, à une hiérarchie) s'est effectué par diminution des valeurs de certaines distances. On peut se poser la question suivante : existe-t-il une ultramétrie plus proche (en un sens à préciser) de la métrique ?

On peut donner l'élément de réponse suivant. On dira qu'une métrique  $d_1$  est inférieure<sup>1</sup> à une métrique  $d_2$  si, pour tout  $x$  et tout  $y$  :

$$d_1(x, y) \leq d_2(x, y)$$

La plus grande ultramétrie inférieure à une métrique  $d$ , au sens précédent, est appelée ultramétrie inférieure maximale ou sous-dominante. C'est elle qui est fournie par l'algorithme du saut minimal.

Pour le démontrer nous allons successivement :

<sup>1</sup> Cette définition permet de munir l'ensemble des métriques définies sur un ensemble  $E$  d'une relation d'ordre partiel.



1. définir, à partir d'une distance  $d$ , une nouvelle distance dite du plus petit saut maximal;
2. montrer que cette distance est une ultramétrie;
3. montrer que cette ultramétrie est la sous-dominante;
4. montrer enfin que cette distance correspond à l'ultramétrie fournie par l'algorithme du saut minimal.

1. *La distance du plus petit saut maximal :*

Soit un ensemble  $E$  muni d'une distance  $d$ . Soit  $x$  et  $y$  deux éléments de  $E$ . Le couple  $(x, y)$  sera appelé *arête* de longueur  $d(x, y)$  du *graphe complet*<sup>1</sup> dont les sommets sont les éléments de  $E$ . Toujours en utilisant le vocabulaire de la théorie des graphes, on appelle *chemin* de  $x$  à  $y$  une succession d'arêtes de types  $(x, t_1), (t_1, t_2), (t_2, t_3), \dots, (t_{k-1}, t_k), (t_k, y)$ , où  $t_1, \dots, t_k$  sont des éléments de  $E$ . Étant donné un chemin de  $x$  à  $y$ , on appelle *saut maximal* la longueur de la plus grande arête du chemin de  $x$  à  $y$ .

A tout chemin joignant  $x$  à  $y$  correspond un saut maximal. L'ensemble des sommets étant fini, il existe un *plus petit saut maximal* sur l'ensemble des chemins allant de  $x$  à  $y$ ; nous le noterons  $d^*(x, y)$ .

2. *Le plus petit saut maximal entre  $x$  et  $y$  est une ultramétrie :*

Il est clair que les deux premiers axiomes d'une distance sont vérifiés par  $d^*$ . Pour vérifier que cette distance est une ultramétrie, considérons trois éléments quelconques  $x, y, z$  de  $E$  (figure 2.2 - 5). Le plus petit saut maximal de  $x$  à  $y$ , en s'astreignant à passer par  $z$  est  $\text{Max} \{d^*(x, z), d^*(z, y)\}$ . Le plus petit saut maximal de  $x$  à  $y$  sans la contrainte de passer par  $z$  ne peut qu'être inférieur ou égal à cette quantité, d'où :

$$d^*(x, y) \leq \text{Max} \{d^*(x, z), d^*(y, z)\}$$

et  $d^*$  est donc bien une ultramétrie.

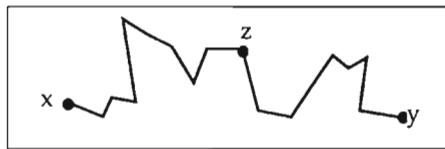


Figure 2.2 - 5  
Chemin de  $x$  à  $y$  contenant  $z$

3. *La distance  $d^*$  est la sous-dominante :*

Pour montrer que  $d^*$  est la sous-dominante, on montrera que  $d^*$  est inférieure à  $d$ , et que  $d^*$  est supérieure à toute ultramétrie inférieure à  $d$ .

<sup>1</sup> L'appellation *graphe complet* est due au fait que tout couple de sommets est joint par une arête.

Tout d'abord, il est clair que l'arête  $(x,y)$  est un chemin particulier allant de  $x$  à  $y$ , donc  $d^*(x,y) \leq d(x,y)$  et  $d^*$  est inférieure à  $d$ .

Soit maintenant  $d_1$  une ultramétrie inférieure à  $d$ . On a évidemment pour tout triplet  $x_1, x_2, x_3$  :

$$d_1(x_1, x_3) \leq \text{Max} \{d_1(x_1, x_2), d_1(x_2, x_3)\}$$

En appliquant de façon successive cette inégalité à un chemin :

$$(x_1, x_2), (x_2, x_3), \dots, (x_{p-1}, x_p)$$

on obtient :

$$d_1(x_1, x_p) \leq \text{Max}\{d_1(x_j, x_{j+1})\}_{j < p}$$

Puisque  $d_1 \leq d$ , on a :

$$d_1(x_1, x_p) \leq \text{Max}\{d(x_j, x_{j+1})\}_{j < p}$$

Cette inégalité est valable pour tout chemin joignant  $x_1$  à  $x_p$ . Pour l'un au moins d'entre eux, on a par définition de  $d^*$  :

$$\text{Max}\{d(x_j, x_{j+1})\}_{j < p} = d^*(x_j, x_{j+1})$$

Cette dernière relation établit l'inégalité annoncée.

4. La distance ultramétrique  $d_u$  produite par l'algorithme du saut minimal n'est autre que la distance  $d^*$  plus petit saut maximal :

Soit  $d_u(x, y)$  la valeur de la distance à l'étape où les points  $x$  et  $y$  sont réunis pour la première fois. Auparavant ces deux points étaient dans des classes distinctes (éventuellement réduites aux points eux-mêmes). Le mode de calcul des distances à chaque agrégation nous assure que  $d_u(x, y)$  est la plus petite distance entre deux éléments appartenant chacun à une classe. Les distances à l'intérieur des classes sont inférieures à  $d_u(x, y)$  puisque l'agrégation est antérieure ; et les distances avec des éléments n'appartenant pas aux deux classes sont supérieures puisque ceux-ci seront agrégés à une étape ultérieure. Les chemins joignant  $x$  et  $y$  auront donc des arêtes internes aux deux classes, de longueur inférieure à  $d_u(x, y)$  et des arêtes externes nécessairement supérieures ou égales à  $d_u(x, y)$ . Ainsi  $d_u(x, y)$  est bien le plus petit saut maximal  $d^*(x, y)$ .

## d – Arbre de longueur minimale : définition et généralités

L'ensemble des  $n$  objets à classer peut être considéré comme un ensemble de points d'un espace. Cette représentation est classique si les objets sont décrits par une série de  $p$  variables : on a  $n$  points dans l'espace  $\mathbb{R}^p$ . On peut alors calculer une distance pour chaque paire de points. Plus généralement, si l'on ne dispose que des valeurs d'un indice de dissimilarité (ne vérifiant pas obligatoirement tous les axiomes d'une distance), on peut représenter les objets par des points (d'un plan par exemple), chaque couple d'objets étant

joint par une ligne continue, à laquelle est attachée la valeur de l'indice de dissimilarité.

On représente ainsi l'ensemble des objets et des valeurs de l'indice par un *graphe complet valué*<sup>1</sup>. Mais si le nombre d'objets dépasse quelques unités, ce type de représentation devient inextricable. On cherchera alors à extraire de ce graphe un *graphe partiel* (ayant les mêmes sommets, mais moins d'arêtes) plus aisé à représenter, et permettant néanmoins de bien résumer les valeurs de l'indice.

Parmi tous les graphes partiels, ceux qui ont une structure d'*arbre*<sup>2</sup> sont particulièrement intéressants, car ils peuvent faire l'objet d'une représentation plane. Un arbre est un *graphe connexe* (il existe un chemin reliant tout couple de sommets) *sans cycle* (un cycle est un chemin partant et aboutissant au même point sans emprunter deux fois la même arête). On peut définir de façon équivalente un arbre à  $n$  sommets soit comme un graphe sans cycle ayant  $n-1$  arêtes, soit comme un graphe connexe ayant  $n-1$  arêtes<sup>3</sup>. La *longueur* d'un arbre sera la somme des "longueurs" (valeurs de l'indice) de ses arêtes. Parmi tous les graphes partiels qui sont des arbres, l'*arbre de longueur minimale* a retenu depuis longtemps l'attention des statisticiens en raison de ses bonnes qualités descriptives, qui ne sont pas étrangères à sa parenté avec les classifications hiérarchiques. Si l'on désire par exemple déceler rapidement sans ordinateur les traits de structure que peut cacher une matrice de corrélations relative à une trentaine de variables, c'est probablement la plus aisée des procédures à mettre en œuvre. Nous allons tout d'abord présenter les algorithmes de recherche de l'arbre de longueur minimale, puis nous montrerons les équivalences avec la classification selon le saut minimal. Nous supposerons que toutes les arêtes du graphe ont des longueurs différentes (valeurs de l'indice ou de la distance) car dans ces conditions l'arbre cherché est unique et ceci simplifie l'exposé des algorithmes.

### e – Arbre de longueur minimale : algorithme de Kruskal (1956)

On range les  $n(n-1)/2$  arêtes dans l'ordre des valeurs croissantes de l'indice. On part des deux premières arêtes, puis on sélectionne successivement toutes les arêtes qui ne font pas de cycle avec les arêtes déjà choisies. On interrompt la procédure dès que l'on a  $n-1$  arêtes. De cette façon, on est sûr d'avoir obtenu un arbre (graphe sans cycle ayant  $n-1$  arêtes).

<sup>1</sup> Les objets à classer sont alors les nœuds du graphe (non orienté); les lignes continues joignant les paires de points sont les arêtes; et les indices, les valuations de ces arêtes.

<sup>2</sup> On ne confondra pas un tel arbre, entendu au sens de la théorie des graphes, et dont les sommets sont les objets à classer, avec l'arbre des parties d'un ensemble (dendrogramme) produit par les techniques de classification hiérarchique, dont les sommets sont des parties (à l'exception des éléments terminaux qui sont les objets à classer eux-mêmes).

<sup>3</sup> On trouvera la démonstration de ces propriétés dans les manuels classiques tels que ceux de Berge (1963, 1973).

Montrons en effet que si  $V_k$  dénote le graphe obtenu à l'étape  $k$ , après avoir sélectionné les arêtes  $v_1, v_2, \dots, v_k$ , alors  $V_{n-1}$  est de longueur minimale. Supposons qu'il existe un arbre distinct  $U$ , de longueur minimale (figure 2.2 - 6). Soit  $v_k$  la première arête sélectionnée dans la construction de  $V_{n-1}$  et qui n'appartienne pas à  $U$  (les arêtes de  $V_{k-1}$  sont donc également des arêtes de  $U$ ). En ajoutant cette arête à  $U$  on crée nécessairement un cycle (car  $U$  est connexe) et un seul (car  $U$  est sans cycle). Il existe donc une arête  $u$  de ce cycle qui n'appartient pas à  $V_{n-1}$  (puisque  $V_{n-1}$  n'a pas de cycle). Alors l'arbre  $U^*$  obtenu à partir de  $U$  en ajoutant  $v_k$  et en supprimant  $u$  est plus court que  $U$ . En effet, le graphe obtenu en ajoutant  $u$  à  $V_{k-1}$  est sans cycle (c'est une partie de  $U$ ); donc  $u$  est plus long que  $v_k$ , par définition de  $v_k$ , et par conséquent  $U^*$  est plus court que  $U$ . Mais ceci contredit la définition de  $U$ . Donc  $V_{n-1}$  est bien de longueur minimale.

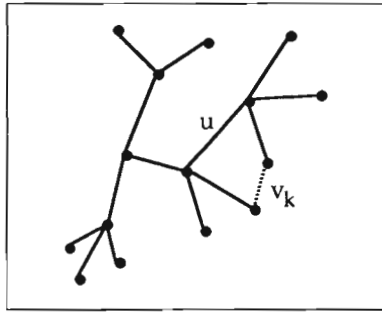


Figure 2.2 - 6  
Représentation de l'arbre  $U$

### f – Arbre de longueur minimale : algorithme de Prim (1957)

On part d'un objet quelconque (sommets du graphe). L'étape 1 consiste à chercher l'objet  $v_1$  le plus proche, c'est-à-dire l'arête la plus courte. L'étape  $k$  consiste à adjoindre au recueil d'arêtes déjà constitué  $V_{k-1}$  la plus courte arête  $v_k$  qui touche un des sommets de  $V_{k-1}$ . Il y a  $n-1$  étapes. Cet algorithme est plus rapide que le précédent. L'arbre obtenu est de longueur minimale car  $V_k$  est à tout moment un arbre de longueur minimale sur les  $k$  sommets concernés.

### g – Arbre de longueur minimale : algorithme de Florek (1951)

A la première étape, on joint chaque sommet à son voisin le plus proche. Cela revient à prendre la plus petite distance dans chaque ligne du tableau des distances. Cette opération rapide produit une forêt  $F_1$  (famille d'arbres, c'est-à-dire simplement : graphe sans cycle). A l'étape  $k$ , chaque arbre de la forêt  $F_{k-1}$  (chaque composante connexe du graphe sans cycle) est joint à son plus proche voisin en prenant comme distance entre arbres la plus petite distance entre un sommet quelconque de l'un et un sommet quelconque de

l'autre. Le processus s'arrête dès que le graphe  $F_k$  est connexe. Cet algorithme est plus rapide à mettre en œuvre manuellement sur des tableaux de distances assez grands. En général, il n'y a que 2 ou 3 étapes.

Montrons que l'on obtient un arbre, ce qui se ramène à prouver que la première étape fournit bien une forêt. Il n'y a pas de sommet isolé car chaque sommet admet effectivement un plus proche voisin. Montrons par l'absurde que l'on ne peut pas créer de cycle. Supposons qu'il en existe un et orientons les arêtes de chaque sommet vers son plus proche voisin. Si les arêtes du cycle sont toutes orientées dans le même sens, le résultat est absurde, car celles-ci seraient nécessairement de plus en plus courtes. Sinon la figure serait également absurde, car deux arêtes partirait d'un même sommet, alors que chaque sommet n'a qu'un seul plus proche voisin.

Il reste à montrer que cet arbre est de longueur minimale. Notons que toute arête tracée à la première étape appartient à l'arbre de longueur minimale  $V$ . En effet, s'il n'en était pas ainsi, il existerait  $y$ , plus proche voisin de  $x$ , tel que l'arête  $(x,y)$  n'appartienne pas à  $V$ . En ajoutant cette arête à  $V$ , on crée un cycle. En supprimant l'autre arête du cycle issue de  $x$ , on obtient un nouvel arbre plus court que  $V$ , ce qui contredit la définition de  $V$ . De la même façon, toute arête tracée à l'étape  $k$  appartient à  $V$ , sachant que la forêt  $F_{k-1}$  est une partie de  $V$ . Le raisonnement est en tout point analogue au précédent.

### **h – Lien entre l'arbre et le saut minimal (Gower et Ross, 1969)**

Soit  $V$  un arbre de longueur minimale construit à partir du tableau des distances entre  $n$  objets.  $V$  étant connexe et n'ayant pas de cycle, il existe un chemin et un seul joignant deux sommets  $x$  et  $y$ . Appelons  $d_v(x, y)$  la longueur de la plus grande arête rencontrée sur ce chemin. Nous allons montrer que  $d_v(x, y)$  n'est autre que  $d^*(x, y)$ , la distance ultramétrique du plus petit saut maximal entre  $x$  et  $y$ .

En effet, soit  $v$  la plus grande arête rencontrée entre  $x$  et  $y$ . La suppression de  $v$  entraîne la division de  $V$  en deux composantes connexes séparées. S'il existe un chemin (n'empruntant pas obligatoirement des arêtes de  $V$ ) de  $x$  à  $y$  dont la plus grande arête est plus courte que  $v$ , il existe une arête  $u$  distincte de  $v$ , et plus courte qui joint les deux composantes connexes. Le fait de remplacer  $v$  par  $u$  donnerait un arbre de longueur inférieure à celle de  $V$ , ce qui contredit la définition de  $V$ . Ainsi  $d_v(x, y)$ , longueur de  $v$ , est bien le plus petit saut maximal.

Le raisonnement fournit un mode de construction de la hiérarchie associée au saut minimal, à partir de l'arbre de longueur minimale  $V$ . Cette construction, descendante, s'opère de la façon suivante. On rompt la plus grande arête de  $V$ ; on obtient ainsi les deux groupes les plus éloignés, l'indice correspondant à leur fusion étant la longueur de cette arête. On rompt ensuite successivement les arêtes par ordre de grandeur décroissantes, ce qui fait descendre dans la hiérarchie jusqu'aux éléments terminaux qui sont les

objets eux-mêmes. La dernière arête rompue correspond aux deux objets agrégés en premier dans l'algorithme ascendant.

On peut représenter simultanément la hiérarchie et l'arbre de longueur minimale en perspective comme le montre la figure 2.2 - 7.

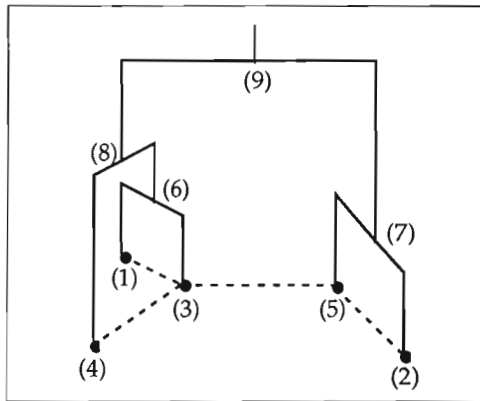


Figure 2.2 - 7  
Représentation simultanée de la hiérarchie  
et de l'arbre de longueur minimale

Quelques informations complémentaires sont apportées à la représentation de la figure 2.2 - 3 (cf. Benzécri et Jambu, 1976). En particulier les positions relatives des points sont mieux respectées. Pour le praticien de l'analyse factorielle, il sera souvent intéressant de porter l'arbre de longueur minimale sur les plans factoriels de façon à remédier, dans une certaine mesure, aux possibles déformations imputables à l'opération de projection.

### 2.2.3 Critère d'agrégation selon la variance

Les techniques de classification selon le saut minimal ont l'avantage de conduire à des calculs simples (pas de recalcul numérique des distances) et possèdent des propriétés mathématiques intéressantes.

Pour certaines applications les résultats sont cependant critiquables. En particulier, le saut minimal a le défaut de produire des "effets de chaîne".

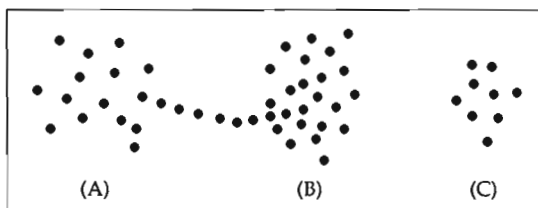


Figure 2.2 - 8  
"Effets de chaîne"

Ainsi pour le nuage de points représenté par la figure 2.2 - 8 les groupes A et B ne seront pas facilement discernables dans l'arbre hiérarchique; de plus, les quelques sommets qui les relient seront agrégés au niveau le plus bas.

D'autres critères d'agrégation donnent éventuellement des résultats plus fiables, par exemple la distance moyenne (cf. également Wishart, 1969).

Les techniques d'agrégation selon la variance cherchent à optimiser, à chaque étape, selon des critères liés à des calculs d'inertie, la partition obtenue par agrégation de deux éléments. Cette technique est particulièrement aisée à mettre en œuvre lorsque l'agrégation est effectuée après une analyse factorielle, les objets à classer étant repérés par leurs coordonnées sur les premiers axes factoriels.

### a – Notations et principe

Nous considérons ici les  $n$  objets à classer comme un nuage de points (le nuage des individus) d'un espace à  $p$  dimensions (espace des variables).

Chaque point  $x_i$  (vecteur à  $p$  composantes) est muni d'une masse  $m_i$ . On note  $m$  la masse totale du nuage :

$$m = \sum_i^n m_i$$

Le carré de la distance entre les points  $x_i$  et  $x_{i'}$  est notée :

$$\|x_i - x_{i'}\|^2 = d^2(x_i, x_{i'})$$

L'inertie totale  $I$  du nuage est la quantité :

$$I = \sum_i^n m_i \|x_i - g\|^2$$

où  $g$  désigne le centre de gravité du nuage :

$$g = \frac{1}{m} \sum_i^n m_i x_i$$

S'il existe une partition de l'ensemble des éléments en  $s$  classes, la  $q^{\text{ième}}$  classe a pour masse :

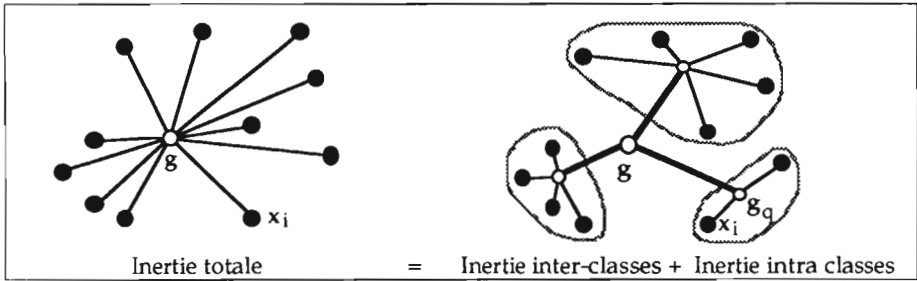
$$m_q = \sum_{i \in q} m_i$$

et pour centre de gravité :

$$g_q = \frac{1}{m_q} \sum_{i \in q} m_i x_i$$

La relation de Huygens fournit une décomposition de la quantité  $I$  en inerties intra-classes et inter-classes suivant la formule :

$$I = \sum_q m_q \|g_q - g\|^2 + \sum_q \sum_{i \in q} m_i \|x_i - g_q\|^2 \quad [2.2 - 1]$$



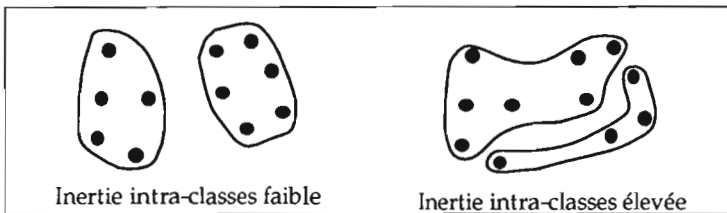
**Figure 2.2 - 9**  
**Décomposition de l'inertie selon la relation de Huygens**

La qualité globale d'une partition est liée à l'homogénéité à l'intérieur des classes (et donc à l'écartement entre les classes).  $I$  étant une quantité constante, il s'agit par conséquent de minimiser la quantité relative à l'inertie intra-classes :

$$I_{intra} = \sum_q \sum_{i \in q} m_i \|x_i - g_q\|^2$$

soit encore à maximiser celle relative à l'inertie inter-classes :

$$I_{inter} = \sum_q m_q \|g_q - g\|^2$$



**Figure 2.2 - 10**  
**Qualité globale d'une partition**

A l'étape initiale, l'inertie intra-classes est nulle et l'inertie inter-classes est égale à l'inertie totale du nuage puisque chaque élément terminal constitue à ce niveau une classe. A l'étape finale, c'est l'inertie inter-classes qui est nulle et l'inertie intra-classes est équivalente à l'inertie totale puisque l'on dispose à ce niveau d'une partition en une seule classe (cf. l'étape 5 de la figure 2.2 - 2). Par conséquent, au fur et à mesure que l'on effectue des regroupements, l'inertie intra-classes augmente et l'inertie inter-classes diminue.

Le principe de l'algorithme d'agrégation selon la variance consiste à rechercher à chaque étape une partition telle que la variance interne de chaque classe soit minimale et par conséquent la variance entre les classes soit maximale.



### b – Perte d'inertie par agrégation de deux éléments : le critère de Ward généralisé

Faire varier le moins possible l'inertie intra-classes à chaque étape d'agrégation revient à rendre minimale la perte d'inertie inter-classes résultant de l'agrégation de deux éléments.

Soit  $x_i$  et  $x_{i'}$  deux éléments de masses  $m_i$  et  $m_{i'}$  appartenant à une partition  $P_s$  à  $s$  classes, que l'on agrège en un seul élément  $x$  de masse  $m_t = m_i + m_{i'}$ , produisant la partition  $P_{s-1}$  à  $s - 1$  classes, avec :

$$x = \frac{m_i x_i + m_{i'} x_{i'}}{m_i + m_{i'}}$$

$x$  est le centre de gravité de  $x_i$  et  $x_{i'}$ .

On peut décomposer l'inertie  $I_{ii'}$  de  $x_i$  et  $x_{i'}$  par rapport à  $g$  suivant la relation de Huygens :

$$I_{ii'} = m_i \|x_i - g\|^2 + m_{i'} \|x_{i'} - g\|^2 = m_i \|x_i - x\|^2 + m_{i'} \|x_{i'} - x\|^2 + m_t \|x - g\|^2$$

Seul le dernier terme subsiste si  $x_i$  et  $x_{i'}$  sont remplacés par leur centre de gravité  $x$ . La perte d'inertie inter-classes  $\Delta I_{ii'}$  due au passage de la partition à  $s$  classes à la partition à  $s - 1$  classes équivaut à :

$$\Delta S = \Delta I_{ii'} = I_{inter}(P_s) - I_{inter}(P_{s-1})$$

et vaut donc :

$$\Delta I_{ii'} = m_i \|x_i - x\|^2 + m_{i'} \|x_{i'} - x\|^2$$

En remplaçant  $x$  par sa valeur en fonction de  $x_i$  et  $x_{i'}$  il vient, tous calculs faits :

$$\Delta I_{ii'} = \frac{m_i m_{i'}}{m_i + m_{i'}} \|x_i - x_{i'}\|^2 = \frac{m_i m_{i'}}{m_i + m_{i'}} d^2(x_i, x_{i'})$$

La stratégie d'agrégation fondée sur le critère de la perte d'inertie minimale, dit critère de Ward généralisé, est donc la suivante : au lieu de chercher les deux éléments les plus proches, on cherchera les éléments  $x_i$  et  $x_{i'}$  correspondant à  $\Delta I_{ii'}$  minimale. Ainsi à chaque étape l'inertie inter-classes augmente de la quantité  $\Delta I_{ii'}$  (et l'inertie intra-classes diminue de cette même quantité). Ceci revient à considérer les  $\Delta I_{ii'}$  comme de nouveaux indices de dissimilarités<sup>1</sup> appelés aussi "indices de niveau".

On vérifie que la somme des indices de niveau dans la hiérarchie est égale à l'inertie totale du nuage  $I$  :

$$\sum_{s=2}^n \Delta S = \sum_{s=2}^n I_{inter}(P_s) - I_{inter}(P_{s-1}) = I \quad [2.2 - 2]$$

<sup>1</sup> Par cette transformation de la matrice des distances, les points les plus légers seront plus facilement agrégés.

Si l'on travaille sur les coordonnées des points, on effectuera les calculs des centres de gravité ( $x$  pour  $x_i$  et  $x_{i'}$ ). Par contre si l'on travaille sur les distances, il est commode de pouvoir calculer les nouvelles distances à partir des anciennes (comme cela était le cas pour les techniques précédentes). Le carré des distances entre un point quelconque  $z$  et le centre de classe  $x$  s'écrit, en fonction des distances à  $x_i$  et  $x_{i'}$  :

$$d^2(x, z) = \frac{1}{m_i + m_{i'}} \left( m_i d^2(x_i, z) + m_{i'} d^2(x_{i'}, z) - \frac{m_i m_{i'}}{m_i + m_{i'}} d^2(x_i, x_{i'}) \right)$$

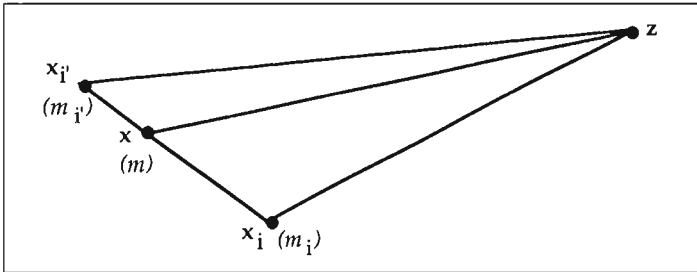


Figure 2.2 - 11  
Théorème de la médiane

Cette formule (théorème de la médiane) s'établit en décomposant l'inertie du doublet  $(x_i, x_{i'})$  par rapport à  $z$  en inertie par rapport à  $x$ , et en inertie de  $x$  par rapport à  $z$  :

$$m_i \|x_i - z\|^2 + m_{i'} \|x_{i'} - z\|^2 = (m_i + m_{i'}) \|x - z\|^2 + \frac{m_i m_{i'}}{m_i + m_{i'}} \|x_i - x_{i'}\|^2$$

L'expression de  $d^2(x, z)$  s'en déduit immédiatement. On réitère le processus sur les éléments restants et le nouvel élément construit par agrégation<sup>1</sup>.

## 2.2.4 Algorithme de recherche en chaîne des voisins réciproques

La principale difficulté dans la construction d'un arbre hiérarchique est le nombre important d'opérations. A chaque étape de l'algorithme est construit un nœud regroupant deux éléments, ce qui nécessite des calculs et des comparaisons de distances entre tous les éléments restant à classer. Le nombre d'opérations à effectuer est de l'ordre de  $n^3$  s'il y a  $n$  objets à classer.

<sup>1</sup> Il existe des variantes de cette méthode qui font appel à des formules de calcul légèrement différentes. On peut par exemple rechercher les classes ayant une inertie interne minimale; on peut aussi utiliser le critère de la variance interne minimale, en désignant par variance l'inertie divisée par la masse. On trouvera des précisions sur ces techniques dans Benzécri (1973).

Les nouveaux algorithmes réunissent à chaque étape non plus deux éléments mais plusieurs couples d'éléments. Ceci réduit considérablement le nombre des opérations qui passe de  $n^3$  à  $n^2$  permettant ainsi la classification de plusieurs milliers d'objets en un temps raisonnable.

Ces algorithmes utilisent le concept de voisins réciproques introduits par McQuitty (1966) : deux éléments  $x_i$  et  $x_{i'}$  sont voisins réciproques si  $x_i$  est le plus proche voisin de  $x_{i'}$  et si  $x_{i'}$  est le plus proche voisin de  $x_i$ .

Ils utilisent également la propriété d'une agrégation hiérarchique selon laquelle, à une étape donnée, deux éléments agrégés pour constituer un nœud sont nécessairement des voisins réciproques (sinon, ils ne constitueraient pas la paire à distance minimale).

Enfin ils utilisent la propriété plus forte (valable seulement si le critère d'agrégation vérifie le critère de la médiane, explicité plus loin) selon laquelle tous les voisins réciproques, à une étape donnée, seront ultérieurement des nœuds de la hiérarchie<sup>1</sup>.

A chaque étape de l'algorithme, au lieu d'agréger seulement les deux plus proches voisins, il y a donc autant de nœuds créés qu'il y a de voisins réciproques. A l'étape finale, tous les éléments sont regroupés en une seule classe et l'arbre est construit.

Le problème de l'algorithme est alors ramené à une recherche efficace des voisins réciproques. Nous allons décrire l'algorithme de cette recherche qui s'effectue en chaîne (Benzécri, 1982c).

## a – Algorithme

Le principe des voisins réciproques peut s'énoncer de la manière suivante : si  $x_i$  est plus proche voisin de  $x_{i'}$  ( $x_i \rightarrow x_{i'}$ ) et si  $x_{i'}$  est plus proche voisin de  $x_i$  ( $x_{i'} \rightarrow x_i$ ) alors  $x_i$  et  $x_{i'}$  sont voisins réciproques ( $x_i \leftrightarrow x_{i'}$ )

Etape 1 : on part d'un objet quelconque  $x_1$  et on cherche son plus proche voisin, noté  $x_2$  puis le plus proche voisin de  $x_2$ , noté  $x_3$ , etc.. On crée ainsi une chaîne d'éléments successifs:

$$x_1 \rightarrow x_2 \rightarrow x_3 \rightarrow \dots \rightarrow x_{i-2} \rightarrow x_{i-1} \rightarrow x_i \rightarrow \dots$$

Une telle chaîne s'arrête nécessairement lorsque deux éléments successifs sont voisins réciproques :

$$\dots \rightarrow x_i \rightarrow \dots \rightarrow x_{k-1} \leftrightarrow x_k$$

La chaîne s'arrêtera ici sur l'élément  $x_k$  si  $x_{k-1}$  est aussi le plus proche voisin de  $x_k$ .  $x_{k-1}$  et  $x_k$  sont voisins réciproques et sont agrégés pour former un nœud.

<sup>1</sup> Le critère de la médiane assure qu'ils resteront toujours voisins réciproques.

Etape 2 : si  $k = 2$  alors la chaîne commence avec un élément qui possède un voisin réciproque:

$$x_1 \leftrightarrow x_2$$

Nous choisissons un nouvel élément à partir duquel une chaîne est construite et qui s'arrête sur de nouveaux voisins réciproques dont l'agrégation fournit un nœud.

Etape 3 : si  $k > 2$ , on continue la recherche des voisins réciproques par extension de la chaîne commençant à l'élément  $x_{k-2}$ . L'algorithme se termine lorsque  $n - 1$  nœuds ont été créés.

### b – Critère de la médiane

Afin de pouvoir utiliser cet algorithme, la chaîne doit pouvoir être prolongée au delà de  $x_{k-2}$  lorsque les voisins réciproques  $x_{k-1}$  et  $x_k$  ont été agrégés. Il est indispensable que cette agrégation ne détruise pas la relation du voisin le plus proche qui existait au préalable entre  $x_{i-1}$  et  $x_i$  avec  $i = 2, 3, \dots, k-2$ . Cette propriété est assurée si le critère d'agrégation utilisé pour construire l'arbre ne crée pas une inversion.

Il n'y a pas inversion si le nœud  $n$ , créé par agrégation de  $a$  et  $b$ , ne peut être plus près d'un quelconque autre élément  $c$  que ne le sont l'élément  $a$  ou l'élément  $b$ . Cette condition<sup>1</sup> dite de "la médiane" s'écrit:

$$\text{si } d(a, b) < \inf \{ d(a, c), d(b, c) \} \text{ alors } \inf \{ d(a, c), d(b, c) \} < d(n, c)$$

Cette propriété est vérifiée par plusieurs critères d'agrégation<sup>2</sup> :

- Saut minimal :  $d(a, b) = \inf \{ d(u, v) \mid u \in a, v \in b \}$
- Saut maximal :  $d(a, b) = \sup \{ d(u, v) \mid u \in a, v \in b \}$
- Distance moyenne :  $d(a, b) = \frac{1}{m_a m_b} \{ \sum_{u \in a} \sum_{v \in b} m_u m_v d(u, v) \}$
- Critère de Ward :  $d(a, b) = \frac{m_a m_b}{m_a + m_b} d(g_a, g_b)$

où  $g_a$  et  $g_b$  sont les centres de gravité des groupes  $a$  et  $b$ .

### 2.2.5 Exemple numérique d'application

L'exemple d'application qui suit comprend deux classifications hiérarchiques effectuées sur les lignes et les colonnes de la table de

<sup>1</sup> Cette condition a été présentée par Bruynooghe (1978) sous le nom d'axiome de réductibilité. Elle permet en effet la mise d'un oeuvre d'un autre algorithme, dit des *voisinages réductibles*, qui permet d'accélérer l'algorithme de base de la classification hiérarchique par l'utilisation de seuils de distances.

<sup>2</sup> On désignera ici à la fois par  $a$  (ou  $b$ ) un élément ou un nœud à une certaine étape de l'agrégation, et l'ensemble des éléments constituant ce nœud.

contingence 1.3 - 10 présentée à la section sur l'analyse des correspondances (cf. § 1.3.8). Les distances entre éléments sont les distances du  $\chi^2$  entre points-profiles et l'agrégation se fait en utilisant le critère de Ward généralisé présenté au paragraphe 2.2.3. Seuls les éléments actifs de l'analyse des correspondances ont été retenus : il s'agit d'une table (8,6) croisant 8 catégories socioprofessionnelle et 6 types de médias, l'unité statistique étant le "contact média".

Comme ce fut le cas pour l'analyse des correspondances de cette même table, la fonction de ce traitement n'est pas la réduction d'un tableau de données trop grand et complexe (fonction principale des techniques d'analyse descriptive multidimensionnelle) mais une présentation pédagogique des différentes étapes de calcul.

### a – Classification des lignes (professions)

Les principales étapes de la classification des lignes sont résumées sur la figure 2.2 - 12, qu'il faut lire de la façon suivante : la première colonne (NUM) donne les numéros des nœuds, qui sont donc des nouveaux éléments à classer et prennent la suite des 8 éléments à classer. La terminologie Ainé et Benjamin (deuxième et troisième colonnes) s'applique aux deux éléments qui sont agrégés à une étape donnée (c'est-à-dire les plus proches à cette étape au sens de l'indice d'agrégation retenu).

CLASSIFICATION HIERARCHIQUE : DESCRIPTION DES 7 NOEUDS (de 9 à 15)						
NUM.	AINÉ	BENJ	EFF.	POIDS	INDICE	HISTOGRAMME DES INDICES DE NIVEAU
9	6	7	2	1927	.00024	*
10	9	5	3	3783	.00038	**
11	2	1	2	789	.00064	****
12	10	4	4	5041	.00208	*****
13	8	11	3	6651	.00276	*****
14	12	13	7	11692	.00493	*****
15	3	14	8	12388	.01125	*****
SOMME	DES	INDICES	=	.02228		

**Figure 2.2 - 12**  
Description des étapes de la classification hiérarchique  
(lignes actives de la table de contingence 1.3 - 10, section 1.3)

On lit ainsi sur la première ligne que le nœud n°9 est formé des éléments terminaux 6 et 7, il est donc formé de 2 éléments (colonne : EFF.) dont le poids total (colonne POIDS) est de 1927. La valeur de l'indice d'agrégation correspondant est de 0.00024. Les valeurs croissantes de l'indice seront illustrées par une esquisse d'histogramme à droite des colonnes numériques<sup>1</sup>. On vérifie que la somme des indices est égale à la somme des valeurs propres issues de l'analyse des correspondances de la même table (tableau 1.3 - 11 du § 1.3.8).

<sup>1</sup> Comme l'indiquait la figure 2.2 - 3, ces histogrammes peuvent donner une idée du nombre de classes d'une bonne partition, qui correspond à un saut important de l'indice.

Le dendrogramme de la figure 2.2 - 13 donne en fait la même information, présentée de façon plus suggestive, car la composition des nœuds à partir des éléments terminaux est maintenant lisible. On note la grande homogénéité des ouvriers (N.Q. et Qual.) et employés (indice très bas), les agriculteurs, petits patrons et inactifs constituant un deuxième groupe moins homogène, alors que les professions intermédiaires occupent une position médiane. Enfin les cadres supérieurs et professions libérales ne se rattachent à l'ensemble des autres catégories que beaucoup plus tard.

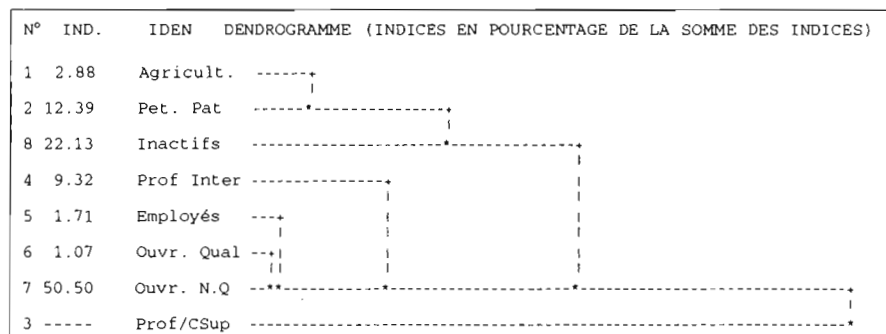


Figure 2.2 - 13  
Dendrogramme

(lignes actives de la table de contingence 1.3 - 10, section 1.3)

On retrouve donc les regroupements visibles sur la figure 1.3 - 23 (section 1.3.8)<sup>1</sup>. Notons ici que le plus grand indice correspond au premier facteur de l'analyse de la section 1.3.8 (opposition des cadres supérieurs et de l'ensemble des catégories), et que le second plus grand indice correspond au second facteur (opposition entre les deux groupes ouvriers/employés et agriculteurs/petits patrons). Cette correspondance entre nœuds et facteurs n'est pas générale, mais fréquente<sup>2</sup>.

## b – Classification des colonnes (médias)

La méthode d'agrégation est la même et conduit évidemment à la même somme des indices (inertie totale). Les règles de lectures des figures 2.2 - 14 et 2.2 - 15 sont les mêmes que précédemment.

Les deux plus grands indices correspondent encore aux principales oppositions visibles sur les deux premiers facteurs de l'analyse des correspondances.

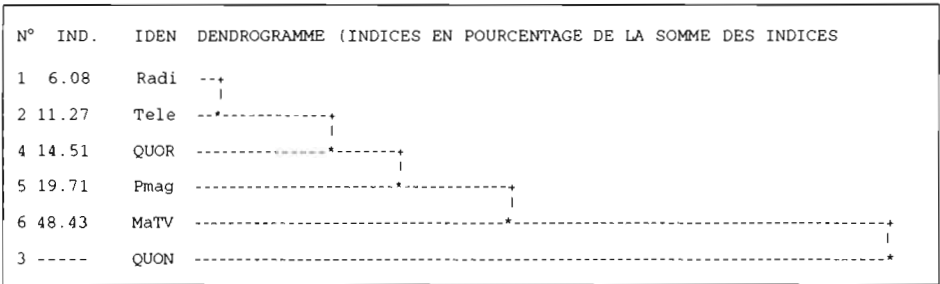
<sup>1</sup> La complémentarité entre les deux approches sera développée section 2.4.

<sup>2</sup> On note également que les deux plus grands indices (0.0112, 0.0049) sont ici inférieurs aux deux plus grandes valeurs propres (0.0139, 0.0072). La section 2.4 précisera quelques relations et inégalités entre ces grandeurs.

La structure observable sur le dendrogramme est celle d'un effet de chaîne, ou de classe absorbante : l'agrégation se fait en ajoutant un élément terminal à la classe de l'étape précédente. Il ne s'agit pas d'un artefact de la méthode<sup>1</sup>. Cela traduit plutôt les diffusions très inégales des différents médias considérés.

CLASSIFICATION HIERARCHIQUE : DESCRIPTION DES 5 NOEUDS (de 7 à 11)						
NUM.	AINE	BENJ	EFF.	POIDS	INDICE	HISTOGRAMME DES INDICES DE NIVEAU
7	2	1	2	7266	.00135	***
8	4	7	3	8933	.00251	*****
9	5	8	4	10236	.00323	*****
10	6	9	5	11950	.00439	*****
11	3	10	6	12388	.01079	*****
SOMME DES INDICES				=	.02228	

**Figure 2.2 - 14**  
Description des étapes de la classification hiérarchique  
(colonnes de la table de contingence 1.3 - 10, section 1.3)



**Figure 2.2 - 15**  
Dendrogramme  
(colonnes de la table de contingence 1.3 - 10, section 1.3)

Notons que si la classification apporte (dans le cas de tableaux en vraie grandeur) certaines informations supplémentaires par rapport à l'analyse des correspondances (les distances sont ici calculées dans tout l'espace), l'absence de représentation simultanée des lignes et des colonnes limite cependant les possibilités d'interprétation.

<sup>1</sup> Contrairement à l'agrégation suivant le saut minimal, le critère de Ward généralisé ne provoque pas facilement d'effets de chaîne.

---

## Classification mixte et description statistique des classes

### 2.3.1 Stratégie de classification mixte

Les algorithmes de classification sont plus ou moins bien adaptés à la gestion d'un nombre important d'objets à classer. La méthode d'agrégation autour des centres mobiles offre des avantages incontestables puisqu'elle permet d'obtenir une partition sur un ensemble volumineux de données à un faible coût, mais elle présente l'inconvénient de produire des partitions dépendant des premiers centres choisis et celui de fixer a priori le nombre de classes. Au contraire, la classification hiérarchique est une famille d'algorithmes que l'on peut qualifier de "déterministes" (i.e. qui donnent toujours les mêmes résultats à partir des mêmes données). De plus, ces algorithmes donnent des indications sur le nombre de classes à retenir mais sont mal adaptés aux vastes recueils de données.

La classification autour des centres mobiles peut en fait être utilisée comme auxiliaire d'autres méthodes de classification. En fournissant des partitions de vastes ensembles de données, elle permet de réduire la dimension de l'ensemble des éléments à classer en opérant des regroupements préalables. De ce fait, l'algorithme de classification qui paraît actuellement bien adapté au partitionnement d'un ensemble comprenant des milliers ou des dizaines de milliers d'individus est un *algorithme mixte*. L'idée repose sur la combinaison des deux techniques de classification présentées précédemment. Cette idée, qui relève du bon sens, a été mise en œuvre spontanément par de nombreux praticiens ; elle se trouve, par exemple, sous le nom de *hybrid clustering* dans Wong (1982).

#### a – Les étapes de l'algorithme

L'algorithme de *classification mixte* procède en trois phases : l'ensemble des éléments à classer subit un partitionnement initial (centres mobiles) de façon à obtenir quelques dizaines, voire quelques centaines de groupes homogènes ; on procède ensuite à une agrégation hiérarchique de ces groupes, dont le dendrogramme suggérera éventuellement le nombre de classes finales à retenir ; et enfin, on optimise (encore par la technique des centres mobiles) la ou les partitions correspondant aux coupures choisies de l'arbre.

La figure 2.3 - 1 schématise les différentes étapes de l'algorithme de classification mixte.



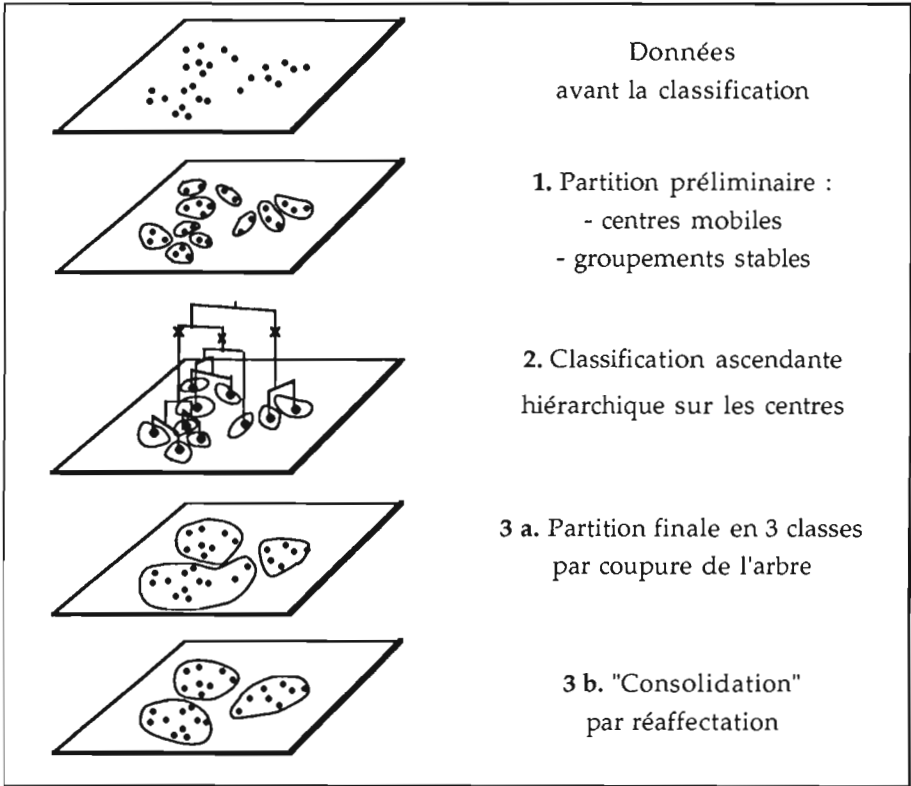


Figure 2.3 - 1  
Schéma de la classification mixte

### 1 - Partitionnement initial

Cette première étape vise à obtenir, rapidement et à un faible coût, une partition des  $n$  objets en  $k$  classes homogènes, où  $k$  est largement plus élevé que le nombre  $s$  de classes désiré dans la population, et largement plus petit que  $n$ . Nous utilisons, pour ce partitionnement initial en quelques dizaines de classes, l'algorithme d'agrégation autour de centres mobiles. Cette procédure augmente l'inertie entre les classes à chaque itération et produit une partition en un nombre fixé au préalable de classes mais qui dépend du choix initial des centres. L'optimalité ne peut être atteinte, mais la partition obtenue peut être améliorée à partir de groupements stables (section 2.1). Ces groupes d'individus ou d'éléments qui apparaissent toujours dans les mêmes classes seront les éléments de base de l'étape suivante.

### 2 - Agrégation hiérarchique des classes obtenues

La seconde étape consiste à effectuer une classification ascendante hiérarchique où les éléments terminaux de l'arbre sont les  $k$  classes de

la partition initiale. Quelques uns de ces groupements peuvent être proches les uns des autres. Ils correspondent à un groupe "réel" qui aurait été coupé artificiellement par l'étape précédente. D'autre part, la procédure crée, en général, plusieurs petits groupes ne contenant parfois qu'un seul élément. Le but de l'étape d'agrégation hiérarchique est de reconstituer les classes qui ont été fragmentées et d'agréger des éléments apparemment dispersés autour de leurs centres d'origine. L'arbre correspondant est construit selon le critère de Ward qui tient compte des masses au moment des choix des éléments à agréger.

### 3 - Partitions finales

La partition finale de la population est définie par coupure de l'arbre de la classification ascendante hiérarchique. L'homogénéité des classes obtenues peut être optimisée par réaffectations.

#### b – Choix du nombre de classes par coupure de l'arbre

Le choix du niveau de la coupure, et ainsi du nombre de classes de la partition, peut être facilité par une inspection visuelle de l'arbre (cf. figures 2.3 - 1 et 2.3 - 2) : la coupure doit être faite après les agrégations correspondant à des valeurs peu élevées de l'indice, qui regroupent les éléments les plus proches les uns des autres, et avant les agrégations correspondant à des valeurs élevées de l'indice, qui dissocient les groupes bien distincts dans la population.

D'une manière générale, plus on agrège des éléments, autrement dit plus on se rapproche du sommet de l'arbre, plus la distance entre les deux classes les plus proches est grande et plus l'indice de niveau est élevé. En coupant l'arbre au niveau d'un saut important de cet indice, on peut espérer obtenir une partition de bonne qualité, car les individus regroupés auparavant étaient proches, et ceux regroupés après la coupure sont nécessairement éloignés, ce qui est la définition d'une bonne partition.

En pratique, la situation n'est pas aussi clairement définie que le montre la figure 2.3 - 2. L'utilisateur pourra choisir entre deux ou trois niveaux de coupure possibles et donc entre deux ou trois partitions finales.

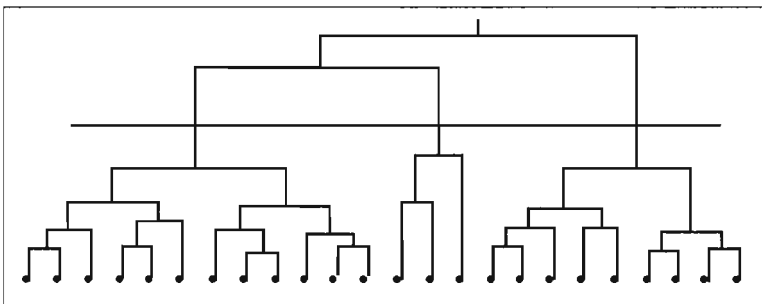


Figure 2.3 - 2  
Coupure visuelle de l'arbre

La coupure de l'arbre peut être facilitée par l'examen de l'histogramme des indices croissants de niveau et l'on coupera au niveau pour lequel cet histogramme marque un palier important. Toute barre de cet histogramme indique la valeur de l'indice d'une agrégation c'est-à-dire la perte d'inertie obtenue en passant d'une partition en  $s$  classes à la partition en  $s - 1$  classes.

La situation idéale est montrée par la figure 2.3 - 3.a où l'on observe un palier évident entre le 4<sup>ème</sup> et le 5<sup>ème</sup> indices suggérant ainsi une bonne partition en cinq classes. La figure 2.3 - 3.b est typique de la situation où il est difficile de décider d'un nombre "réel" de groupes dans la population. Mais une telle partition, en  $s$  classes par exemple, n'est pas la meilleure possible, car l'algorithme de classification hiérarchique n'a pas la propriété de donner à chaque étape une partition optimale.

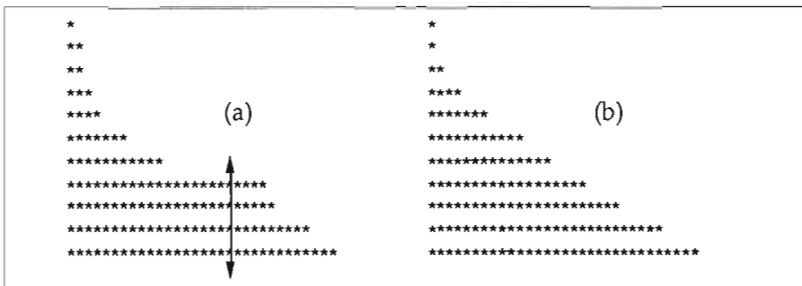


Figure 2.3 - 3  
Histogrammes des indices de niveau

Cela tient en particulier à la contrainte imposée pour la construction de l'arbre : la partition en  $s$  classes contient la partition en  $s - 1$  classes et est contenue dans la partition en  $s + 1$  classes, l'algorithme ne remettant jamais en cause une agrégation effectuée. On peut donc améliorer encore la partition obtenue si on s'affranchit de cette contrainte.

### c – Procédure de consolidation

Pour améliorer la partition obtenue, on utilise de nouveau une procédure d'agrégation autour des centres mobiles dont on sait qu'elle ne peut qu'augmenter l'inertie entre les classes à chaque itération.

Au départ, les centres de classes sont ceux obtenus par coupure de l'arbre. A la première itération, on affecte les éléments à leur centre de gravité le plus proche, ce qui crée de nouvelles classes dont on calcule les centres. A la deuxième itération et aux suivantes, on réaffecte les éléments à leur centre le plus proche. Après un certain nombre d'itérations, il n'y a plus de réaffectation à opérer et le calcul est arrêté. En pratique, la procédure est arrêtée dès que l'inertie entre les classes cesse de croître de façon sensible.

Cette procédure de consolidation a pour effet d'optimiser, par réaffectation, la partition obtenue par coupure de l'arbre hiérarchique. Malgré la relative complexité de la procédure, on ne peut toujours pas être assuré d'avoir

trouvé la "meilleure partition en  $k$  classes" mais on s'en approche vraisemblablement dans beaucoup de situations courantes.

### 2.3.2 Description statistique des classes

Les éléments d'une même classe se ressemblent vis-à-vis de l'ensemble des critères choisis pour les décrire. Il reste maintenant à préciser quels sont les critères qui sont à l'origine des regroupements observés. On procède à une description automatique des classes qui constitue en pratique une indispensable étape de toute procédure de classification<sup>1</sup>.

Les aides à l'interprétation des classes sont généralement fondées sur des comparaisons de moyennes ou de pourcentages à l'intérieur des classes avec les moyennes ou les pourcentages obtenus sur l'ensemble des éléments à classer<sup>2</sup>. Pour sélectionner les variables continues ou les modalités des variables nominales les plus caractéristiques de chaque classe, on mesure l'écart entre les valeurs relatives à la classe et les valeurs globales. Ces statistiques peuvent être converties en un critère appelé *valeur-test* permettant d'opérer un tri sur les variables, et de désigner ainsi les variables les plus caractéristiques (cf. Morineau, 1984).

Parmi les variables figurent également celles qui n'ont pas contribué à la construction des classes mais qui peuvent participer à leur description sur le même principe que les variables supplémentaires dans une analyse factorielle. Ces variables permettent *a posteriori* d'identifier et de caractériser les regroupements établis à partir des variables actives.

#### a – Valeurs-test pour les variables continues

Pour caractériser une classe par les variables continues, on compare  $\bar{X}_k$ , la moyenne d'une variable  $X$  dans la classe  $k$ , à la moyenne générale  $\bar{X}$  et on évalue l'écart en tenant compte de la variance  $s_k^2(X)$  de cette variable dans la classe. La valeur-test est ici simplement la quantité :

$$t_k(X) = \frac{\bar{X}_k - \bar{X}}{s_k(X)}$$

avec :

$$s_k^2(X) = \frac{n - n_k}{n - 1} \frac{s^2(X)}{n_k}$$

<sup>1</sup> Il existe également des possibilités purement graphiques de représentation des classes (graphiques de densité, arbres, dendrogrammes) qui peuvent compléter les descriptions statistiques de ce paragraphe. Sur ce point, cf. Ohsumi (1988).

<sup>2</sup> Ce qui est équivalent à comparer moyennes (ou pourcentages) dans la classe et hors de la classe.

où  $s^2(X)$  est la variance empirique de la variable  $X$ . On reconnaît en  $s_k^2(X)$  la variance d'une moyenne dans le cas d'un tirage sans remise des  $k$  éléments concernés.

### *Interprétation en termes de probabilités (variables supplémentaires)*

Sous l'hypothèse "nulle" d'un tirage au hasard sans remise des  $n_k$  individus de la classe  $k$ , la variable  $\bar{X}_k$  représentant la moyenne dans la classe a pour espérance  $\bar{X}$  et pour variance théorique  $s_k^2(X)$ .

La valeur-test  $t_k(X)$  suit donc approximativement une loi de Laplace-Gauss centrée et réduite (théorème de la limite centrale). Elle évalue la distance entre la moyenne dans la classe et la moyenne générale en nombre d'écart-types d'une loi normale.

Il va de soi que cette interprétation n'a de sens que pour une variable  $X$  supplémentaire, c'est-à-dire n'ayant pas participé à la construction des classes. On ne peut en effet stipuler une indépendance entre les classes d'une partition et une des variables ayant servi à définir cette partition.

On calcule ensuite la probabilité que la variable dépasse la valeur absolue de la différence observée. Plus la valeur-test est forte (plus la probabilité est faible), plus l'hypothèse d'avoir les  $n$  valeurs de la variable  $X$  tirées au hasard parmi les valeurs possibles est discutable. Dans ce cas, la moyenne dans la classe diffère de la moyenne générale, et la variable est caractéristique de la classe. Le classement des variables par probabilités de dépassement croissantes est le même que le classement par valeurs-test décroissantes. Du point de vue de la désignation des variables les plus caractéristiques, les deux informations sont équivalentes.

### *Extension aux variables actives*

S'il n'est pas licite d'interpréter de façon probabiliste les valeurs-test calculées sur les variables actives, il est possible de les utiliser pour obtenir un *classement* de celles-ci en vue de caractériser chaque classe. Les valeurs absolues des valeurs-test constituent alors de simples mesures de similarité entre variables et classes.

## **b – Valeurs-test pour les variables nominales**

Une modalité (ou catégorie) d'une variable nominale est considérée comme caractéristique de la classe si son *abondance* dans la classe est jugée significativement supérieure à ce qu'on peut attendre compte tenu de sa présence dans la population. En notant  $n_{kj}$  le nombre d'individus ayant la modalité  $j$  parmi les  $n_k$  individus de la classe  $k$ ,  $n_j$  le nombre d'individus ayant la modalité  $j$  et  $n$  l'effectif total, l'abondance de la modalité  $j$  est définie, en premier lieu, en comparant son pourcentage dans la  $k$  ème classe :

$$\frac{n_{kj}}{n_k} \text{ à son pourcentage dans la population } \frac{n_j}{n}.$$

La valeur-test prend en compte tous les éléments du tableau 2.3 - 1.

**Tableau 2.3 - 1**  
**Modalités de variables nominales et classes d'individus**

	classe k	autres classes	population
modalité j	$n_{kj}$	*	$n_j$
autres modalités	*	*	*
population	$n_k$	*	$n$

Sous l'hypothèse "nulle"<sup>1</sup> où les  $n_k$  individus de la classe  $k$  sont tirés au hasard sans remise parmi la population des  $n$  individus, le pourcentage d'individus de la classe  $k$  ayant la modalité  $j$  d'une part, et le pourcentage d'individus ayant la modalité  $j$  dans la population d'autre part, devraient coïncider aux fluctuations aléatoires près :

$$\frac{n_{kj}}{n_k} = \frac{n_j}{n}$$

C'est l'hypothèse d'indépendance sous laquelle le nombre  $N$  d'individus de la classe  $k$  ayant la modalité  $j$  est une variable aléatoire qui suit une loi hypergéométrique dont les trois paramètres apparaissent dans les marges du tableau 2.3 - 1. On calcule donc la probabilité d'obtenir une valeur  $N$  supérieure à  $n_{kj}$  :

$$p_k(j) = Prob(N \geq n_{kj})$$

Plus cette probabilité<sup>2</sup>  $p_k(j)$  est faible, plus l'hypothèse d'un tirage au hasard est difficile à accepter. On se sert de cette probabilité pour ranger les modalités caractéristiques de la classe (la plus caractéristique correspondant à la plus petite probabilité).

Cette probabilité est souvent très faible. Il est commode de lui substituer la valeur  $t_k(N)$  de la variable de Laplace-Gauss correspondant à la même probabilité. C'est la *valeur-test*. Elle mesure l'écart entre la proportion dans la classe et la proportion générale, en nombre d'écart-types d'une loi normale. La valeur-test, pour une modalité d'une variable nominale, est

<sup>1</sup> Comme dans le cas des variables continues, cette hypothèse nulle n'a de sens que pour des variables nominales supplémentaires. Mais les valeurs-test que l'on va calculer pourront encore jouer le rôle d'indices de similarités entre modalités actives et classes et donc servir à ranger ces modalités par ordre d'intérêt pour chaque classe.

<sup>2</sup> Si l'on désigne par  $C_a^b$  le nombre de parties distinctes de  $b$  éléments que l'on peut extraire d'un ensemble de  $a$  éléments, la probabilité  $Prob(N = x)$  s'écrit ici :

$$Prob(N = x) = \frac{C_{n_j}^x C_{n-n_j}^{n_k-x}}{C_n^{n_k}} \text{ et la probabilité } p_k(j) \text{ vaut alors : } p_k(j) = \sum_{x=n_{kj}}^{x=n_k} Prob(N = x).$$

donc un critère statistique associé à la comparaison des effectifs dans le cadre d'une loi hypergéométrique<sup>1</sup>.

### c – Variables caractéristiques d'une classe

La valeur-test revient à effectuer un changement de mesure en transformant la probabilité d'une distribution quelconque en nombre d'écart-types d'une loi normale centrée réduite.

Que ce soit pour la recherche des variables continues ou des modalités des variables nominales caractéristiques d'une classe, la valeur absolue de la valeur-test est l'analogie de la valeur absolue d'une variable normale centrée réduite<sup>2</sup>.

Les variables sont d'autant plus intéressantes que les valeurs-test associées sont fortes en valeur absolue. On peut alors ranger ces variables suivant les valeurs-test décroissantes et ne retenir que les éléments les plus significatifs, ce qui permet de caractériser très rapidement les classes.

En sélectionnant, pour chaque classe, les variables les plus caractéristiques, et en calculant leur moyenne ou leur pourcentage dans la classe, on constitue ainsi le "profil-type" de la classe. Rappelons que la valeur-test ne correspond à un vrai test d'hypothèse<sup>3</sup> que si la variable à laquelle elle est associée est supplémentaire.

Mentionnons enfin, comme cela a été fait à la section 1.2 à propos de l'analyse en composantes principales, que le fait de calculer simultanément plusieurs valeurs-test met l'utilisateur dans une situation de "comparaisons multiples", qui impose de prendre des seuils de signification plus sévères que ceux mis en œuvre lors d'un test unique.

<sup>1</sup> Notons qu'une estimation approchée de la valeur-test peut être obtenue de façon plus simple en prenant en compte l'espérance de  $N$  :

$$E(N) = n_k \frac{n_j}{n} \quad \text{et la variance de } X \quad s_k^2(N) = n_k \frac{n - n_k}{n - 1} \frac{n_j}{n} \left( 1 - \frac{n_j}{n} \right),$$

et en calculant la quantité  $t_k(N) = \frac{N - E(N)}{s_k(N)}$  qui donne directement la variable centrée,

réduite et normale si l'on peut appliquer l'approximation de Laplace-Gauss de la loi hypergéométrique. Cette approximation est suffisante dans les applications qui ne mettent pas en jeu des effectifs faibles.

<sup>2</sup> Dans le cadre de tests classiques, on dira qu'elle est significative au seuil usuel 5% si elle dépasse la valeur 1,96 : l'hypothèse "nulle" est rejetée et la moyenne ou la proportion d'une variable sur la population globale et celle dans la classe diffèrent significativement.

<sup>3</sup> Ici on a l'hypothèse qu'une variable continue ou une modalité d'une variable nominale est indépendante de la partition.

---

## Complémentarité entre analyse factorielle et classification

Les méthodes factorielles (notamment l'analyse des correspondances multiples) sont particulièrement bien adaptées à l'exploration de grands tableaux de données individuelles tels que ceux produits par les enquêtes. Mais elles ne suffisent pas toujours à fournir une vue satisfaisante de l'ensemble des données. Non seulement les visualisations ne véhiculent qu'une partie de l'information, mais elles sont parfois elles-mêmes trop complexes pour être interprétées facilement.

Dans ces circonstances, les techniques de classification peuvent compléter et nuancer les résultats des analyses factorielles. La complémentarité entre analyse factorielle et classification concerne la compréhension de la structure des données et celle des aides pratiques dans la phase d'interprétation des résultats.

Dans une première partie, on justifiera cette utilisation conjointe du point de vue de l'utilisateur confronté à un ensemble complexe de données. Puis on examinera quelques aspects techniques et théoriques de cette complémentarité.

### 2.4.1 Utilisation conjointe de l'analyse factorielle et de la classification

Face à de très grands tableaux de données, il est indispensable de disposer d'une vue d'ensemble de la base d'information. De ce point de vue, les méthodes factorielles sont certainement les techniques exploratoires les mieux adaptées.

#### a – Nécessité... et insuffisance des méthodes factorielles

Mais, les représentations graphiques issues des méthodes factorielles présentent certains inconvénients, dont certains sont d'ailleurs interdépendants :

##### - 1- Difficultés d'interprétation

Il est toujours difficile d'interpréter les axes ou plans factoriels au delà du plan principal. Le plan (3,4), engendré par les axes factoriels 3 et 4, décrit des proximités qui sont des termes correctifs par rapport aux proximités principales observées sur les deux premiers axes. L'interprétation de ces proximités est donc assez délicate.



- *2- Compression excessive et déformations*

Les visualisations sont limitées à deux, ou en général à très peu de dimensions, alors que le nombre d'axes "significatifs" peut être bien supérieur. Cette compression excessive de l'espace peut entraîner des distorsions fâcheuses et des superpositions de points occupant des positions distinctes dans l'espace.

*3 - Manque de robustesse*

Les visualisations peuvent manquer de robustesse. Un point-profil aberrant peut notablement influencer le premier facteur et par là toutes les dimensions suivantes, puisque ces dimensions sont reliées au premier axe à travers la contrainte d'orthogonalité des axes.

*4 - Graphiques factoriels inextricables*

Les visualisations peuvent concerner des centaines de points et donner lieu à des graphiques chargés ou illisibles.

Pour remédier à ces lacunes, montrons, point par point, quels peuvent être les apports d'une classification menée simultanément.

- *Difficultés d'interprétation et compression excessive des données (points 1 et 2) :*

On complète l'analyse factorielle par une classification réalisée sur l'espace tout entier ou sur un sous-espace défini par les premiers facteurs les plus significatifs. Les classes prennent en compte la dimension réelle du nuage de points. Elles corrigent donc certaines déformations dues à l'opération de projection.

Une classe peut aussi être typique d'un axe de rang élevé et aider à l'interprétation de ce sous-espace particulier difficilement observable autrement.

- *Robustesse imparfaite (point 3) :*

La plupart des algorithmes de classification, et particulièrement les algorithmes d'agglomération, sont localement robustes au sens où les parties basses des dendrogrammes produits (nœuds correspondant aux plus petites distances) sont indépendantes des éventuels points marginaux isolés.

- *Allègement et description automatique des sorties graphiques (point 4) :*

Lorsqu'il y a trop de points-individus sur un plan factoriel, il paraît utile de procéder à des regroupements d'individus en familles homogènes. Il faut donc à ce stade faire appel aux capacités de gestion et de calcul de l'ordinateur pour compléter, aider et clarifier la présentation des résultats. Les classes peuvent être utilisées pour aider l'interprétation des plans factoriels en identifiant des zones bien décrites. Il est en effet plus

facile de décrire des classes qu'un espace continu, même à deux dimensions. La notion de classe est élémentaire et accessible à l'intuition. Les descriptions de ces classes peuvent être fondées sur d'élémentaires comparaisons de moyennes ou de pourcentages. Les nombreux points sont ainsi remplacés par quelques centres de gravité de classes. Comme les algorithmes utilisés pour ces regroupements fonctionnent de la même façon que les points soient situés dans un espace à deux ou à dix dimensions, on allège les sorties graphiques tout en améliorant la qualité de la représentation (points 1 et 2 ci-dessus).

*Mais les méthodes factorielles sont nécessaires, malgré leurs insuffisances : la faculté descriptive des axes, les descriptions sous forme de continuum géométrique restent irremplaçables.*

La classification ne réussit pas toujours à montrer l'importance de certaines tendances ou de facteurs latents continus. Pour observer l'organisation spatiale des classes, le positionnement des classes sur les axes factoriels s'avère indispensable. La classification peut évidemment aider à découvrir l'existence de groupes d'individus. L'analyse factorielle peut mettre en avant des facteurs latents inattendus. La découverte de tels phénomènes ou dimensions cachées est l'objectif de ces deux familles de méthodes et certainement le plus ambitieux. Leur utilisation complémentaire est souvent indispensable pour atteindre cet objectif.

### **b – Mise en œuvre pratique dans le cas de la classification mixte**

Pour décrire un ensemble de données de grande taille, principale circonstance dans laquelle l'usage complémentaire des techniques factorielles et de classification est utile, la mise en œuvre conjointe de ces techniques s'opère de la façon suivante.

#### **- Étape 1 : L'analyse factorielle**

L'analyse factorielle est utilisée comme une étape préalable à la classification pour deux raisons : pour son pouvoir de description, présenté dans les chapitres précédents, et pour son pouvoir de filtrage, qui permettra éventuellement de travailler sur des coordonnées factorielles moins nombreuses que les variables de départ.

#### **- Étape 2 : Classification à partir des facteurs**

Il est équivalent d'effectuer une classification des individus sur un ensemble de  $p$  variables ou sur l'ensemble des  $p$  facteurs. Mais on peut aussi ne prendre en compte qu'un sous-espace factoriel de dimension  $q$  ( $q < p$ ) et réaliser une classification sur les  $q$  premiers axes. Cela présente l'avantage d'éliminer des fluctuations aléatoires qui constituent en général l'essentiel de la variance recueillie dans les directions des  $p - q$  derniers axes (variations non systématiques contenues dans les données). Le fait d'abandonner les derniers facteurs revient à effectuer une sorte de "lissage" des données, ce qui en général

améliore la partition en produisant des classes plus homogènes. Les distances entre points sont calculées dans l'espace des premiers axes factoriels avec la distance euclidienne usuelle. Le calcul est simple et la classification peut être menée sur des grands ensembles d'individus<sup>1</sup>. La difficulté réside parfois dans le choix du nombre d'axes à retenir (cf. § 4.2.3).

### Étape 3 : Description automatique des classes

Une fois les individus regroupés en classes, on a vu (§ 2.3.2) qu'il est facile d'obtenir une description automatique de ces classes. On calcule, pour les variables numériques comme pour les variables nominales, des statistiques d'écart entre les valeurs internes à la classe et les valeurs globales. Les valeurs-test permettent de les ranger par ordre d'intérêt.

### Étape 4 : Positionnement des classes dans le plan factoriel

La division en classes opère un découpage plus ou moins arbitraire d'un espace continu. L'analyse en axe principaux préalable permet alors de visualiser les positions relatives des classes dans l'espace et peut mettre en évidence certaines "trajectoires" masquées par la discontinuité des classes. Il est intéressant de projeter les centres de gravité des classes au sein des variables ou des modalités actives sur le premier plan factoriel (figure 2.4 - 1).

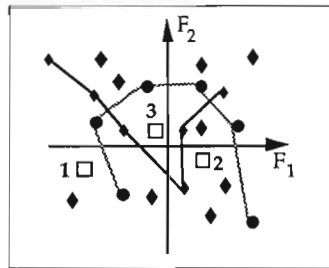


Figure 2.4 - 1  
Positions relatives des classes dans l'espace factoriel

Le support visuel permet d'apprécier les distances entre les classes. Par ailleurs, la position de chaque individu repéré par le numéro de sa classe permet de représenter la densité et la dispersion des classes dans le plan (cf. figure 2.4 - 2).

<sup>1</sup> Une technique de classification hiérarchique tel que l'algorithme des voisins réciproques (et particulièrement l'algorithme de recherche en chaîne) peut être réalisée sans garder la matrice des distances en mémoire centrale. Les distances entre couples de points sont recalculées à la demande dans l'espace réduit des  $q$  premiers facteurs. La mise en mémoire de la matrice  $(n, q)$  construite à partir des  $q$  principales coordonnées des  $n$  observations est souvent beaucoup moins encombrante que le tableau des  $n(n-1)$  distances.

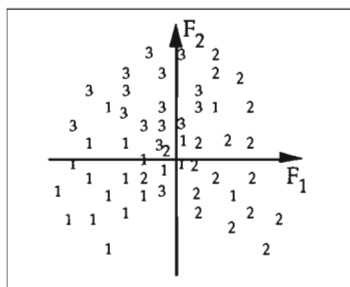


Figure 2.4 - 2  
Densité et dispersion des classes dans l'espace factoriel

L'utilisation conjointe de l'analyse factorielle et de la classification permet de se prononcer non seulement sur la réalité des classes, mais également sur leurs positions relatives, leur forme, leur densité et leur dispersion. Les deux techniques se valident mutuellement.

### c – Autres aspects de la complémentarité

A propos des liens entre les méthodes d'analyse par axes principaux et les méthodes de classification, il faudrait évoquer des méthodes que l'on peut qualifier d'hybrides, c'est-à-dire qui produisent simultanément des axes et des classes. Ainsi, le lien existant entre le haut de l'arbre et les premiers axes factoriels peut suggérer d'utiliser ceux-ci pour construire un arbre à partir des plus grands indices (classification descendante ou divisive, cf. par exemple Reinert, 1986). On peut également chercher des axes principaux susceptibles de représenter au mieux une classification (Art et *al.*, 1982 ; Gnanadesikan et *al.*, 1982). Certaines de ces méthodes (projections révélatrices, analyses de contiguïté) seront brièvement présentées au § 3.7.4 du chapitre 3. Dans un autre esprit, van Buuren et Heiser (1989), pour classer des individus décrits par des variables nominales, cherchent simultanément des classes et un codage des variables qui optimise un critère de qualité de la classification.

## 2.4.2 Aspects techniques et théoriques de la complémentarité

La complémentarité entre l'analyse des correspondances et la classification ascendante hiérarchique présente des avantages pratiques pour l'utilisateur. On examinera dans ce paragraphe certains aspects plus techniques de cette complémentarité.

### a – Classification des lignes ou colonnes d'un tableau de contingence

La classification ascendante hiérarchique agrège des groupes d'éléments suivant différents critères d'agrégation. Parmi ceux-ci, le critère de Ward

généralisé apparaît compatible avec l'analyse des correspondances puisqu'il est fondé sur une notion d'inertie similaire. On a montré en particulier (cf. § 2.2.3.b formule [2.2 - 2]) que la somme des valeurs propres (inertie totale du nuage) est égale la somme des indices de niveau. Aussi, malgré des classes sphériques que ce critère a tendance à produire, il y a une certaine cohérence à utiliser le critère d'inertie de Ward sur un tableau de coordonnées factorielles elles-mêmes issues d'un calcul d'inertie. Si l'arbre de la classification est construit sur les  $q$  premiers axes factoriels, on vérifiera que la somme des indices de niveau est égale à la somme des  $q$  premières plus grandes valeurs propres retenues.

Une propriété importante de l'analyse des correspondances va dans le sens d'une bonne compatibilité avec la classification : l'équivalence distributionnelle (cf. § 1.3.2.f et 1.3.3.a) qui garantit la stabilité des résultats quand on regroupe les éléments ayant des profils semblables.

Agréger les lignes et les colonnes d'un tableau de contingence est naturel dans le sens où il s'agit de remplacer des classes par des classes au lieu de remplacer des individus par des groupes d'individus ou des variables par des groupes de variables<sup>1</sup>.

### b – Un exemple de coïncidence entre les deux approches

Considerons la table de contingence  $K_{IJ}$  (tableau 2.4 - 1). Elle a, nous allons le vérifier, la propriété de donner des résultats similaires lorsqu'elle subit une analyse des correspondances et une classification hiérarchique utilisant le critère d'agrégation de Ward (cf. § 2.2.3.b).

Tableau 2.4 - 1  
Table de Contingence  $K_{IJ}$

	COL7	COL2	COL3	COL4	COL5	COL6	COL1	COL8
LIG1	2	18	12	12	2	2	30	2
LIG4	2	12	21	27	2	2	12	2
LIG5	14	2	2	2	24	20	2	14
LIG2	2	30	12	12	2	2	18	2
LIG6	14	2	2	2	20	24	2	14
LIG7	23	2	2	2	14	14	2	21
LIG3	2	12	27	21	2	2	12	2
LIG8	21	2	2	2	14	14	2	23

En fait, un réarrangement des lignes et des colonnes montre que cette table n'est pas anodine. Elle contient de forts traits structuraux (tableau 2.4 - 2). Elle est symétrique et semble formée de blocs et de sous-blocs particuliers. Ce réarrangement, on va le voir, est un sous-produit de l'analyse des correspondances.

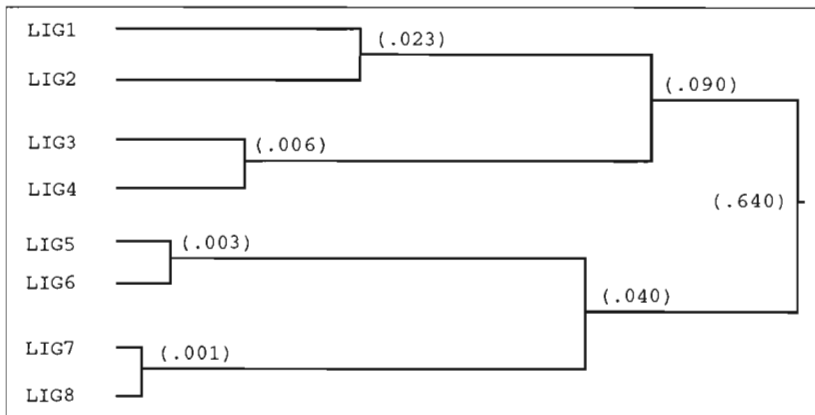
<sup>1</sup> La classification des éléments d'un tableau de contingence fondée sur le regroupement de catégories homogènes a été abordée par Benzécri (1973), Jambu et Lebeaux (1978), Govaert (1984), Cazes (1986), Gilula (1986), Escoufier (1988), Greenacre (1988).

**Tableau 2.4 - 2**  
**Table de Contingence  $K_{IJ}$  réordonnée**

	COL1	COL2	COL3	COL4	COL5	COL6	COL7	COL8
LIG1	30	18	12	12	2	2	2	2
LIG2	18	30	12	12	2	2	2	2
LIG3	12	12	27	21	2	2	2	2
LIG4	12	12	21	27	2	2	2	2
LIG5	2	2	2	2	24	20	14	14
LIG6	2	2	2	2	20	24	14	14
LIG7	2	2	2	2	14	14	23	21
LIG8	2	2	2	2	14	14	21	23

Cette table de contingence fait en fait partie d'une famille plus large de tableaux décrits dans Benzécri (1973, vol. 2, chapitre 11) qui seront brièvement évoqués plus bas.

Une classification ascendante hiérarchique utilisant le critère de Ward produit le dendrogramme représenté sur la figure 2.4 - 3, où les indices de niveaux figurent entre parenthèses près des nœuds correspondants.



**Figure 2.4 - 3**  
**Esquisse du dendrogramme décrivant la classification**  
**hiérarchique de la table de contingence (8,8)  $K_{IJ}$**

Les valeurs propres issues de l'analyse des correspondances de  $K_{IJ}$  figurent dans le tableau 2.4 - 3. Elles coïncident avec les indices d'agrégation.

**Tableau 2.4 - 3**  
**Valeurs propres issues de l'analyse des correspondances de  $K_{IJ}$**

$\lambda_1 =$	.640	(80.0 % de la trace)
$\lambda_2 =$	.090	(11.0 %)
$\lambda_3 =$	.040	( 5.0 %)
$\lambda_4 =$	.023	( 3.0 %)
$\lambda_5 =$	.006	( .7 %)
$\lambda_6 =$	.003	( .4 %)
$\lambda_7 =$	.001	( .1 %)

Le tableau 2.4 - 4 donne les coordonnées factorielles des points lignes (qui sont les mêmes que celles des points colonnes au signe près, puisque la matrice de départ est symétrique). La façon dont sont organisés ces vecteurs propres permet de comprendre le processus de construction de la table de contingence : on part des facteurs structurés de cette façon et on utilise la formule de reconstitution des données.

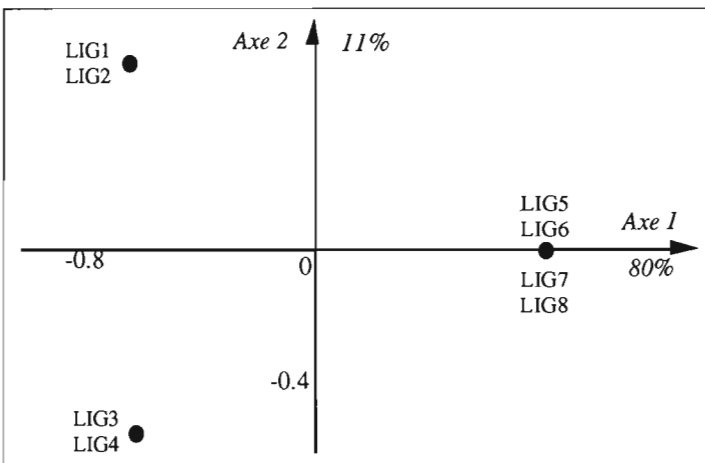
Chaque vecteur oppose deux blocs. Il est orthogonal au vecteur précédent et les coordonnées sont égales à l'intérieur de chaque bloc. Tous les vecteurs sont centrés et orthogonaux à la première bissectrice.

La figure 2.4 - 4 donne la représentation des points-profils dans le plan des deux premiers axes factoriels.

**Tableau 2.4 - 4**  
Coordonnées factorielles issues  
de l'analyse des correspondances de  $K_{ij}$

Axes	1	2	3	4	5	6	7
LIGNE1	-.80	.42	0.00	.30	0.00	0.00	0.00
LIGNE2	-.80	.42	0.00	-.30	0.00	0.00	0.00
LIGNE3	-.80	-.42	0.00	0.00	-.15	0.00	0.00
LIGNE4	-.80	-.42	0.00	0.00	.15	0.00	0.00
LIGNE5	.80	0.00	-.28	0.00	0.00	.10	0.00
LIGNE6	.80	0.00	-.28	0.00	0.00	-.10	0.00
LIGNE7	.80	0.00	.28	0.00	0.00	0.00	.06
LIGNE8	.80	0.00	.28	0.00	0.00	0.00	-.06

On constate que cette figure bi-dimensionnelle permet de distinguer les deux grands blocs (axe 1), puis, à l'intérieur de l'un d'eux, deux sous-blocs (axe 2), mais qu'elle est moins riche d'information que la figure 2.5 - 3, elle aussi bidimensionnelle.



**Figure 2.4 - 4**  
Premier plan factoriel de l'analyse de  $K_{ij}$

La figure 2.4 - 3 du dendrogramme a en effet l'avantage de montrer simultanément tous les blocs et tous les niveaux de la hiérarchie.

Décrivons brièvement ces tableaux de *correspondances hiérarchiques* en renvoyant à Benzécri (1973, op. cit.) et Cazes (1984, 1986 a) pour une présentation systématique et des généralisations de ces notions.

D'une manière générale, dans une *hiérarchie binaire*  $H$  sur un ensemble  $I$  à  $n$  éléments chaque élément non terminal  $h \in H$  peut être partitionné de façon unique en deux éléments  $a(h)$  et  $b(h)$  :

$$h = a(h) \cup b(h) \quad \text{avec } a(h) \in H \text{ et } b(h) \in H$$

On suppose cette hiérarchie indicée (cf. § 2.2.1.c). On suppose également que l'indice  $\lambda(h)$  prend ses valeurs dans  $(0,1)$  et qu'il est nul pour les éléments terminaux. Chaque élément  $i \in I$  est d'autre part muni d'une masse  $p_i$  strictement positive avec :

$$\sum_{i=1}^n p_i = 1$$

Pour chaque nœud  $h$  de la hiérarchie, on peut associer une fonction sur  $I$  à valeurs réelles  $f_h$ , de moyenne nulle, c'est-à-dire telle que :

$$\sum_{i=1}^n p_i f_h(i) = 0$$

Cette fonction est nulle en dehors de  $h$  ( $i \notin h \Rightarrow f_h(i) = 0$ ) et constante sur chacun des deux nœuds  $a(h)$  et  $b(h)$  qui constituent  $h$ .

Ces constantes sont définies par les formules suivantes, en notant  $p_h, p_a$  et  $p_b$  les masses respectives des éléments  $h, a(h)$  et  $b(h)$  :

$$f_h(i) = \sqrt{\frac{p_b}{p_h p_a}} \quad \text{pour } i \in a(h)$$

$$f_h(i) = -\sqrt{\frac{p_a}{p_h p_b}} \quad \text{pour } i \in b(h)$$

Si l'on munit l'espace des fonctions  $f_h$  du produit scalaire :

$$\langle f_{h'}, f_h \rangle = \sum_{i=1}^n p_i f_{h'}(i) f_h(i)$$

On vérifie facilement que les fonctions  $f_h$  sont de norme (ou de variance) 1 et que les  $n-1$  fonctions correspondant aux nœuds de la hiérarchie constituent une base orthonormée de l'ensemble des fonctions sur  $I$ .

La formule de reconstitution des données en analyse des correspondances (cf. § 1.3.3.h) permet alors de générer un tableau de correspondances symétrique  $C$  de terme général  $c_{ii'}$  :

$$c_{ii'} = p_i p_{i'} \left( 1 + \sum_{h=1}^{n-1} \sqrt{l_h} f_h(i) f_h(i') \right)$$



les  $n-1$  nœuds repérés par  $h$  étant supposés numérotés par ordre d'indices d'agrégation  $\lambda_h$  décroissants. La table de contingence  $\mathbf{K}_{JJ}$  ci-dessus a été générée<sup>1</sup> de cette façon.

### 2.4.3 Valeurs propres et indices de niveau

Hormis des cas très particuliers, comme ceux constitués par les correspondances hiérarchiques étudiées au paragraphe précédent, les relations entre analyse des correspondances et classification opérés sur une même table de contingence sont difficiles à étudier.

Dans le cas de la classification hiérarchique utilisant le critère de Ward, on peut mettre en évidence certaines inégalités et étudier certaines structures particulières.

#### a – Quelques inégalités

Notons tout d'abord que pour une table de contingence quelconque (si l'on excepte les tables symétriques), la classification hiérarchique donnera des indices différents selon que l'on agrège les lignes et les colonnes, alors que l'analyse des correspondances ne fournit qu'une série de valeurs propres.

La plus grande valeur propre issue de l'analyse des correspondances est supérieure ou égale au plus grand indice d'agrégation (lignes ou colonnes) donné par la classification. Cet indice est en effet une mesure de variance externe (dite variance "inter", par opposition à la variance "intra", mesurant la dispersion à l'intérieur des groupes) entre les deux derniers groupes agrégés. Cette variance externe est inférieure à la variance totale mesurée sur la droite qui joint les centres de gravités des deux groupes, elle-même inférieure à la meilleure variance totale possible sur une droite quelconque, ce qui est la définition de la plus grande valeur propre<sup>2</sup>.

Plus généralement, Benzécri et Cazes (1978) ont montré que la somme des  $r$  plus grandes valeurs propres est supérieure ou égale à la somme des  $r$  plus grands indices d'agrégation.

Enfin, ces auteurs ont donné un intéressant contre-exemple montrant qu'il n'existe pas de borne inférieure positive pour le quotient entre le plus grand indice d'agrégation et la plus grande valeur propre : on peut trouver des distributions de densité telles que le plus grand indice soit une fraction arbitrairement petite de la plus grande valeur propre.

<sup>1</sup> On trouvera la preuve de la non-négativité des termes  $c_{ii'}$  dans Benzécri (1973, Tome IIB, Chapitre 11).

<sup>2</sup> Notons bien, sur la figure 2.4 - 4 précédente, le cas de coïncidence pour lequel les variances "intra" sur l'axe sont nulles, et pour lequel le meilleur axe factoriel est précisément celui qui relie les deux centres de classes.

### b – Le cas des tables de contingence structurées par blocs

Cette structure déjà évoquée en section 1.3.4 (cf. les figures 1.3 - 15 et 1.3 - 16) est aisément reconnue par l'analyse des correspondances car  $k$  blocs engendrent  $k$  valeurs propres égales à 1 (y compris la valeur propre triviale, qui correspond au cas usuel d'un seul bloc).

Cette structure n'est cependant pas systématiquement reconnue par la classification hiérarchique utilisant le critère de Ward, comme l'ont montré par un contre-exemple Kharchaf et Rousseau (1988, 1989).

### c – Une étude empirique du lien entre valeurs propres et indices

Ces inégalités et contre-exemples ne donnent que peu d'information sur les liaisons entre valeurs propres et indices, et les liaisons fonctionnelles du paragraphe 2.4.2 ne concernent que des cas d'école. Les liaisons stochastiques entre indices et valeurs propres (dans le cas d'une famille de tables de contingences aléatoires) sont certainement trop complexes pour faire l'objet d'une étude analytique.

**Tableau 2.4 - 5**  
Moyennes et écart-types des valeurs propres  
et des indices d'agrégation.  
1000 tables de contingences pseudo-aléatoires (8,8).  
Pour chaque table,  $k = 1000$ .

Identificateur	Moyenne	Ecart type	Ecart type de la moyenne
<b>Valeurs propres</b>			
VP1 *	.02130 *	.00560 *	.00018
VP2 *	.01282 *	.00353 *	.00011
VP3 *	.00772 *	.00234 *	.00007
VP4 *	.00442 *	.00156 *	.00005
VP5 *	.00214 *	.00100 *	.00003
VP6 *	.00070 *	.00050 *	.00002
VP7 *	.00010 *	.00014 *	.00000
<b>Indices des lignes (INL<sub>i</sub>) et des colonnes (INC<sub>i</sub>)</b>			
INL1 *	.01692 *	.00452 *	.00014
INL2 *	.01063 *	.00289 *	.00009
INL3 *	.00733 *	.00197 *	.00006
INL4 *	.00537 *	.00148 *	.00005
INL5 *	.00391 *	.00117 *	.00004
INL6 *	.00280 *	.00090 *	.00003
INL7 *	.00183 *	.00074 *	.00002
INC1 *	.01679 *	.00450 *	.00014
INC2 *	.01061 *	.00291 *	.00009
INC3 *	.00739 *	.00202 *	.00006
INC4 *	.00535 *	.00151 *	.00005
INC5 *	.00396 *	.00118 *	.00004
INC6 *	.00280 *	.00091 *	.00003
INC7 *	.00182 *	.00075 *	.00002

Une exploration par simulation pourra cependant donner une idée des liens stochastiques existant entre indices et valeurs propres.

Pour procéder à cette exploration, des tables de contingence à 8 lignes et 8 colonnes ont été simulées sous l'hypothèse d'indépendance selon un schéma de remplissage multinomial (les marges théoriques sont supposées égales, l'effectif total de chaque table simulée est  $k = 1\ 000$ ).

1 000 simulations ont été réalisées, donnant lieu chacune à une analyse des correspondances, et à deux classifications hiérarchiques (selon le critère de la variance) : une sur les lignes et une sur les colonnes.

Le tableau 2.4 - 5 donne les moyennes des 7 valeurs propres, des 7 indices-lignes et des 7 indices colonnes, calculées sur 1 000 observations. Les indices d'agrégation des lignes suivent évidemment la même loi que ceux des colonnes, cette propriété permettant de vérifier la cohérence de la simulation<sup>1</sup>.

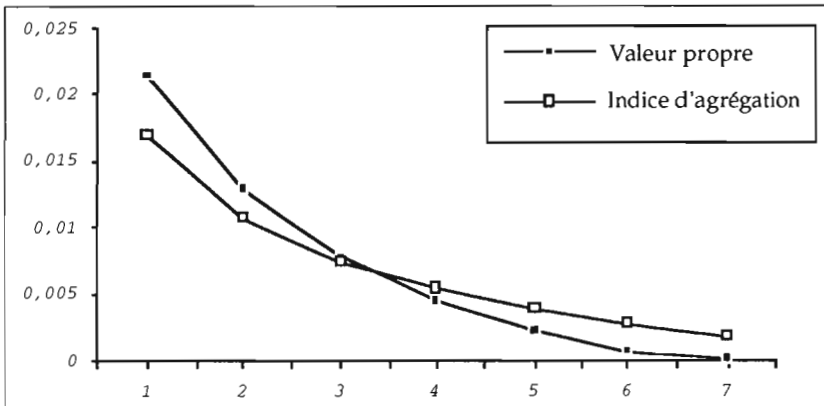


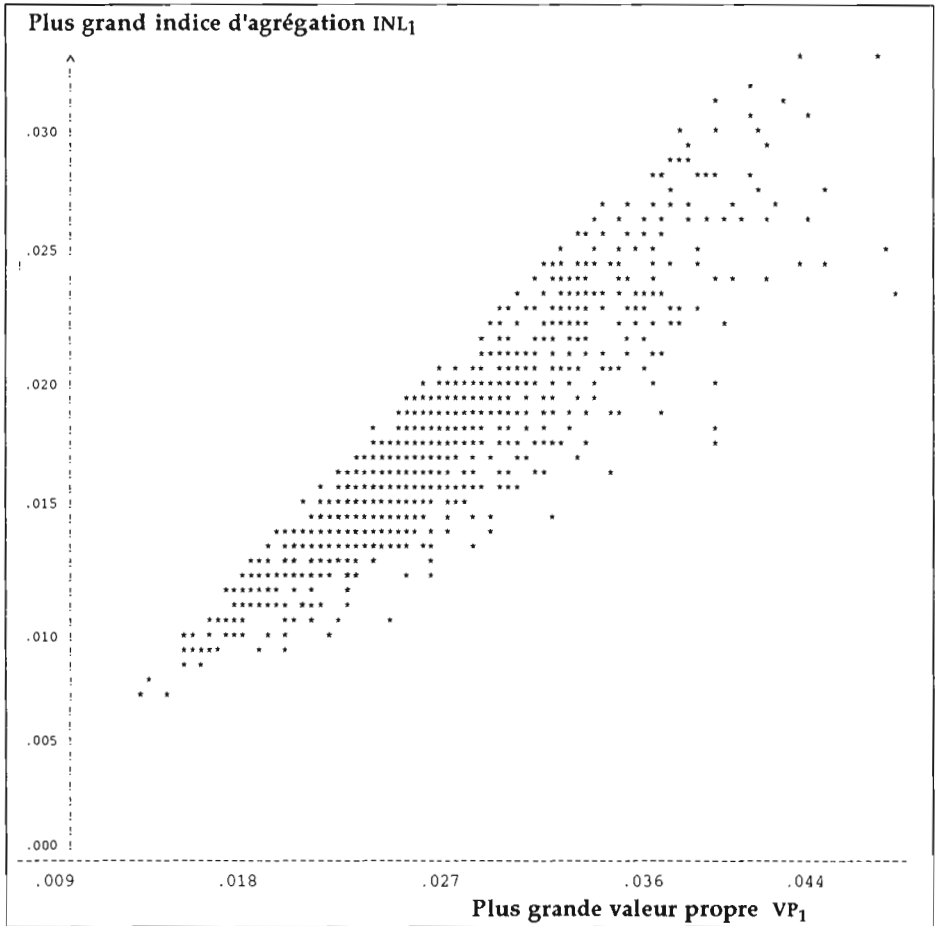
Figure 2.4 - 5  
Séquences des valeurs propres et des indices

La figure 2.4 - 5, qui représente graphiquement les éléments de la première colonne du tableau 2.4 - 5, met en évidence l'intervalle de variation plus réduit des indices dans l'hypothèse d'indépendance des lignes et des colonnes. Il est intéressant de compléter ces mesures de niveau par une analyse des corrélations entre valeurs propres et indices.

La figure 2.4 - 6 présente le diagramme de la distribution jointe de la première valeur propre  $\lambda_1 = VP1$  et du plus grand indice-ligne de classification INL1, chacun des 1000 couples (VP1, INL1) correspondant à une même matrice pseudo-aléatoire.

<sup>1</sup> Remarquons que ces résultats sont cohérents avec le test usuel d'indépendance du  $\chi^2$  (la somme  $t$  des différentes valeurs propres vaut 0,0492, le  $\chi^2$  usuel moyen valant  $1000 \times t = 49,2$  pour 49 degrés de liberté. Les approximations connues de la loi des valeurs propres (loi des valeurs propres d'une matrice de Wishart (7,7)) sont également vérifiées ici (cf. § 4.1.2).

Le coefficient de corrélation<sup>1</sup> entre  $VP_1$  et  $INL_1$  est de 0.91. La contrainte théorique  $INL_1 < \lambda_1$  définit de façon claire le demi-plan contenant le nuage de 1000 points. On voit que les écarts entre valeurs propres et indices peuvent être notables, ceux-ci pouvant parfois être de 30% inférieurs à celles-là.



**Figure 2.4 - 6**  
**Corrélation entre la plus grande valeur propre  $VP_1$**   
**et le plus grand indice d'agrégation  $INL_1$ .**

(Chacun des 1000 points correspond à une matrice pseudo-aléatoire (8,8))

L'étude du système complexe des corrélations entre valeurs propres et indices sera l'occasion de présenter ci-dessous une application méthodologique de l'analyse en composantes principales.

<sup>1</sup> Le coefficient de corrélation entre  $VP_1$  et  $INC_1$  a la même valeur.

Les résultats de cette expérience peuvent être présentés dans un tableau  $X$  ayant 1 000 lignes (les 1 000 tableaux simulés) et 21 colonnes (les 7 valeurs propres et les  $2 \times 7 = 14$  indices d'agrégation).

On a choisi ici de procéder à une analyse en composantes principales avec comme variables actives les 7 premières colonnes, les indices étant projetés en variables illustratives. On privilégie donc la structure des corrélations interne à l'ensemble des valeurs propres, et l'on situe ensuite les indices par rapport à cette structure.

La figure 2.4 - 7 représente le premier plan factoriel ainsi obtenu, qui correspond à environ 60% de la variance totale.

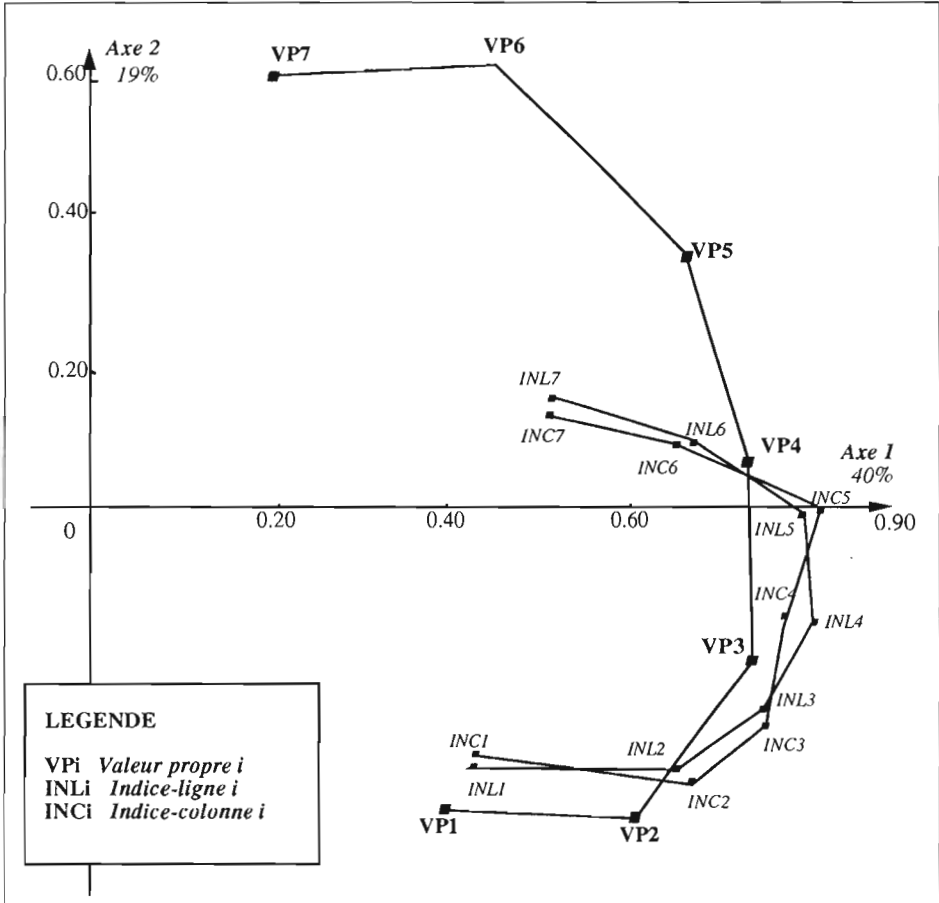


Figure 2.4 - 7

**Structure des corrélations entre valeurs propres et indices**

Plan Principal d'une analyse en composantes principales de la matrice (1000,7) contenant les 1000 observations (en lignes) des 7 valeurs propres VP1,...VP7.

Les 7 indices-ligne INL1, ...INL7 et les 7 indices-colonnes INC1,...INC7 sont projetés en éléments supplémentaires dans ce plan.

Les principaux éléments d'interprétation sont les suivants :

- On note tout d'abord que le premier facteur est un facteur de taille : tous les points-variables sont situés dans le demi-plan des points dont les coordonnées sont positives sur cet axe. Brièvement, cela signifie qu'il y a des tables pour lesquelles toutes les valeurs propres sont grandes, et d'autres pour lesquelles elles sont toutes petites, et que ce facteur d'échelle général est la principale source de variabilité<sup>1</sup>.
- On note ensuite, en remarquant que les 7 valeurs propres forment une trajectoire régulière, qu'il existe une corrélation entre valeurs propres consécutives : la première valeur propre est plus liée à la deuxième qu'à la troisième, etc. Inversement, les couples de valeurs propres de rangs éloignés (1,6), (1,7), (2,7) sont corrélés négativement.
- Les indices lignes et colonnes, sensibles eux aussi à l'effet-taille, ont des trajectoires très voisines, dont les positions et la forme font penser à celles des valeurs propres, avec cependant un décalage très net du côté des plus grandes valeurs propres. Les premiers (plus grands) indices et les premières valeurs propres sont donc fortement liés entre eux (on a vu précédemment que  $VP_1$  et  $INL_1$  avait un coefficient de corrélation de 0.91), mais les derniers indices sont peu corrélés avec les dernières valeurs propres.

En fait, cette structure est en accord avec l'expérience des praticiens de ces méthodes. Il existe très souvent des relations entre les premiers facteurs et les ultimes nœuds du dendrogramme correspondant aux plus grandes valeurs de l'indice.

En revanche, les premiers nœuds du dendrogramme fournissent souvent de précieuses informations sur des groupements ou des structures locales qui correspondent à des facteurs de rangs moyens, mais non aux derniers facteurs. Ceux-ci rendent plutôt compte d'un bruit rarement identifiable.

#### 2.4.4 La complémentarité en pratique : un exemple

Cet exemple d'enchaînement résume certaines étapes d'une application "en vraie grandeur". Il est extrait de traitements de l'enquête sur les conditions de vie et aspirations des Français<sup>2</sup>

L'objectif poursuivi ici est double : donner une description d'ensemble des principales attitudes et opinions relevées dans le système d'enquêtes précité ; montrer dans quel cadre factuel s'inscrivent les attitudes et opinions.

<sup>1</sup> Ce premier facteur est évidemment très lié à la somme  $t$  des valeurs propres, donc au  $\chi^2$  (ici  $\chi^2 = 1\,000t$ ).

<sup>2</sup> Cf. Lebart et Houzel (1981), Babeau et Lebart (1984), Lebart (1987b) pour des informations générales sur cette enquête.

Le fichier partiel correspondant à cette application comprend 14 variables nominales actives et en fait plusieurs centaines de variables nominales supplémentaires. Les 14 000 individus correspondent à 7 vagues de 2000 individus (de 1978 à 1984), chaque vague étant représentative de la population de résidents métropolitains âgés de 18 ans ou plus. Un des intérêts de cet exemple est que les structures observées pourront être validées par les échantillons indépendants annuels. Il s'agit d'une situation exceptionnellement favorable pour éprouver la stabilité des résultats d'une analyse exploratoire.

**14 questions actives pour décrire les perceptions des conditions de vie et du cadre de vie (60 modalités)**

*Deux questions sur la perception de l'évolution des conditions de vie*  
*Trois questions sur le thème «Famille»*  
*Trois questions sur l'environnement physique et technologique*  
*Trois questions sur la santé et l'institution médicale*  
*Une question sur l'attitude vis-à-vis des équipements collectifs*  
*Deux questions sur la justice et la société*

**a – Les étapes**

L'enchaînement de méthodes décrit ici est une formulation plus détaillée de la procédure d'utilisation conjointe des méthodes factorielles et de la classification exposée au § 2.4.1.b. Cette procédure est présentée du point de vue du praticien.

**- Etape 1 : Analyse factorielle**

L'étape 1 (analyse factorielle), comprend les trois phases suivantes :

**- Choix d'un thème actif**

Choisir un thème, c'est-à-dire une batterie homogène de variables actives, c'est adopter un point de vue particulier pour la description. On peut décrire les individus du point de vue de leurs caractéristiques de base, mais aussi à partir d'un thème particulier de l'enquête par exemple les habitudes de consommation, les durées d'activité (budgets-temps), les contacts-médias, les déplacements, etc. Ici, le thème choisi est : la perception des conditions de vie et du cadre de vie (cf. encadré ci-dessus).

**- Description graphique de la population**

Les graphiques résultant des analyses factorielles (ici : correspondances multiples) fournissent une description de l'échantillon des individus interrogés. La proximité entre individus est fonction de la similitude des réponses aux questions du thème actif.

**- Positionnement des éléments illustratifs sur les plans factoriels**

On s'intéresse aux questions ne faisant pas partie du thème actif pour aider à interpréter les proximités entre individus. Lorsque la lecture des résultats est gênée par l'abondance des éléments illustratifs, les seuls

éléments pertinents pour l'interprétation seront sélectionnés par leurs valeurs-test. Ceci permet d'envisager des explorations systématiques, avec de nombreux croisements de variables.

Comme au § 2.4.1 b, les trois phases suivantes sont :

- *Etape 2 : Partition de l'ensemble des individus*
- *Etape 3 : Descriptions statistiques du contenu de chaque classe*
- *Etape 4 : Positionnement des centres des classes en éléments supplémentaires dans les plans factoriels*

Cet enchaînement est souvent utilisé sous le nom de *thémascope*. C'est donc un outil qui permet de décrire un thème (actif), multidimensionnel par nature, en utilisant la conjonction des deux techniques disponibles (réduction de dimension d'une part, regroupement d'autre part). Il situe ensuite ce thème dans le contexte global de l'enquête, grâce aux techniques de projection de variables supplémentaires sur les plans factoriels et de description automatique des classes. La sélection automatique des éléments les plus significatifs sur les plans factoriels et lors de la description des classes fournit au lecteur une information filtrée et lisible.

#### **b – L'espace des variables actives (Figure 2.4 - 8)**

La figure 2.4 - 8 est l'esquisse du premier plan factoriel d'une analyse des correspondances multiples du tableau (14 000, 60). Les 14 réponses aux questions actives (60 modalités) répartissent les individus interrogés de façon continue dans l'espace. Il n'existe pas de regroupement très net d'individus dans ce continuum, mais il est toujours possible de le découper en grandes zones de la façon la moins arbitraire possible ; les cloisons entoureront ainsi les régions de forte densité et seront disposées de façon à ce que la dispersion des individus soit minimale à l'intérieur des zones. C'est l'arbre hiérarchique de la figure 2.4 - 9 qui est schématiquement tracé sur le plan factoriel (coupure correspondant à 8 classes). Pour limiter le nombre de graphiques, le résultat de l'*étape 4* figure d'emblée sur la figure.

#### **c – Exemples de description automatique de trois classes**

On va maintenant illustrer la description automatique des classes (cf. § 2.3.2) en caractérisant de façon plus détaillée trois classes (ou zones) sélectionnées parmi les huit précédentes. On distinguera successivement les opinions et perceptions (éléments actifs, et pour certains d'entre-eux, supplémentaires), puis les caractéristiques de base (éléments toujours supplémentaires dans cette analyse).

Chaque pourcentage interne à la zone sera suivi, entre parenthèses, du pourcentage moyen dans l'ensemble de la population. Les valeurs-test (cf. § 2.3.2.b) qui ont permis de sélectionner et de classer ces variables caractéristiques sont des fonctions de l'écart entre ces deux pourcentages.



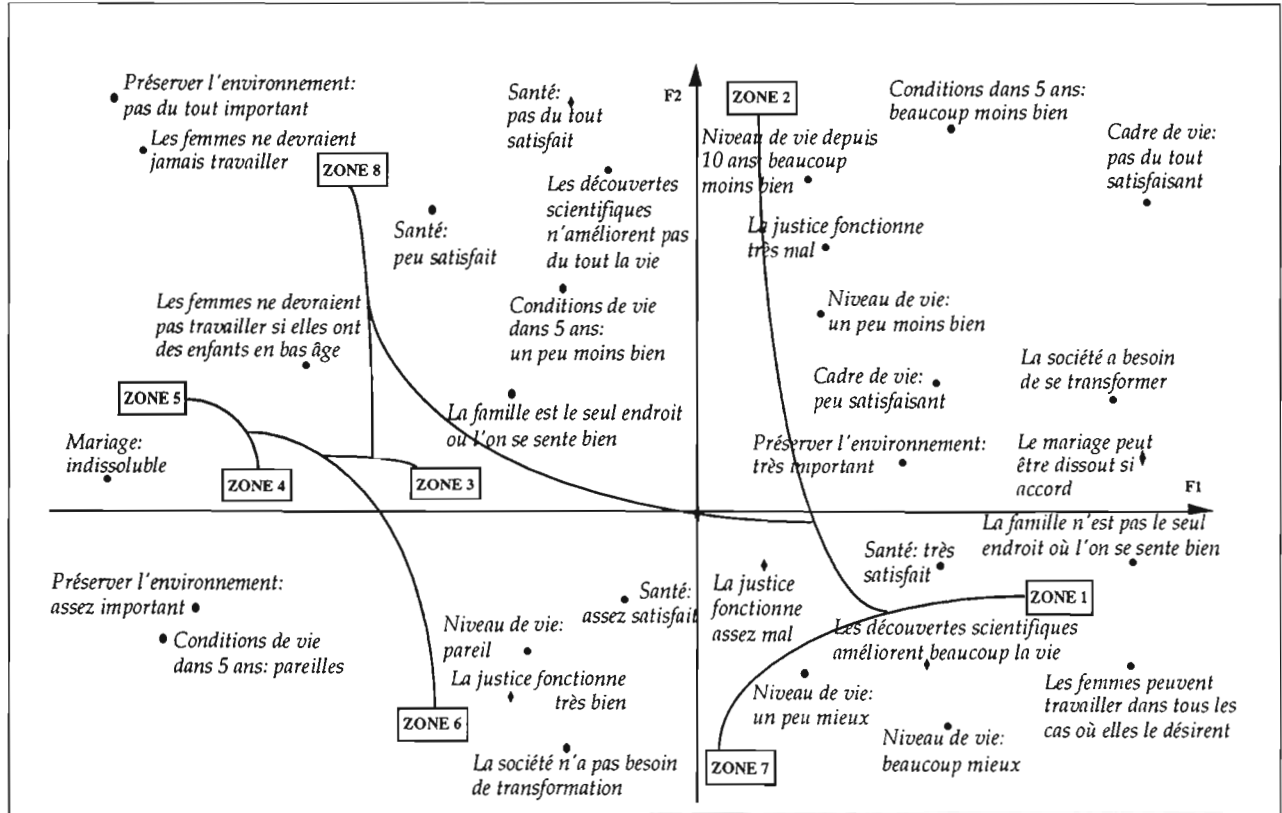
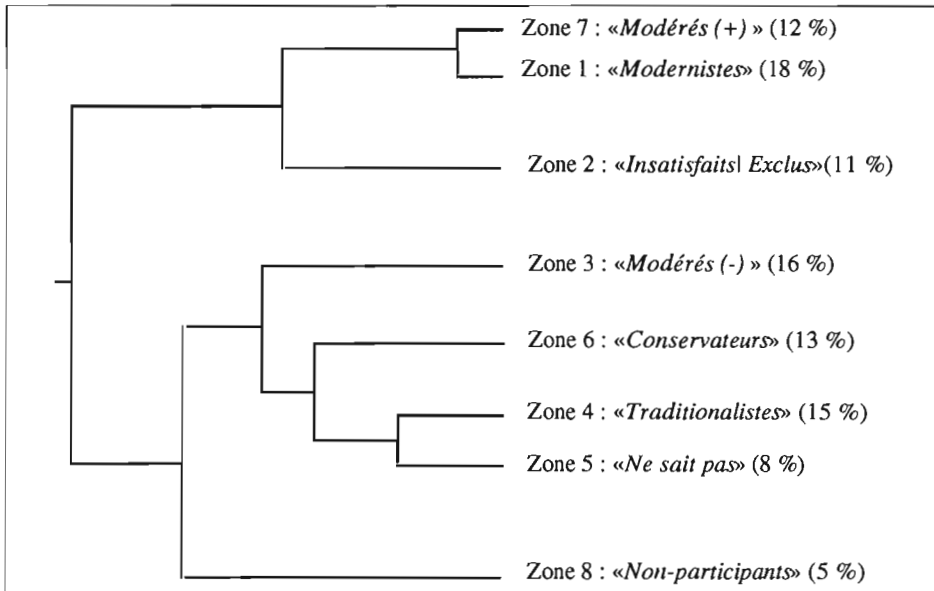


Figure 2.4 - 8 : Visualisation plane de l'espace des opinions et positionnement des zones

On représente ici les proximités statistiques existant entre une trentaine de modalités de réponses aux questions actives choisies parmi les plus caractéristiques. Les centres des zones sont positionnés comme des modalités supplémentaires.



**Figure 2.4 - 9**  
**Classification hiérarchique des 14 000 individus en 8 zones**  
**Guide de lecture du dendrogramme :**

L'algorithme de classification mixte de la section 2.3 permet de mettre en évidence huit zones<sup>1</sup>, positionnées en éléments supplémentaires sur le plan factoriel de la figure 2.4 - 8, et comme éléments terminaux du dendrogramme de la figure 2.4 - 9. Cette figure permet donc de compléter la figure 2.4 - 8. Ainsi, contrairement à ce que l'on observe sur la figure 2.4 - 8 qui ne donne qu'une approximation plane de l'espace, et donc qui déforme les distances, la zone 2 est, d'après le dendrogramme, plus proche des zones 1 et 7 que de la zone 8.

### Description de la zone 1 (Modernistes) [droite de la figure 2.4 - 8]

Cette zone stable représentant en moyenne 18% des personnes interrogées se distinguent par une certaine distance vis-à-vis de la famille traditionnelle.

#### **Variables actives**

- 87% pensent que «la famille n'est pas le seul endroit où l'on se sent bien et détendu» (ce pourcentage n'est que de 35% pour l'ensemble de la population)
- 84% déclarent que «le mariage est une union qui peut être dissoute sur simple accord» (35%)
- 83% estiment que «les femmes devraient travailler dans tous les cas où elles le désirent» (37%)
- 86% jugent que «préserver l'environnement est une chose très importante» (65%)

<sup>1</sup> On parle de zones et non de classes ou de groupes pour rappeler qu'il s'agit de portions d'espace et non d'entités sociologiques ou de catégories ayant une existence indépendante de la batterie des questions actives utilisées ici. Les libellés de ces zones sont purement mnémotechniques.

**Variables supplémentaires (signalétique) : jeunes, instruits, parisiens**

- 52% n'ont jamais eu d'enfant (28%)
- 32% habitent la région parisienne (15%)
- 78% ont moins de 40 ans (47%)
- 67% sont des locataires (51%)
- 20% sont diplômés d'université ou de grande école (8%)

**Autres variables supplémentaires : Spécificités de comportement**

- 31% se couchent après 23 h (13%)
- 35% fréquentent un cinéma (17%)
- 57% participent aux activités d'au moins une association (44%)

**Description de la zone 2 (Insatisfaits / exclus) [haut de la figure 2.4 - 8]**

Cette zone est probablement la seule à mériter le statut de «classe» au sens statistique du terme dans la mesure où elle réapparaît chaque année (de 1978 à 1985) avec un effectif remarquablement constant qui oscille entre 9% et 13%.

**Opinions et perceptions : niveau et cadre de vie non satisfaisants**

- 69% pensent que leur «niveau de vie personnel va beaucoup moins bien» (13%)
- 62% estiment que leurs «conditions de vie vont beaucoup se détériorer au cours des cinq prochaines années» (12%)
- 61% considèrent que «la justice fonctionne très mal» (26%)
- 85% déclarent «s'imposer régulièrement des restrictions» (61%)
- 17% ne sont «pas du tout satisfaits de leur cadre de vie quotidien» (5%) ; 21% en sont «peu satisfaits» (14%)
- 90% pensent que «la société a besoin de se transformer» (74%)

**Variables supplémentaires (signalétique) : des ressources faibles <sup>1</sup>**

- 38% souffrent d'un handicap, d'une infirmité ou d'une maladie chronique (26%)
- 38% n'ont aucun élément de patrimoine (27%)
- 15% sont chômeurs (en 1983 et 84) (6%)
- 53% sont locataires (44%)
- 22% habitent en HLM ou ILN (16%) 9% sont séparés ou divorcés (5%)

**Autres variables supplémentaires :**

- 55% ont déclaré «avoir souffert de nervosité au cours des quatre dernières semaines» (37%),
- 28% ont dit avoir souffert d'«état dépressif» (15%),
- 38% d'«insomnie» (25%),
- 49% de «mal au dos» (38%),
- 45% s'estiment «beaucoup inquiets de l'éventualité du chômage» (25%).

**Description de la zone 5 (réponses "ne-sait-pas") [gauche de la figure 2.4 - 8]**

Cette zone *a priori* peu intéressante du point de vue des opinions exprimées joue cependant un rôle méthodologique important.

<sup>1</sup> Cette zone n'a pas de caractéristiques socio-démographiques aussi typées que la zone 1. Elle constitue avant tout une classe de personnes aux ressources faibles, au niveau de vie bas, qui subissent des tensions où font face à des difficultés variées. On a affaire ici typiquement une «classe polythétique», c'est-à-dire une classe qui peut être définie non par une combinaison fixe d'attributs, mais par la possession d'un certain nombre d'attributs dans une liste : il y a dans ce cas cumul de handicaps d'origines variées.

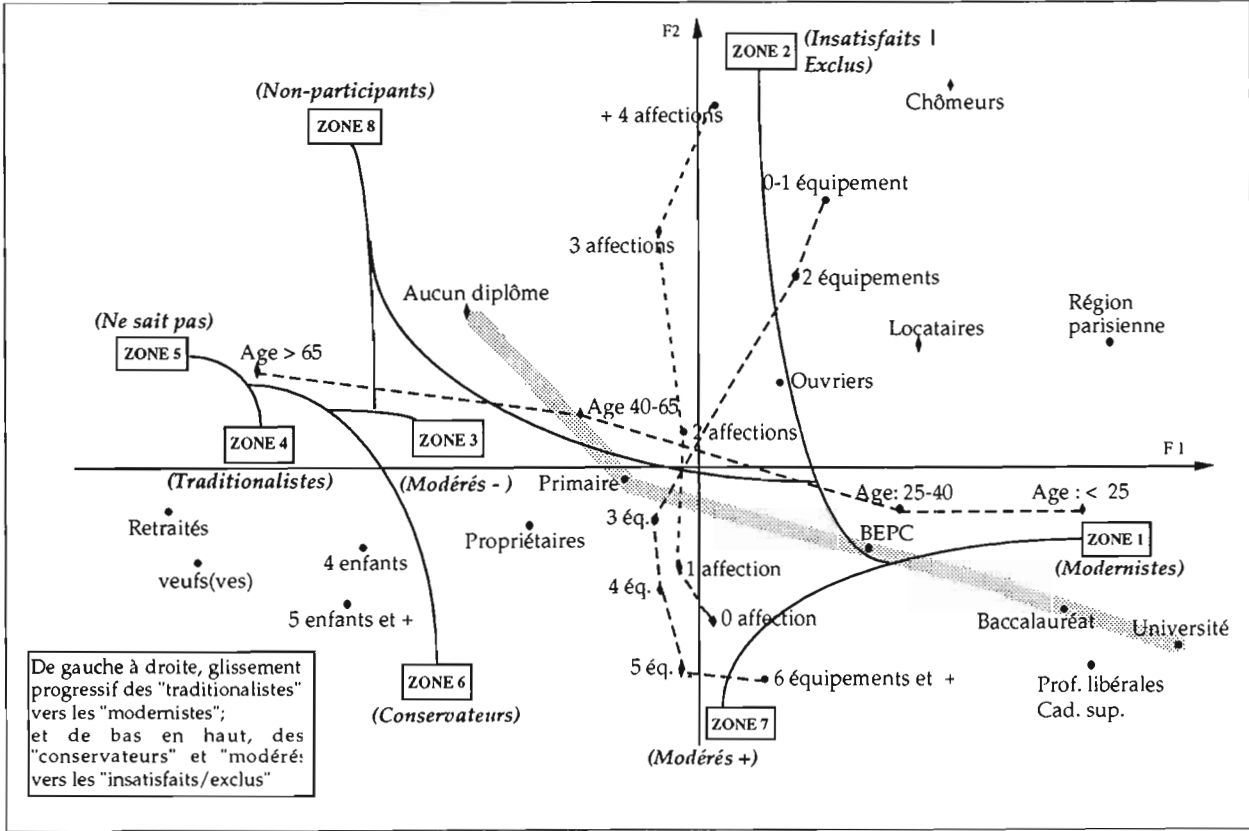


Figure 2.4 - 10 : Projection de quelques caractéristiques (en supplémentaires) sur le plan principal de la figure 2.4 - 8.

Alors que les refus ou les dissimulations entachent la qualité des enquêtes socio-économiques usuelles, les réponses du type «ne sait pas» viennent s'ajouter aux défections précédentes dans le cas des mesures de perceptions ou d'opinions.

**Variables actives**

- 65% répondent NSP (pour «ne sait pas») à la question «la société a-t-elle besoin de se transformer ?» (9%)
- 53% répondent NSP à la question sur «le fonctionnement de la justice» (7%) ; 8% refusent de répondre à cette question (2%)

**Variables supplémentaires (signalétique) : femmes âgées peu instruites**

- 67% sont des femmes (53%)
- 46% n'ont aucun diplôme (26%)
- 43% habitent des communes de moins de 2 000 habitants (29%)
- 75% n'appartiennent à aucune association (56%)<sup>1</sup>.

**d – Projection de variables signalétiques (en supplémentaires) sur le plan principal de la figure 2.4 - 8 (figure 2.4 - 10)**

Les descriptions zones par zones donnent déjà une idée de l'«ancrage factuel» des perceptions, mais un positionnement direct des caractéristiques de base a le mérite de montrer à quel point l'espace des perceptions est un continuum<sup>2</sup>. Les modalités des différentes variables s'ordonnent en effet régulièrement dans le plan de la figure 2.4 - 10.

Il n'y a pas de discontinuité entre les «traditionalistes» âgés, ruraux, peu instruits situés dans la partie gauche de la figure 2.4 - 10 et les «modernistes» jeunes, instruits, urbains, situés à l'extrémité droite de l'axe horizontal.

Il y a de même une certaine continuité entre les «conservateurs» et les «modérés + » d'âge moyen situés dans la partie basse de la figure 2.4 - 10 et les insatisfaits dans la partie haute. *Le nombre d'équipements et d'éléments de patrimoine* jalonne régulièrement cette direction verticale, tout comme *le nombre d'affections déclarées* (petites affections au cours des quatre dernières semaines), indicateur dont les liens avec l'insatisfaction sont connus.

<sup>1</sup> Le fait qu'il s'agisse surtout de femmes âgées peu instruites habitant en milieu rural, alors que les questions «non répondues» sont peut-être les plus politiques de ce questionnaire (les transformations de la société, la justice) confirme les résultats de travaux de méthodologie d'enquête (cf. par exemple Michelat et Simon, 1985).

<sup>2</sup> L'étude complète comporte une description beaucoup plus détaillée de l'ensemble des classes, une étude de l'évolution des trajectoires des points-modalités et des classes dans les plans factoriels au cours du temps, et l'utilisation systématique de croisements de variables supplémentaires (cf. Lebart, 1986; 1988). La sélection automatique des catégories provenant de croisements de variables supplémentaires par leurs valeurs-test (avec des seuils prenant en compte les comparaisons multiples) est un outil efficace de détection d'interactions.

## Chapitre 3

---

# **LIENS AVEC LES MÉTHODES EXPLICATIVES USUELLES, MÉTHODES DÉRIVÉES**



## Introduction

Ce chapitre fait le lien entre les *approches exploratoires* présentées dans les chapitres 1 et 2 et les *approches inférentielles et confirmatoires* qui constituent le volet le plus ample et le plus classique de la science statistique.

Rappelons brièvement les caractéristiques de ces deux familles de méthodes, qui correspondent à des approches complémentaires.

- *La statistique descriptive et exploratoire* : elle permet par des résumés et des graphiques plus ou moins élaborés de décrire des ensembles de données statistiques, d'établir des relations entre les variables sans faire jouer de rôle privilégié à une variable particulière. Les conclusions ne portent dans cette phase de travail que sur les données étudiées, sans être inférées à une population plus large. L'analyse exploratoire s'appuie essentiellement sur des notions élémentaires telles que des indicateurs de moyenne et de dispersion, sur des représentations graphiques et sur les techniques descriptives multidimensionnelles abordées dans la première partie (analyse en composantes principales, analyse des correspondances, classification).

- *La statistique inférentielle et confirmatoire* : elle permet de valider ou d'infirmer, à partir de tests statistiques ou de modèles probabilistes, des hypothèses formulées a priori (ou après une phase exploratoire), et d'extrapoler, c'est-à-dire d'étendre certaines propriétés d'un échantillon à une population plus large. Les conclusions obtenues à partir des données vont au delà de ces données. La statistique confirmatoire fait surtout appel aux méthodes dites explicatives<sup>1</sup> et prévisionnelles destinées, comme leurs noms l'indiquent, à expliquer puis à prévoir, suivant des règles de décision, une variable privilégiée à l'aide d'une ou de plusieurs variables explicatives (régressions multiples et logistiques, analyse de la variance, analyse discriminante, segmentation, etc.).

Les démarches sont complémentaires, l'exploration et la description devant en général précéder les phases explicatives et prédictives. En effet, une exploration préliminaire est souvent utile pour avoir une première idée de la nature des liaisons entre variables, et pour traiter avec prudence les variables corrélées et donc redondantes qui risquent de charger inutilement les modèles.

---

<sup>1</sup> La statistique n'explique rien mais fournit des éléments potentiels d'explication. Aussi le terme de variable explicative ou variable à expliquer n'est sans doute pas le plus judicieux. On dit aussi indépendante et dépendante, ou exogène et endogène. Ces deux derniers termes sont peut être les plus adéquats mais ne sont pas assez évocateurs. L'adjectif indépendant est, en revanche, source de confusions.



Cependant, les démarches elles-mêmes ne sont pas toujours faciles à discerner, à identifier. L'exploration pure est très rare, et correspond à une situation limite et irréaliste, un peu comme les gaz parfaits en physique... car il existe toujours des informations et des connaissances *a priori* sur le tableau de données, et donc des hypothèses générales, des attentes de la part de l'utilisateur<sup>1</sup>.

D'où l'intérêt d'éclaircir cette relation entre *instruments d'observation et modèles*, en insistant sur l'insertion, théorique et pratique, des outils exploratoires dans l'arsenal des techniques statistiques disponibles.

### *Les méthodes explicatives usuelles*

Les méthodes explicatives présentées dans les sections 3.1 à 3.5 recouvrent les utilisations les plus courantes. Elles comprennent l'analyse canonique, la régression linéaire et ses variantes, l'analyse discriminante, les modèles log-linéaires, les méthodes de segmentation par arbre binaire.

- Parce que l'*analyse canonique* joue un rôle théorique important dans les méthodes multidimensionnelles et permet de jeter un pont entre les formalismes des méthodes explicatives et descriptives, nous commencerons ce chapitre par exposer ses principes (section 3.1). On verra que l'analyse canonique, qui étudie les liaisons entre deux groupes de variables, contient comme cas particuliers la *régression multiple* si l'un des deux groupes est réduit à une seule variable  $y$  numérique, l'*analyse discriminante* lorsque les variables de l'un des deux groupes sont les variables indicatrices d'une partition des individus (ce qui revient à dire que la variable  $y$  est nominale), enfin l'*analyse des correspondances* si les deux groupes sont constitués par les variables indicatrices des deux partitions.

- La *régression multiple* (section 3.2) se situe directement dans le cadre théorique du *modèle linéaire*, lorsque la variable à expliquer  $y$  est une variable continue (ou numérique). Les variables explicatives sont généralement continues. Lorsque les variables explicatives sont toutes nominales, on parle plutôt d'*analyse de la variance*, alors qu'on réserve le nom d'*analyse de covariance* au cas mixte (variables explicatives nominales et continues).

- L'*analyse factorielle discriminante* (section 3.3) est, schématiquement, l'analogue de la régression multiple lorsque  $y$  est nominale. Dans ce cas, la variable à expliquer définit les classes d'une partition a priori de la population. L'objet est alors d'étudier les liaisons entre les variables explicatives et les classes de cette partition. On définit ainsi des fonctions discriminantes qui vont permettre, dans une phase décisionnelle, d'affecter

---

<sup>1</sup> Les instruments d'observation correspondent d'ailleurs eux-mêmes à des modèles généraux : ainsi, les axes factoriels de l'analyse en composantes principales sont proches de ceux de l'analyse factorielle classique des psychologues (cf. section 3.2.9) qui représentent les variables latentes d'un modèle a priori. Inversement, la régression multiple, méthode explicative par excellence, peut aussi être utilisée pour explorer des structures de corrélation.

de nouveaux individus à ces classes. D'autres méthodes de discrimination sont brièvement évoquées.

- Bien qu'extérieurs à ce cadre formel général, les *modèles log-linéaires* (section 3.4) sont utilisés dans des circonstances voisines. Ce sont des techniques d'analyse des tableaux de contingence multidimensionnels qui se rapprochent de la régression multiple dans leur problématique. Les modèles log-linéaires peuvent d'ailleurs être considérés comme une extension du modèle logistique également abordé dans cette section.

- Les techniques de *segmentation par arbre binaire* (section 3.5) sont intéressantes à présenter dans le cadre de ce chapitre pour diverses raisons. D'une part, elles s'appliquent à toutes les variables quel que soit leur statut ou leur nature, et d'autre part elles intègrent simultanément la phase explicative et de décisionnelle. Elles constituent de ce fait une méthode de prévision à part entière, très accessible, dont les résultats sont faciles à communiquer.

### *Les analyses de données structurées*

Les sections 3.6 à 3.8 contiennent une série de présentations, souvent brèves, de méthodes qui occupent une position intermédiaire entre les outils purement exploratoires des deux premiers chapitres et les méthodes à vocation plus explicative présentées dans les sections précédentes.

Les méthodes exploratoires de base posent un modèle très général qui distingue, pour chaque application, deux familles d'éléments : les éléments actifs (variables ou individus, ligne ou colonnes) qui servent à établir des espaces de visualisation complétés par des classifications, et les éléments supplémentaires, qui jouent un rôle passif, et interviennent *a posteriori* pour illustrer, identifier, caractériser les représentations obtenues à partir des éléments actifs.

En général, le tableau des éléments actifs est amorphe et homogène : il ne doit pas exister de structure *a priori* (dépendance fonctionnelle, relations comptables, etc.) entre les variables et les individus, et les distances entre éléments doivent avoir un sens pour l'utilisateur.

Or, il est fréquent que le tableau des données actives soit déjà structuré. C'est le cas par exemple des données géographiques ou temporelles où la structure intervient au niveau des observations (individus voisins ou consécutifs). Il peut exister des groupes d'individus ou des groupes de variables connus *a priori*. Le tableau peut ne pas se ramener de façon univoque à la forme rectangulaire (tables de contingences multiples, séries chronologiques de tableaux).

Il est souvent possible d'aborder ces problèmes dans le cadre du modèle exploratoire de base, mais la tentation est forte, dans le cas où les applications se présentent de façon répétitive, de proposer des variantes adaptées aux types de tableaux ou de structures rencontrés. Il reste que l'on doit envisager une économie de l'analyse des données, en ce sens que la

panoplie des méthodes disponibles ne peut s'accroître indéfiniment, sous peine de voir le rendement de ces méthodes décroître<sup>1</sup>.

A propos des méthodes de classification pour lesquelles il estime le nombre de publications à près de mille par an, Cormack (1971) remarque que "lorsque la technique (de classification) échoue, la réaction de l'auteur est de modifier la technique, au lieu d'utiliser une technique plus *standard* ou de remettre en question tout le traitement". Cette attitude comporte un certain danger. Si la panoplie des techniques est très étendue, le risque d'adéquation accidentelle de la technique aux données est augmenté. Ce problème est récurrent lorsqu'il s'agit d'articuler exploration et inférence, et se rapproche du problème plus classique des comparaisons multiples, déjà évoqué à propos de la description des classes par les valeurs-test, et dont on reparlera à propos du modèle log-linéaire. Un défi auquel est confrontée la statistique multidimensionnelle est précisément la gestion de cette diversification, nécessaire pour la recherche, mais source de difficultés au niveau des applications en vraie grandeur. Précisons, dans ce contexte méthodologique, quelles sont les méthodes d'analyses de données structurées qui feront l'objet des trois dernières sections de ce chapitre.

Les méthodes d'*analyses partielles ou projetées* (section 3.6) concernent les situations pour lesquelles les individus ou observations (lignes d'un tableau  $X$  d'ordre  $(n, p)$ ) peuvent être décrits par  $p$  variables (colonnes de  $X$ ) mais peuvent aussi être dépendants de  $q$  variables : colonnes d'un tableau  $Z$  d'ordre  $(n, q)$  dont on désirerait, dans la mesure du possible, soit prendre en compte, soit éliminer l'effet.

Les techniques d'*analyses locales*, mettant en jeu des *structures de graphes* (section 3.7) sont appropriées lorsqu'il existe des informations *a priori* ou externes sur les couples d'individus ou d'observations (existence d'une relation binaire symétrique ou structure de graphe non orienté décrivant des proximités temporelles ou géographiques). Sera évoqué ici le cas d'une variable nominale externe (partition *a priori* des individus donnant lieu à des analyses dites *intra* et *inter*), qui entre à la fois dans le cadre des sections 3.6 et 3.7.

Enfin les méthodes de traitement de *tableaux multiples* ou de *groupes de variables*<sup>2</sup> (section 3.8), qui correspondent à une famille quasi-illimitée de techniques, seront évoquées au travers d'une sélection des approches qui nous paraissent les plus utiles en pratique : analyse procrustéenne, méthode STATIS, analyse factorielle multiple, analyse canonique généralisée.

<sup>1</sup> Faut-il, pour un utilisateur dont la recherche statistique n'est pas l'activité principale, investir dans une méthode complexe qui ne servira qu'une fois? Vaut-il mieux utiliser une méthode de description un peu grossière, mais parfaitement dominée conceptuellement, en raison d'expériences accumulées, qu'une méthode plus subtile dont les résultats laissent perplexes? Le temps disponible, les possibilités de formation, les budgets d'acquisition de logiciels ne sont pas des ressources inépuisables.

<sup>2</sup> Notons que la section 3.6 traite un cas particulier de tableaux multiples: le couple  $(X, Z)$  est en effet un tableau avec deux groupes de variables.

## Analyse Canonique

La méthode d'analyse canonique développée par Hotelling (1936) constitue un cadre théorique général important dont la régression multiple et l'analyse discriminante, qui seront exposées plus loin, ainsi que l'analyse des correspondances, sont des cas particuliers. Sous sa forme générale, l'analyse canonique ne présente cependant qu'un intérêt assez limité pour les applications, car elle conduit à de grandes difficultés d'interprétation.

L'analyse canonique cherche à synthétiser les interrelations existant entre *deux groupes* de variables, en mettant en évidence les combinaisons linéaires des variables du premier groupe les plus *corrélées* à des combinaisons linéaires des variables du second groupe.

### 3.1.1 Formulation du problème et notations

Le tableau de données  $R$ , à  $n$  lignes et  $p+q$  colonnes, est partitionné en deux sous-tableaux  $X$  et  $Y$ , ayant respectivement  $p$  et  $q$  colonnes.

$$R = [X, Y]$$

Les lignes représentent les individus ou observations : les  $p$  premières colonnes sont les variables du premier groupe et les  $q$  suivantes sont celles du second groupe.

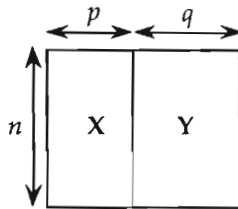


Figure 3.1 - 1  
Tableau des données  $R$

Nous supposons, sans perte de généralité, que les variables sont *centrées*, ce qui signifie que chaque colonne de  $R$  est telle que la somme de ses éléments vaut 0.

Alors la matrice des *covariances expérimentales* des  $p+q$  variables s'écrit :

$$V(R) = \frac{1}{n} R'R$$

Elle a pour terme général :

$$v_{jj'} = \frac{1}{n} \sum_i r_{ij} r_{ij'}$$

soit, en faisant apparaître les blocs :

$$V(\mathbf{R}) = \frac{1}{n} \begin{bmatrix} \mathbf{X}'\mathbf{X} & \mathbf{X}'\mathbf{Y} \\ \mathbf{Y}'\mathbf{X} & \mathbf{Y}'\mathbf{Y} \end{bmatrix}$$

Considérons l'individu  $i$ , caractérisé par la  $i^{\text{ème}}$  ligne de  $\mathbf{R}$  :

$$(x_{i1}, x_{i2}, \dots, x_{ip}, y_{i1}, y_{i2}, \dots, y_{iq})$$

Soient  $\mathbf{a}$  et  $\mathbf{b}$  deux vecteurs à  $p$  et  $q$  composantes, définissant deux combinaisons linéaires  $a(i)$  et  $b(i)$  :

$$a(i) = \sum_{j=1}^p a_j x_{ij} \qquad b(i) = \sum_{j=1}^q b_j y_{ij} .$$

Les  $n$  valeurs de  $a(i)$  pour tous les individus  $i$  sont les composantes de  $\mathbf{Xa}$ . De même, les  $n$  valeurs de  $b(i)$  sont les composantes de  $\mathbf{Yb}$ . Les vecteurs  $\mathbf{Xa}$  et  $\mathbf{Yb}$  représentent aussi deux points de  $\mathbb{R}^n$  appartenant aux sous-espaces  $V_X$  et  $V_Y$  engendrés par les colonnes de  $\mathbf{X}$  et  $\mathbf{Y}$ .

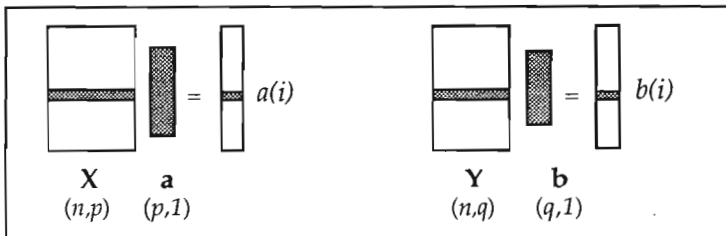


Figure 3.1 - 2  
Variables canoniques  $a(i)$  et  $b(i)$

Nous nous proposons de chercher les deux combinaisons linéaires  $a(i)$  et  $b(i)$  les plus corrélées sur l'ensemble des valeurs de  $i$ . Puisque les variables initiales sont centrées, leurs combinaisons linéaires sont également centrées.

Comme le coefficient de corrélation ne dépend pas de l'échelle des variables, nous imposerons aux deux combinaisons linéaires d'avoir une *variance unité*. La variance de l'ensemble des valeurs de  $a(i)$  pour  $i = 1, 2, \dots, n$  sera notée  $\text{var}(\mathbf{a})$ ; elle s'écrit :

$$\text{var}(\mathbf{a}) = \frac{1}{n} \sum_{i=1}^n a^2(i) = \frac{1}{n} (\mathbf{Xa})' \mathbf{Xa} = \frac{1}{n} \mathbf{a}' \mathbf{X}' \mathbf{Xa}$$

de la même façon :

$$\text{var}(\mathbf{b}) = \frac{1}{n} \mathbf{b}' \mathbf{Y}' \mathbf{Yb}$$

Dans ces conditions, le coefficient de corrélation entre les combinaisons linéaires  $a(i)$  et  $b(i)$  s'identifie avec la covariance :

$$\text{cov}(\mathbf{a}, \mathbf{b}) = \frac{1}{n} \sum_{i=1}^n a(i)b(i)$$

soit :

$$\text{cov}(\mathbf{a}, \mathbf{b}) = \frac{1}{n} \mathbf{a}' \mathbf{X}' \mathbf{Y} \mathbf{b}$$

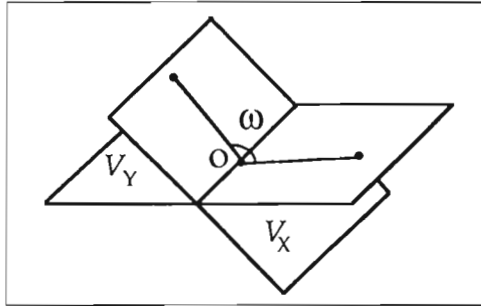


Figure 3.1 - 3  
Représentation géométrique des sous-espaces  $V_X$  et  $V_Y$

Finalement le problème de la recherche de la corrélation maximale s'écrit, après s'être affranchi des coefficients  $\frac{1}{n}$  (rappelons que  $\mathbf{X}$  et  $\mathbf{Y}$  sont centrés) :

- trouver  $\mathbf{a}$  et  $\mathbf{b}$  qui rendent maximal :  $\mathbf{a}' \mathbf{X}' \mathbf{Y} \mathbf{b}$
- avec les contraintes : 
$$\begin{cases} \mathbf{a}' \mathbf{X}' \mathbf{X} \mathbf{a} = 1 \\ \mathbf{b}' \mathbf{Y}' \mathbf{Y} \mathbf{b} = 1 \end{cases}$$

Les données étant centrées, le coefficient de corrélation n'est autre que le cosinus de l'angle entre les sous-espaces  $V_X$  et  $V_Y$ . La recherche des coefficients  $\mathbf{a}$  et  $\mathbf{b}$  revient donc à minimiser l'angle  $\omega$  entre les sous-espaces  $V_X$  et  $V_Y$ .

On appellera *variables canoniques* le couple  $(\mathbf{a}, \mathbf{b})$  ayant respectivement  $p$  et  $q$  composantes.

### 3.1.2 Les variables canoniques

#### a – Calcul des variables canoniques

La démonstration est analogue à celle rencontrée lors de l'analyse générale (§ 1.1.7). Deux multiplicateurs de Lagrange  $\lambda$  et  $\mu$  interviennent. Il faut rendre maximal :

$$\mathcal{L} = \mathbf{a}' \mathbf{X}' \mathbf{Y} \mathbf{b} - \lambda (\mathbf{a}' \mathbf{X}' \mathbf{X} \mathbf{a} - 1) - \mu (\mathbf{b}' \mathbf{Y}' \mathbf{Y} \mathbf{b} - 1)$$

L'annulation des dérivées de ce lagrangien par rapport aux vecteurs  $\mathbf{a}$  et  $\mathbf{b}$  conduit au système :

$$\begin{cases} \mathbf{X}'\mathbf{Y}\mathbf{b} - 2\lambda \mathbf{X}'\mathbf{X}\mathbf{a} = 0 \\ \mathbf{Y}'\mathbf{X}\mathbf{a} - 2\mu \mathbf{Y}'\mathbf{Y}\mathbf{b} = 0 \end{cases}$$

Prémultiplions les membres de ces deux relations respectivement par  $\mathbf{a}'$  et  $\mathbf{b}'$ . En tenant compte des contraintes :

$$\mathbf{a}'\mathbf{X}'\mathbf{X}\mathbf{a} = \mathbf{b}'\mathbf{Y}'\mathbf{Y}\mathbf{b} = 1$$

Elles se simplifient en :

$$\begin{cases} \mathbf{a}'\mathbf{X}'\mathbf{Y}\mathbf{b} = 2\lambda \\ \mathbf{b}'\mathbf{Y}'\mathbf{X}\mathbf{a} = 2\mu \end{cases}$$

Par conséquent  $\lambda = \mu$ . Nous poserons dorénavant :

$$\beta = 2\lambda$$

On remarquera que  $\beta$  est la valeur du *coefficient de corrélation maximal* recherché. Le système précédent s'écrit alors :

$$\begin{cases} \mathbf{X}'\mathbf{Y}\mathbf{b} = \beta \mathbf{X}'\mathbf{X}\mathbf{a} & [3.1 - 1] \\ \mathbf{Y}'\mathbf{X}\mathbf{a} = \beta \mathbf{Y}'\mathbf{Y}\mathbf{b} & [3.1 - 2] \end{cases}$$

La résolution est immédiate quand les matrices  $\mathbf{X}'\mathbf{X}$  et  $\mathbf{Y}'\mathbf{Y}$  sont *inversibles*. En reportant la valeur de  $\mathbf{a}$  tirée de [3.1 - 1] dans la relation [3.1 - 2] par exemple, on obtient :

$$\mathbf{Y}'\mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{Y}\mathbf{b} = \beta^2 \mathbf{Y}'\mathbf{Y}\mathbf{b} \quad [3.1 - 3]$$

Ceci montre que  $\mathbf{b}$  est *vecteur propre* de la matrice :

$$\mathbf{M} = (\mathbf{Y}'\mathbf{Y})^{-1}\mathbf{Y}'\mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{Y}$$

relatif à la plus grande *valeur propre* notée  $\beta^2$ , carré du coefficient de corrélation entre les combinaisons linéaires  $\mathbf{a}$  et  $\mathbf{b}$  et carré du cosinus maximum entre les sous-espaces  $V_X$  et  $V_Y$ . Cette valeur  $\beta^2$  est la première *racine canonique*, ou carré du premier *coefficient de corrélation canonique* entre les deux variables.

De façon analogue, on calcule  $\mathbf{a}$  à partir de la relation [3.1 - 1] ou en considérant directement  $\mathbf{a}$  comme vecteur propre de :

$$\mathbf{N} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{Y}(\mathbf{Y}'\mathbf{Y})^{-1}\mathbf{Y}'\mathbf{X} \quad [3.1 - 4]$$

Si  $\mathbf{X}$  est de plein rang, alors  $\mathbf{X}'\mathbf{X}$  est inversible et la relation [3.1 - 1] permet d'écrire :

$$\mathbf{a} = \frac{1}{\beta} (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{Y}\mathbf{b}$$

Un raisonnement analogue à celui fait lors de l'analyse générale nous permettrait de généraliser le résultat à la recherche des  $r$  variables

canoniques,  $r$  étant le plus petit des deux entiers  $p$  et  $q$  : les  $r$  vecteurs propres successifs, dans l'ordre des valeurs propres décroissantes, correspondent aux couples de combinaisons linéaires de chaque ensemble les plus corrélées entre elles, les combinaisons linéaires successives relatives à un même ensemble étant assujetties à être non corrélées.

### b – Interprétation géométrique

Les relations [3.1 - 1] et [3.1 - 2] peuvent s'écrire :

$$\mathbf{a} = \frac{1}{\beta} (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}'\mathbf{Y}\mathbf{b}, \quad \text{et} \quad \mathbf{b} = \frac{1}{\beta} (\mathbf{Y}'\mathbf{Y})^{-1} \mathbf{Y}'\mathbf{X}\mathbf{a}$$

Prémultipliant les deux membres de chacune d'elles respectivement par  $\mathbf{X}$  et  $\mathbf{Y}$  on obtient :

$$\mathbf{X}\mathbf{a} = \frac{1}{\beta} \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}'\mathbf{Y}\mathbf{b} \quad [3.1 - 5]$$

$$\mathbf{Y}\mathbf{b} = \frac{1}{\beta} \mathbf{Y}(\mathbf{Y}'\mathbf{Y})^{-1} \mathbf{Y}'\mathbf{X}\mathbf{a} \quad [3.1 - 6]$$

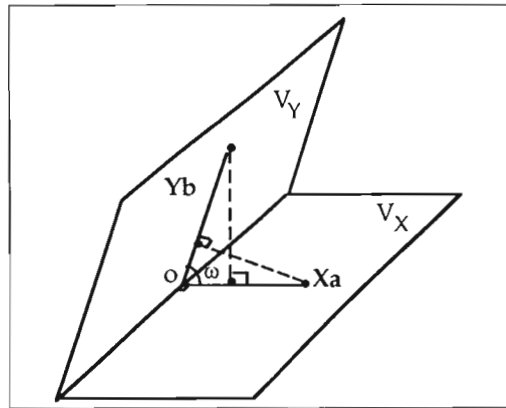


Figure 3.1 - 4  
Interprétation géométrique de l'analyse canonique

Les matrices symétriques et idempotentes :

$$\mathbf{P}_X = \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}' \quad \text{et} \quad \mathbf{P}_Y = \mathbf{Y}(\mathbf{Y}'\mathbf{Y})^{-1} \mathbf{Y}'$$

sont les *opérateurs de projection orthogonale* respectivement sur les sous-espaces  $V_X$  et  $V_Y$ .

Autrement dit les relations [3.1 - 5] et [3.1 - 6] expriment que chacun des vecteurs  $\mathbf{X}\mathbf{a}$  et  $\mathbf{Y}\mathbf{b}$  est colinéaire à la projection de l'autre.

Les vecteurs  $\mathbf{X}\mathbf{a}$  et  $\mathbf{Y}\mathbf{b}$  étant unitaires, les formules montrent en effet que :

$$\beta = \cos(\omega) = \cos(\mathbf{X}\mathbf{a}, \mathbf{Y}\mathbf{b})$$



Il apparaît que la première racine canonique  $\beta^2$  est le carré du cosinus du plus petit angle<sup>1</sup> entre les sous-espaces  $V_X$  et  $V_Y$ .

### c – Cas de matrices non inversibles

Examinons le cas où les matrices  $X'X$  ou  $Y'Y$  sont singulières. Prenons  $Y'Y$  pour fixer les idées. Cela signifie que la matrice  $Y$  d'ordre  $(n, q)$  a un rang inférieur à  $q$  : soit  $q - s$  son rang.

Il y a deux façons de procéder pour résoudre le système des équations matricielles [3.1 - 1] et [3.1 - 2] :

- on prend dans  $\mathbb{R}^n$  une base du sous-espace  $V_Y$  à  $q - s$  dimensions engendrée par  $Y$ , base décrite par les  $q - s$  colonnes d'une matrice<sup>2</sup>  $\hat{Y}$ ; à  $Yb$  on substitue dans les calculs  $\hat{Y}\hat{b}$  où  $\hat{b}$  est un vecteur à  $q - s$  composantes. La matrice  $\hat{Y}'\hat{Y}$  est maintenant inversible.
- Comme cela est fréquent dans le cas du modèle linéaire général, on construit une matrice  $Y_0$  de plein rang d'ordre  $(n, q)$ , telle que  $V_Y \subset V_{Y_0}$ . Pour retrouver le sous-espace  $V_Y$ , il est alors nécessaire d'imposer à  $b$  une contrainte, à savoir :  $Y_0b$  devra appartenir à  $V_Y$ . Si  $Y_1$  désigne une matrice d'ordre  $(n, s)$ , telle que  $Y_1'Y = 0$  et que  $Y_1b \in V_{Y_0}$ , la contrainte sur  $b$  s'écrira :

$$Y_1'Y_0b = 0$$

#### Remarque :

Cette situation se présentera également en *analyse discriminante* dans un contexte simple : la matrice  $Y$  d'ordre  $(n, q)$  est singulière, alors que la matrice initiale  $Y_0$  (avant centrage) est de plein rang. Ceci résulte du fait que le sous-espace  $V_{Y_0}$  engendrée par  $Y_0$  contient le vecteur  $e_n$  de  $\mathbb{R}^n$  dont toutes les composantes valent 1. On travaillera alors avec la matrice  $Y_0$  sachant que  $b$  est assujéti à vérifier :

$$e_n'Y_0b = 0$$

relation qui s'écrit :

$$\sum_{j=1}^q y_{.j}b_j = 0$$

( $y_{.j}$  désignant la somme de la colonne  $j$  de la matrice  $Y_0$ ).

<sup>1</sup> Notons que ces considérations géométriques nous auraient permis d'écrire *directement* les formules [3.1 - 5] et [3.1 - 6], et donc de procéder au calcul des variables canoniques : on remplace, par exemple dans la relation [3.1 - 6],  $Xa$  par sa valeur tirée de la relation [3.1 - 5].

<sup>2</sup> On choisira de préférence une base orthogonale, obtenue, par exemple, par le procédé d'orthogonalisation de Gram-Schmidt, ou une base issue d'une analyse générale de  $Y$ .

### 3.1.3 Liens avec l'analyse des correspondances

L'analyse canonique contient comme cas particulier l'analyse des correspondances et peut se généraliser au cas de plus de deux variables nominales.

En reprenant les notations de l'analyse des correspondances multiples (section 1.4), le tableau de données  $\mathbf{R} = [\mathbf{Z}_1, \dots, \mathbf{Z}_q, \dots, \mathbf{Z}_s]$  à  $n$  lignes et  $p$  colonnes est le tableau disjonctif complet juxtaposant  $s$  sous-tableaux. Chaque sous-tableau  $\mathbf{Z}_q$  correspond à une question  $q$  totalisant  $p_q$  modalités de réponses et engendre, dans l'espace  $\mathbb{R}^p$ , un sous-espace  $V_{Z_q}$  à  $p_q$  dimensions<sup>1</sup>.

#### a – Le cas de l'analyse des correspondances simples

L'analyse des correspondances du tableau de contingence croisant deux variables  $q$  et  $q'$  revient à étudier les positions relatives des sous-espaces  $V_{Z_q}$  et  $V_{Z_{q'}}$ . C'est l'analyse canonique du tableau  $[\mathbf{Z}_q, \mathbf{Z}_{q'}]$ .

Soit  $\phi_q$  le vecteur dont les  $p_q$  composantes sont les coordonnées d'un point  $\mathbf{m}_q$  de  $V_{Z_q}$  dans la base définie par les colonnes de  $\mathbf{Z}_q$ . Les coordonnées de  $\mathbf{m}_q$  dans  $\mathbb{R}^n$  sont les composantes de  $\mathbf{m}_q = \mathbf{Z}_q \phi_q$ .

Le carré de la distance de ce point  $\mathbf{m}_q$  à l'origine, selon la norme euclidienne usuelle, n'est autre que :

$$\phi_q' \mathbf{Z}_q' \mathbf{Z}_q \phi_q = \phi_q' \mathbf{D}_q \phi_q$$

Les relations de double transition [1.4 - 7] et [1.4 - 8] s'écrivent ici (en omettant l'indice  $\alpha$  de l'axe pour alléger les notations) :

$$\begin{cases} \phi_q = \frac{1}{\sqrt{\lambda}} \mathbf{D}_q^{-1} \mathbf{Z}_q' \mathbf{Z}_{q'} \phi_{q'} \\ \phi_{q'} = \frac{1}{\sqrt{\lambda}} \mathbf{D}_{q'}^{-1} \mathbf{Z}_{q'}' \mathbf{Z}_q \phi_q \end{cases}$$

On en déduit le système suivant :

$$\begin{cases} \mathbf{Z}_q \phi_q = \frac{1}{\sqrt{\lambda}} \mathbf{Z}_q \mathbf{D}_q^{-1} \mathbf{Z}_q' \mathbf{Z}_{q'} \phi_{q'} \\ \mathbf{Z}_{q'} \phi_{q'} = \frac{1}{\sqrt{\lambda}} \mathbf{Z}_{q'} \mathbf{D}_{q'}^{-1} \mathbf{Z}_{q'}' \mathbf{Z}_q \phi_q \end{cases}$$

soit :

<sup>1</sup> Rappelons que les  $s$  sous-espaces ont en commun au moins la première bissectrice. Le rang de  $\mathbf{R}$  est donc au plus égal à  $p - s + 1$ .

$$\mathbf{m}_q = \frac{1}{\sqrt{\lambda}} \mathbf{P}_q \mathbf{m}_{q'} \quad [3.1 - 7]$$

$$\mathbf{m}_{q'} = \frac{1}{\sqrt{\lambda}} \mathbf{P}_{q'} \mathbf{m}_q \quad [3.1 - 8]$$

où :

$$\mathbf{P}_q = \mathbf{Z}_q (\mathbf{Z}'_q \mathbf{Z}_q)^{-1} \mathbf{Z}'_q \quad \text{et} \quad \mathbf{P}_{q'} = \mathbf{Z}_{q'} (\mathbf{Z}'_{q'} \mathbf{Z}_{q'})^{-1} \mathbf{Z}'_{q'}$$

Les matrices  $\mathbf{P}_q$  et  $\mathbf{P}_{q'}$  représentent respectivement les opérateurs projection sur les sous-espaces  $V_{Z_q}$  et  $V_{Z_{q'}}$ .

Les relations [3.1 - 7] et [3.1 - 8] expriment que la projection orthogonale de  $\mathbf{m}_q$  sur  $V_{Z_{q'}}$  est colinéaire à  $\mathbf{m}_{q'}$  (et semblablement pour  $\mathbf{m}_{q'}$  sur  $V_{Z_q}$ ).

Présentée comme la recherche des plus petits angles entre deux sous-espaces  $V_{Z_q}$  et  $V_{Z_{q'}}$ , l'analyse canonique ne se généralise pas facilement au cas de plus de deux questions<sup>1</sup>.

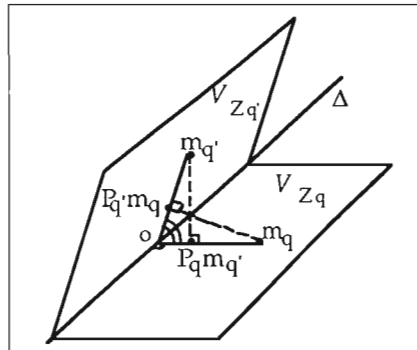


Figure 3.1 - 2  
Projections sur  $V_{Z_q}$  et  $V_{Z_{q'}}$ .

Mais une autre formulation va permettre de présenter l'analyse des correspondances multiples comme une analyse canonique généralisée particulière.

## b – L'analyse des correspondances multiples

L'analyse canonique du tableau  $[\mathbf{Z}_q, \mathbf{Z}_{q'}]$  peut aussi se formuler de la façon suivante :

<sup>1</sup> On reviendra sur ce lien entre analyse de correspondances et analyse canonique au paragraphe 3.3.4.b, à propos de l'analyse factorielle discriminante, qui est elle aussi une analyse canonique particulière.

trouver deux points  $\mathbf{m}_q$  et  $\mathbf{m}_{q'}$  tels que la somme des carrés de leurs distances à l'origine soit constante :

$$\varphi'_q \mathbf{D}_q \varphi_q + \varphi'_{q'} \mathbf{D}_{q'} \varphi_{q'} = 2n \quad [3.1 - 9]$$

et tels que la distance à l'origine du point  $\mathbf{m} = \mathbf{m}_q + \mathbf{m}_{q'}$  soit maximale.

En effet, cette distance a pour carré :

$$\|\mathbf{m}\|^2 = \varphi'_q \mathbf{D}_q \varphi_q + \varphi'_{q'} \mathbf{D}_{q'} \varphi_{q'} + 2\varphi'_q \mathbf{Z}'_q \mathbf{Z}_{q'} \varphi_{q'}$$

soit :

$$\|\mathbf{m}\|^2 = 2n \left( 1 + \frac{1}{n} \varphi'_q \mathbf{Z}'_q \mathbf{Z}_{q'} \varphi_{q'} \right)$$

Rendre maximale  $\|\mathbf{m}\|^2$  avec la contrainte [3.1 - 9], ou avec les deux contraintes :

$$\varphi'_q \mathbf{D}_q \varphi_q = \varphi'_{q'} \mathbf{D}_{q'} \varphi_{q'} = n$$

conduit au même résultat<sup>1</sup>.

Avec la contrainte unique [3.1 - 9], le problème se généralise aisément au cas de plus de deux questions.

On désigne par  $\varphi_1, \dots, \varphi_q, \dots, \varphi_s$  respectivement les vecteurs des composantes de  $s$  points  $\mathbf{m}_1, \dots, \mathbf{m}_q, \dots, \mathbf{m}_s$  dans les bases  $\mathbf{Z}_1, \dots, \mathbf{Z}_q, \dots, \mathbf{Z}_s$  et soit  $\mathbf{m} = \mathbf{m}_1 + \mathbf{m}_q + \mathbf{m}_s$ .

On cherchera à rendre maximale la quantité :

$$\|\mathbf{m}\|^2 = \sum_{q \in S} \sum_{q' \in S} \varphi'_q \mathbf{Z}'_q \mathbf{Z}_{q'} \varphi_{q'}$$

avec la contrainte :

$$\sum_{q \in S} \varphi'_q \mathbf{D}_q \varphi_q = sn$$

Si  $\Phi$  désigne le vecteur à  $p$  composantes défini par :

$$\Phi' = \{\varphi'_1, \dots, \varphi'_q, \dots, \varphi'_s\}$$

le problème revient à rendre maximal :

$$\Phi' \mathbf{B} \Phi$$

avec la contrainte :

$$\Phi' \mathbf{D} \Phi = sn$$

où l'on rappelle que  $\mathbf{B}$  est le tableau de contingence de Burt obtenu à partir du tableau disjonctif complet.

<sup>1</sup> En effet, les multiplicateurs de Lagrange relatifs à ces deux dernières contraintes sont égaux.

Les facteurs  $\Phi$  cherchés sont donc les vecteurs propres de  $\mathbf{D}^{-1}\mathbf{B}$  relatifs aux plus grandes valeurs propres.

Il s'agit d'une généralisation simple de l'analyse canonique au cas de plus de deux ensembles : elle conduit à une diagonalisation de matrice symétrique, opération classique et maîtrisée<sup>1</sup>.

Les autres méthodes (introduction de  $s$  contraintes au lieu d'une seule) demandent des algorithmes itératifs assez coûteux et ne conduisent pas à des règles d'interprétation simples.

---

<sup>1</sup> Cette extension de l'analyse canonique sera présentée à nouveau dans un cadre plus général au paragraphe 3.8.5.

## Régression multiple, modèle linéaire

La régression multiple vise à expliquer ou prédire une variable continue (dite variable dépendante ou à expliquer ou encore endogène) à l'aide d'un ensemble de variables dites explicatives (ou exogènes). On réserve en général le nom de régression multiple au cas où les variables explicatives sont continues. Lorsque celles-ci sont des variables nominales, on parle d'analyse de la variance et pour un ensemble de variables mixtes, d'analyse de la covariance. La théorie statistique qui englobe ces diverses techniques constitue le *modèle linéaire*.

La régression constitue sans doute la méthode statistique la plus utilisée bien que sa portée et ses limites ne soient pas toujours bien connues. De ce fait, elle n'est pas toujours pratiquée à bon escient. La littérature sur la régression et le modèle linéaire est extrêmement abondante. C'est en économétrie, champ d'application privilégié du modèle linéaire, que l'on trouve les premiers manuels généraux en langue française exposant les méthodes et les principaux types de résultats (Malinvaud, 1964; Fourgeaud et al., 1978). On citera également l'ouvrage de Tomassone et al. (1983), exposé complet, simple et opérationnel sur tous les aspects de la régression. Pour un exposé plus concis, on renverra à Saporta (1990). Mais ces quelques titres<sup>1</sup> ne sauraient rendre justice de la profusion des excellents manuels sur ce sujet.

### 3.2.1 Formulation du problème : le modèle linéaire

On dispose d'un ensemble de  $n$  observations sur lesquelles ont été effectuées  $p+1$  mesures des variables  $y, x_1, x_2, \dots, x_p$ . On veut expliquer ou prévoir  $y$  à l'aide des variables explicatives ou prédicteurs,  $x_1, x_2, \dots, x_p$ , lesquels sont supposés connus sans erreur.

---

<sup>1</sup> La littérature en anglais sur le modèle linéaire est particulièrement vaste : on trouvera une bibliographie commentée (déjà ancienne) de plusieurs centaines d'articles et ouvrages dans Harter (1974-1975). Searle (1971) et Seber (1977) traitent de façon extensive les problèmes d'analyse de la variance et de covariance; Theil (1971) situe le modèle linéaire dans un cadre économétrique général; l'ouvrage de Rao (1973), réédition d'un manuel classique, est consacré à l'opération d'induction statistique sur le modèle linéaire. Un autre manuel classique est l'ouvrage de Draper et Smith (1981). Mosteller et Tukey (1977), Besley et al (1980), Atkinson (1985) présentent des points de vue un peu plus modernes, incluant diverses méthodes de sélection de variables, alors que Chatterjee et Price (1991) insistent sur la mise en oeuvre pratique.

Supposons par exemple qu'une personne désire acquérir un magasin ayant une surface  $S$  dans une zone où la population environnante est  $P$ . Des études antérieures montrent que le chiffre d'affaires d'un magasin dépend linéairement de la surface et de la population, et les données relatives à 30 magasins du même type sont disponibles. Quel chiffre d'affaires peut espérer l'acheteur ? Le chiffre d'affaires est la variable à prévoir et les variables explicatives ou prédicteurs sont la population et la surface. Ce type de problème trouve une solution dans le cadre de la régression, technique de prévision linéaire, qui consiste tout d'abord à procéder à une estimation d'un modèle, puis à utiliser le modèle estimé pour le calcul de la valeur attendue.

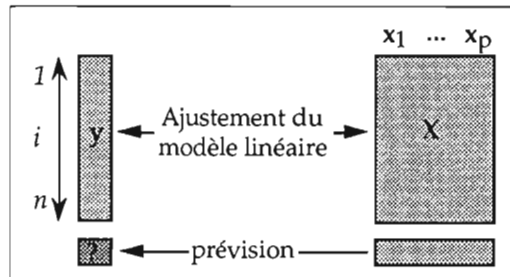


Figure 3.2 - 1  
Prévision linéaire

On cherche à approcher  $y$  par une combinaison linéaire des variables explicatives  $x_1, x_2, \dots, x_p$ . Pour cela, on pose le modèle<sup>1</sup> :

$$y_i = \alpha_0 + \alpha_1 x_{i1} + \alpha_2 x_{i2} + \dots + \alpha_p x_{ip} + \varepsilon_i$$

où  $\alpha_0, \alpha_1, \alpha_2, \dots, \alpha_p$  sont les coefficients inconnus du modèle. Le terme constant  $\alpha_0$  peut être considéré comme coefficient d'une variable explicative particulière artificielle  $x_0$  dont les valeurs  $x_{i0}$  seraient toujours égales à 1.  $\varepsilon_i$  est le résidu représentant l'écart entre la valeur observée  $y_i$  et la partie "expliquée" de l'observation  $(\alpha_0 + \alpha_1 x_{i1} + \alpha_2 x_{i2} + \dots + \alpha_p x_{ip})$ .

On suppose dans la plupart des spécifications du modèle que tous les résidus  $\varepsilon_i$  sont des quantités aléatoires indépendantes.

Ce modèle s'exprime sous forme matricielle :

$$\underset{(n,1)}{y} = \underset{(n,p+1)}{X} \underset{(p+1,1)}{\alpha} + \underset{(1,n)}{\varepsilon}$$

<sup>1</sup> La linéarité des relations par rapport aux coefficients  $\alpha_0, \alpha_1, \alpha_2, \dots, \alpha_p$  peut n'apparaître qu'après transformations des données. Par exemple :

$$y = \alpha_3 x_1^{\alpha_1} x_2^{\alpha_2} (1 + \varepsilon)$$

deviendra linéaire après la transformation logarithmique :

$$\log(y) = \alpha_1 \log(x_1) + \alpha_2 \log(x_2) + \log(\alpha_3) + \log(1 + \varepsilon)$$

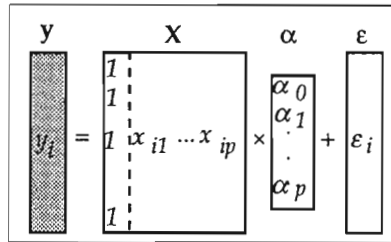


Figure 3.2 - 2  
Schématisation du modèle linéaire

On dispose, pour évaluer les coefficients inconnus du modèle, d'un système de  $n$  équations linéaires ayant  $n + p + 1$  inconnues. Le système admet donc une infinité de solutions.

Soient  $a_0, a_1, a_2, \dots, a_p$  les coefficients correspondant à une des solutions possibles. On cherchera la solution qui minimise globalement, suivant un critère à définir, l'ensemble des écarts à la linéarité, c'est-à-dire :

$$\left\{ \begin{array}{l} \text{choisir } (a_0, a_1, a_2, \dots, a_p) \text{ qui minimisent l'ensemble des } e_i \\ \text{avec } e_i = y_i - (a_0 + a_1 x_{i1} + a_2 x_{i2} + \dots + a_p x_{ip}) \end{array} \right.$$

Parmi les critères possibles de minimisation, citons la méthode des moindres carrés  $\min(\sum e_i^2)$  (norme dite "L<sub>2</sub>") celle des moindres valeurs absolues  $\min(\sum |e_i|)$  (norme dite "L<sub>1</sub>") , celle du minimax  $\min(\max_{(i)} e_i)$  (norme dite "L<sub>∞</sub>")<sup>1</sup>. Le critère des moindres carrés s'avère conduire à des calculs algébriques simples, se prêter à une interprétation géométrique claire, et donner lieu à des interprétations statistiques intéressantes<sup>2</sup>.

### 3.2.2 Ajustement par la méthode des moindres-carrés

On appelle ajustement du modèle linéaire toute solution du système d'équations :

$$y_i = a_0 + a_1 x_{i1} + a_2 x_{i2} + \dots + a_p x_{ip} + e_i \quad (i = 1, 2, \dots, n)$$

ce qui correspond sous forme matricielle à :

$$\underset{(n,1)}{\mathbf{y}} = \underset{(n,p+1)}{\mathbf{X}} \underset{(p+1,1)}{\mathbf{a}} + \underset{(1,n)}{\mathbf{e}}$$

<sup>1</sup> Plus généralement, la norme L<sub>k</sub> correspond au critère  $\min(\sum |e_i|^k)$

<sup>2</sup> La norme L<sub>1</sub>, qui privilégie moins les écarts importants, est à la base de méthodes de régression plus *robustes* (cf. Huber, 1981; 1987). Sur le rôle de cette norme en analyse descriptive des données, cf. Fichet (1987), et Le Calvé (1987). L'utilisation de la norme L<sub>1</sub> dans le cas de la régression linéaire remonte à Laplace (1793). Une étude historique de l'utilisation des normes L<sub>1</sub> et L<sub>∞</sub> a été réalisée par Farebrother (1987).



Pour la  $i^{\text{ème}}$  observation, la valeur prédite par le modèle est :

$$\tilde{y}_i = a_0 + a_1x_{i1} + a_2x_{i2} + \dots + a_px_{ip}$$

le résidu du modèle correspondant vaut donc<sup>1</sup> :

$$e_i = y_i - \tilde{y}_i$$

D'une manière générale, on cherche  $\tilde{y}$  le plus proche possible de  $y$  :

$$\tilde{y} = Xa = a_0x_0 + a_1x_1 + a_2x_2 + \dots + a_px_p$$

L'ajustement par la méthode des moindres carrés est celui qui fournit les coefficients  $a_0, a_1, a_2, \dots, a_p$  conduisant au minimum de la somme des carrés des écarts :

$$\min(\sum e_i^2)$$

Dans la suite, nous allons supposer que les variables sont centrées, ce qui implique  $a_0 = 0$ . Une des propriétés de la régression multiple est que les estimations des coefficients autres que  $a_0$  sont les mêmes, que les variables soient centrées *a priori* ou pas.

### a – Calcul et propriétés de l'ajustement des moindres-carrés

Il s'agit de déterminer le vecteur  $a$  des coefficients qui minimise :

$$e'e = \sum e_i^2 = \|y - \tilde{y}\|^2$$

Le vecteur de coefficients  $a$  doit vérifier la condition d'extremum<sup>2</sup> :

$$X'Xa = X'y \quad [3.2 - 1]$$

qui est un système de  $p$  équations à  $p$  inconnues.

Si  $n$  est supérieur ou égal à  $p$  (plus d'équations que d'inconnues) et si  $X$  est de plein rang (c'est-à-dire de rang  $p$ ), alors  $X'X$  est inversible.

On tire de la relation [3.2 - 1] la solution :

$$a = (X'X)^{-1}X'y \quad [3.2 - 2]$$

<sup>1</sup> Le vocabulaire et les notations distinguent les résidus définis par le modèle théorique  $\varepsilon_i = y_i - \sum_k \alpha_k x_{ik}$  et les écarts définis par un ajustement  $e_i = y_i - \sum_k a_k x_{ik}$

<sup>2</sup> La quantité scalaire  $e'e$  étant une fonction des inconnues ( $a_1, a_2, \dots, a_p$ ), une condition nécessaire d'extremum est l'annulation des dérivées partielles premières, soit :

$$\frac{\partial}{\partial a} (e'e) = \underset{(p,1)}{0}$$

on a : 
$$e'e = (y - Xa)'(y - Xa) = y'y - 2a'X'y + a'X'Xa$$

d'où : 
$$\frac{\partial}{\partial a} (e'e) = -2X'y + 2X'Xa$$

on en tire la condition d'extremum :  $X'Xa = X'y$

Le vecteur  $\mathbf{a}$  est le vecteur des coefficients de régression multiple<sup>1</sup>.

Il reste à vérifier que l'extremum atteint par  $\mathbf{e}'\mathbf{e}$  est bien un minimum.

Soit  $\tilde{\mathbf{a}}$  une autre solution et  $\tilde{\mathbf{e}}$  le vecteur correspondant des écarts :

$$\tilde{\mathbf{e}} = \mathbf{y} - \mathbf{X}\tilde{\mathbf{a}} = (\mathbf{y} - \mathbf{X}\mathbf{a}) + (\mathbf{X}\mathbf{a} - \mathbf{X}\tilde{\mathbf{a}}) = \mathbf{e} + \mathbf{X}(\mathbf{a} - \tilde{\mathbf{a}})$$

et

$$\tilde{\mathbf{e}}'\tilde{\mathbf{e}} = \mathbf{e}'\mathbf{e} + 2(\mathbf{a} - \tilde{\mathbf{a}})'\mathbf{X}'(\mathbf{y} - \mathbf{X}\mathbf{a}) + (\mathbf{a} - \tilde{\mathbf{a}})'\mathbf{X}'\mathbf{X}(\mathbf{a} - \tilde{\mathbf{a}})$$

Dans le membre de droite, le terme central est nul d'après [3.2 - 1]; il reste donc :

$$\tilde{\mathbf{e}}'\tilde{\mathbf{e}} = \mathbf{e}'\mathbf{e} + (\mathbf{X}(\mathbf{a} - \tilde{\mathbf{a}}))'(\mathbf{X}(\mathbf{a} - \tilde{\mathbf{a}}))$$

Il est clair que le dernier terme est une somme de carrés et ne peut être que positif ou nul. Par conséquent  $\mathbf{e}'\mathbf{e}$  est bien la plus petite somme de carrés d'écarts.

### b – Approche géométrique dans $\mathbb{R}^n$

Les propriétés algébriques de l'ajustement vont nous permettre d'interpréter géométriquement l'opération effectuée.

Plaçons-nous dans l'espace  $\mathbb{R}^n$  où  $n$  est le nombre des observations effectuées sur  $p+1$  variables :  $y, x_1, x_2, \dots, x_p$ .

La recherche de  $\mathbf{y}$  comme combinaison linéaire des  $x_1, x_2, \dots, x_p$  revient à définir  $\tilde{\mathbf{y}}$  dans le sous-espace engendré par les variables explicatives  $V_X$ . La technique d'ajustement des moindres-carrés consiste alors à approcher  $\mathbf{y}$  par sa projection orthogonale  $\tilde{\mathbf{y}}$  sur le sous-espace  $V_X$ .

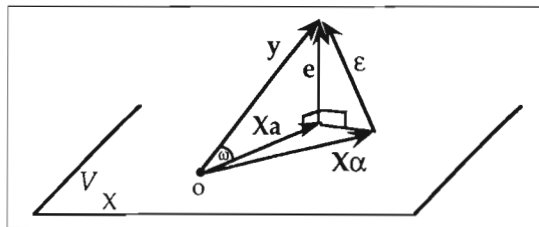


Figure 3.2 - 3  
Projection de  $\mathbf{y}$  sur  $V_X$

En remplaçant  $\mathbf{a}$  par sa valeur obtenue dans [3.2 - 2], on obtient :

$$\tilde{\mathbf{y}} = \mathbf{X}\mathbf{a} = \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y} = \mathbf{P}_X\mathbf{y}$$

<sup>1</sup> La régression simple correspond au modèle  $y = \alpha x + \varepsilon$  (une seule variable explicative,  $y$  et  $x$  centrés). La formule [3.2 - 2] devient  $a = x'y/x'x$  ou  $a = cov(x,y)/var(x)$ .

avec :

$$\mathbf{P}_X = \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}' \quad [3.2-3]$$

où la matrice  $\mathbf{P}_X$  désigne l'opérateur de projection orthogonale<sup>1</sup> sur  $V_X$ . Comme le montre la figure 3.2 - 3, le modèle théorique  $\mathbf{y} = \mathbf{X}\boldsymbol{\alpha} + \boldsymbol{\varepsilon}$  définit une décomposition de  $\mathbf{y}$  en deux termes inconnus, l'un  $\mathbf{X}\boldsymbol{\alpha}$  dans  $V_X$  et l'autre  $\boldsymbol{\varepsilon}$  dans  $\mathbb{R}^n$ . La technique des moindres-carrés propose pour solution la décomposition  $\mathbf{y} = \mathbf{X}\mathbf{a} + \mathbf{e}$  qui minimise la "longueur" de  $\mathbf{e}$  en projetant orthogonalement  $\mathbf{y}$  en  $\mathbf{X}\mathbf{a}$  sur  $V_X$  et  $\boldsymbol{\varepsilon}$  en  $\mathbf{e}$  sur le sous-espace orthogonal à  $V_X$  dans  $\mathbb{R}^n$ . Les deux vecteurs  $\mathbf{X}\mathbf{a}$  et  $\mathbf{e}$  sont orthogonaux.

### c – Le coefficient de corrélation multiple

Remarquons que les variables étant centrées, les longueurs dans l'espace  $\mathbb{R}^n$  s'interprètent en termes de variances. Le théorème de Pythagore appliqué au triangle rectangle de la figure 3.2 - 3 dont les côtés sont  $\mathbf{e}$  et  $\mathbf{X}\mathbf{a}$  et l'hypoténuse  $\mathbf{y}$ , peut s'écrire :

$$\mathbf{y}'\mathbf{y} = \mathbf{e}'\mathbf{e} + \mathbf{a}'\mathbf{X}'\mathbf{X}\mathbf{a}$$

En divisant par  $n$  chacun de ces termes, on obtient la relation :

$$\frac{1}{n} \sum (y_i)^2 = \frac{1}{n} \sum (y_i - \bar{y})^2 + \frac{1}{n} \sum (\tilde{y}_i)^2$$

*variance*
*variance*
*variance*  
*totale*
*résiduelle*
*expliquée*

Afin d'avoir une idée globale de la qualité de l'ajustement, on définit le coefficient de corrélation multiple  $R$  comme le cosinus de l'angle  $\omega$  entre  $\mathbf{y}$  et  $\mathbf{X}\mathbf{a}$  qui n'est autre que le coefficient de corrélation entre les valeurs initiales et les valeurs ajustées :

$$R = \text{cor}(\mathbf{y}, \tilde{\mathbf{y}}) = \text{cor}(\mathbf{y}, \mathbf{X}\mathbf{a}).$$

Son carré peut s'exprimer sous différentes formes :

$$R^2 = \frac{\text{cov}^2(\mathbf{y}, \tilde{\mathbf{y}})}{\text{var}(\mathbf{y})\text{var}(\tilde{\mathbf{y}})} = \frac{\text{var}(\tilde{\mathbf{y}})}{\text{var}(\mathbf{y})} = \frac{\sum (\tilde{y}_i)^2}{\sum (y_i)^2} = \frac{\text{variance expliquée}}{\text{variance totale}}.$$

De façon explicite en fonction des données initiales  $\mathbf{X}$  et  $\mathbf{y}$ ,  $R^2$  s'écrit :

$$R^2 = \frac{\mathbf{a}'\mathbf{X}'\mathbf{X}\mathbf{a}}{\mathbf{y}'\mathbf{y}} = \frac{\mathbf{y}'\mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y}}{\mathbf{y}'\mathbf{y}}$$

Ce coefficient décrit donc le partage de la variance totale en variance "expliquée" et "résiduelle" :

<sup>1</sup> Cet opérateur, symétrique et idempotent, a déjà été rencontré à propos de l'analyse canonique (cf. § 3.1.2.b).

$$\left\{ \begin{array}{l} \text{variance expliquée} \quad R^2 \text{var}(\mathbf{y}) = \text{var}(\bar{\mathbf{y}}) \\ \text{variance résiduelle} \quad (1 - R^2) \text{var}(\mathbf{y}) = \text{var}(\mathbf{e}) \\ \hline \text{variance totale} \quad \text{var}(\mathbf{y}) = \text{var}(\bar{\mathbf{y}}) + \text{var}(\mathbf{e}) \end{array} \right.$$

Ainsi, en minimisant  $\sum e_i^2$ , on maximise  $R^2$ . En d'autres termes, l'ajustement des moindres-carrés détermine la combinaison linéaire des variables explicatives ayant une corrélation maximale<sup>1</sup> avec la variable à expliquer  $\mathbf{y}$ .

### 3.2.3 Lien avec l'analyse canonique

La régression multiple est un cas particulier de l'analyse canonique quand la matrice  $\mathbf{Y}$  n'a qu'une colonne  $\mathbf{y}$  ( $q = 1$ ), et donc le sous-espace  $V_{\mathbf{Y}}$  est réduit à une droite. La variable canonique  $\mathbf{b}$  n'a alors qu'une composante notée  $b$ . Le produit  $\mathbf{y}'\mathbf{y}$  étant maintenant un scalaire, la relation [3.1-3] (cf. § 3.1.2.a) devient :

$$\beta^2 = \frac{\mathbf{y}'\mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y}}{\mathbf{y}'\mathbf{y}}$$

L'unique racine canonique  $\beta^2$  est le carré du coefficient de corrélation multiple entre la colonne  $\mathbf{y}$  et les colonnes de  $\mathbf{X}$  c'est-à-dire entre la variable à expliquer et les variables explicatives.

Compte tenu de la relation [3.1-1], la variable canonique  $\mathbf{a}$  s'écrit :

$$\mathbf{a} = \frac{b}{\beta} (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}'\mathbf{y}$$

Cette relation montre que le vecteur  $\mathbf{a}$  est proportionnel (au coefficient  $\frac{b}{\beta}$  près) au vecteur des coefficients de la régression multiple expliquant la variable  $\mathbf{y}$  par les  $p$  variables colonnes de  $\mathbf{X}$ .

Le coefficient  $\frac{b}{\beta}$  est d'ailleurs facile à calculer puisque, d'après la contrainte de normalisation,  $b = \frac{1}{\sqrt{\mathbf{y}'\mathbf{y}}}$ .

<sup>1</sup> On remarquera par ailleurs que l'introduction dans le modèle d'une nouvelle variable explicative quelconque ne peut que diminuer la somme des carrés des écarts et par conséquent augmenter  $R$ . En ajoutant en effet une dimension à  $V_{\mathbf{X}}$ , on ne peut que diminuer la distance de  $\mathbf{y}$  à ce sous-espace. Dans ces conditions, la valeur prise par  $R$  ne peut être un critère absolu pour apprécier la qualité de l'ajustement.

### 3.2.4 Qualité de l'ajustement

Jusqu'à présent, on s'est borné à résoudre un problème purement numérique d'ajustement, avec une mesure globale de qualité fournie par le coefficient de corrélation multiple. Il s'agit maintenant de tester la qualité de cet ajustement et la signification statistique des coefficients de régression, ce qui nécessite de faire des hypothèses sur  $y$  et  $\varepsilon$ .

#### a – Spécification du modèle

On suppose que le résidu  $\varepsilon_i$  est l'effet résultant d'un grand nombre de causes non identifiées, et à ce titre, on le considérera comme une perturbation *aléatoire*. Ce point de vue étendu aux  $n$  relations du modèle introduit un vecteur aléatoire de résidus  $\varepsilon$  (ayant  $n$  composantes) et, par cet intermédiaire, définit  $y = X\alpha + \varepsilon$  comme vecteur aléatoire.

Le tableau 3.2 - 1 résume les caractéristiques des différents éléments du modèle :

Tableau 3.2 - 1  
Caractéristiques des éléments du modèle

$y = X\alpha + \varepsilon$	Observé	Non observable
Aléatoire	$y$ ( $n,1$ )	$\varepsilon$ ( $n,1$ )
Non aléatoire	$X$ ( $n,p$ )	$\alpha$ ( $p,1$ )

On supposera que les résidus  $\varepsilon_i$  ont une espérance nulle, qu'ils ont tous même variance  $\sigma^2$  et sont deux à deux non corrélés :

$$E(\varepsilon) = \begin{matrix} 0 \\ (1,n) \end{matrix} \quad \text{et} \quad \text{Var}(\varepsilon) = E(\varepsilon\varepsilon') = \begin{matrix} \sigma^2 I \\ (n,n) \end{matrix}$$

ce qui implique les relations :

$$E(y) = X\alpha \quad \text{et} \quad \text{Var}(y) = \text{Var}(\varepsilon) = \begin{matrix} \sigma^2 I \\ (n,n) \end{matrix} \quad [3.2 - 4]$$

Sous ces hypothèses, les coefficients de régression  $a_k$ , ( $k=1, \dots, p$ ), fournis par la technique des moindres-carrés sont les meilleurs estimateurs<sup>1</sup> des coefficients inconnus  $\alpha_k$ .

<sup>1</sup> Il s'agit plus précisément d'estimateurs à *variance minimale* sur l'ensemble des estimateurs linéaires, cette propriété étant connue sous le nom de théorème de Gauss-Markov. On renvoie aux ouvrages cités au début de ce chapitre pour plus de détails sur ce théorème et ses généralisations.

### b – Moyenne et variance des coefficients

Le vecteur  $\mathbf{a} = (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}'\mathbf{y}$  des coefficients de régression étant une fonction de  $\mathbf{y}$ , est lui même un vecteur aléatoire. La formule [3.2 - 4] nous montre immédiatement que son espérance mathématique s'écrit :  $E(\mathbf{a}) = \alpha$ .

Un calcul élémentaire<sup>1</sup> montre que la matrice des covariances des coefficients s'écrit :

$$\mathbf{V}(\mathbf{a}) = \sigma^2 (\mathbf{X}'\mathbf{X})^{-1}$$

Notons que  $\sigma^2$  est la variance théorique des résidus et n'est donc pas connue. On peut estimer  $\sigma^2$  par  $s^2$ , la variance empirique des écarts calculés après l'ajustement.

Si l'on désigne par  $\mathbf{V}$  la matrice des covariances empiriques des variables explicatives supposées centrées ( $\mathbf{V} = \frac{1}{n} \mathbf{X}'\mathbf{X}$ ), on a la relation :

$$\mathbf{V}(\mathbf{a}) = \frac{\sigma^2}{n} \mathbf{V}^{-1}$$

On remarque la dualité qui existe entre les variables explicatives et les coefficients de ces variables dans le modèle. Des variables explicatives non corrélées (matrice  $\mathbf{V}$  diagonale) conduiront à des coefficients de régression non-corrélés. Ce lien entre structure des prédicteurs et structure des coefficients sera précisé dans le paragraphe 3.2.5 consacré à la régression sur composantes principales.

### c – Tests sous l'hypothèse de normalité des résidus

Les résultats précédents (coefficient de corrélation multiple, matrices des covariances des coefficients) permettent d'imaginer des procédures de validation sous des hypothèses assez générales. Le fait de spécifier la loi des résidus autorise des épreuves de validation classiques que l'on rappelle ici, sans démonstration.

#### 1- Test sur les coefficients de régression

Pour savoir si une variable explicative  $x_k$  a une influence réelle sur la variable à expliquer  $y$ , on procède à un test d'hypothèse sur le coefficient de régression  $\alpha_k$ .

<sup>1</sup> La variance de  $\mathbf{a}$  s'écrit  $\mathbf{V}(\mathbf{a}) = E[(\mathbf{a} - \alpha)(\mathbf{a} - \alpha)']$ .

Or,  $\mathbf{a} - \alpha = (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}'\mathbf{y} - \alpha$

d'où :  $\mathbf{a} - \alpha = (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}'(\mathbf{X}\alpha + \varepsilon) - \alpha$

soit :  $\mathbf{a} - \alpha = (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}'\varepsilon$

On obtient donc :  $E[(\mathbf{a} - \alpha)(\mathbf{a} - \alpha)'] = (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}' E(\varepsilon\varepsilon') \mathbf{X} (\mathbf{X}'\mathbf{X})^{-1}$

Finalement :  $\mathbf{V}(\mathbf{a}) = \sigma^2 (\mathbf{X}'\mathbf{X})^{-1}$

L'hypothèse nulle ( $H_0$ ) est l'éventuelle non-influence qui se traduit par :

$$(H_0) \quad \alpha_k = 0 \quad (\text{les autres coefficients sont quelconques})$$

On écrit alors la statistique de Student :

$$t = \frac{a_k}{s_k}$$

où  $s_k$  est l'estimation de l'écart-type du  $k^{\text{ième}}$  coefficient de régression  $a_k$  :

$$s_k = \sqrt{\frac{\|y - Xa\|^2}{n-p}} a_{kk}, \text{ où } a_{kk} \text{ désigne le } k^{\text{ème}} \text{ élément diagonal de } (X'X)^{-1}.$$

Si ( $H_0$ ) est vraie, la statistique suit une loi de Student à  $(n-p)$  degrés de liberté<sup>1</sup>. Soit  $p_c$  la probabilité tirée de la distribution de Student correspondant à la valeur  $t_c$  prise par  $t$  :

$$p_c = P(|t| \geq t_c)$$

Si cette probabilité est jugée "trop faible", on rejette<sup>2</sup> l'hypothèse ( $H_0$ ). On peut étendre la procédure de ce test à une combinaison linéaire quelconque des coefficients.

## 2- Test sur un sous-ensemble de coefficients

On vient de voir comment tester l'un après l'autre la nullité de chaque coefficient. Cependant, les réponses à des questions telles " $\alpha_1 = 0$  sans rien supposer sur  $\alpha_2$ "? puis " $\alpha_2 = 0$  sans rien supposer sur  $\alpha_1$ ?" ne déterminent pas la réponse à cette autre question : " $\alpha_1 = 0$  et simultanément  $\alpha_2 = 0$ ?" D'où l'utilité de savoir tester la *nullité simultanée* de plusieurs coefficients de régression.

On se place ici, sans perte de généralité, dans le cas où les  $q$  coefficients sont les premiers des  $p$  coefficients. L'hypothèse  $H_0$  se traduit par :

$$-(H_0) \quad \alpha_1 = \alpha_2 = \dots = \alpha_q = 0 \quad (\text{les autres } \alpha_k \text{ quelconques})$$

$$-(H_1) \quad \text{un au moins des } q \text{ premiers } \alpha_k \text{ n'est pas nul}$$

Convenons de noter  $X_{H_0}$  les  $p-q$  dernières colonnes de  $X$  et  $\alpha_{H_0}$  les  $p-q$  dernières composantes de  $\alpha$ . L'écriture matricielle des modèles sera :

$$\begin{cases} \text{modèle (complet) sous } H_1 : & y = X\alpha + \varepsilon \\ \text{modèle (réduit) sous } H_0 : & y_0 = X_{H_0}\alpha_{H_0} + \varepsilon \end{cases}$$

<sup>1</sup> Le modèle contient  $p+1$  coefficients à estimer : le terme constant et les coefficients des  $p$  variables explicatives.

<sup>2</sup> On effectue par exemple le test au seuil de confiance  $0,05$  : si  $p_c < 0,05$  on rejette l'hypothèse selon laquelle la variable  $x_k$  n'a pas d'influence réelle (avec moins de 5 chances sur 100 de se tromper) ; alors que si  $p_c \geq 0,05$ , on ne peut pas rejeter cette hypothèse.

On considère la statistique  $F$  qui suit une loi de Fisher<sup>1</sup> à  $q$  et  $n - p$  degrés de liberté :

$$F = \frac{\left( \|y - \bar{y}_0\|^2 - \|y - \bar{y}\|^2 \right) / q}{\|y - \bar{y}\|^2 / (n - p)} \quad [3.2-5]$$

On note les sommes des carrés des écarts :

$$S_0 = \|y - \bar{y}_0\|^2 \quad \text{et} \quad S_1 = \|y - \bar{y}\|^2$$

Si la différence entre les deux quantités  $S_0$  et  $S_1$  est grande ( $F$  grand) alors l'effet des  $q$  premières variables est important et on devra rejeter l'hypothèse nulle; les  $q$  variables  $x_1, \dots, x_q$  ont simultanément une influence sur  $y$ . On effectue donc deux ajustements successifs<sup>2</sup> pour calculer d'une part  $S_1$  sur le modèle complet et d'autre part  $S_0$  sur le modèle pour lequel sont exclues les  $q$  variables explicatives en cause.

### 3.2.5 Régression régularisée et régression sur composantes principales

On a vu que la structure du tableau à  $n$  lignes et  $p$  colonnes  $X$  des variables explicatives (structure décrite par la matrice des covariances) avait des répercussions sur la qualité des coefficients de régression (§ 3.2.4.b). Le calcul des coefficients de régression requiert une matrice  $X'X$  inversible et donc des vecteurs  $x_1, x_2, \dots, x_p$  linéairement indépendants.

Si les variables explicatives sont fortement corrélées (autrement dit si certains des vecteurs  $x_1, x_2, \dots, x_p$  ont des directions voisines) alors l'inversion de la matrice  $X'X$  est difficile. Le vecteur  $a$  dont les composantes sont les coordonnées de la projection de  $y$  dans la base de  $V_X$  formée par  $x_1, x_2, \dots, x_p$  est mal spécifié. Les résultats de la régression seront instables<sup>3</sup>.

<sup>1</sup> Le principe de tous ces tests est très simple : les statistiques  $F$  sont des quotients de  $\chi^2$  indépendants. Les  $\chi^2$  sont indépendants car ils correspondent à des composantes normales orthogonales du vecteur résiduel (ici : côté de l'angle droit du triangle rectangle  $(y, \bar{y}, \bar{y}_0)$  dont l'hypothénuse est  $(y, \bar{y}_0)$ ).

<sup>2</sup> D'un point de vue numérique on peut passer d'une somme de carrés à l'autre sans être obligé de refaire un ajustement complet.

<sup>3</sup> La décomposition en éléments propres de  $X'X$  s'écrit :  $X'X = U'\Lambda U = \sum_{\alpha=1}^p \lambda_\alpha u_\alpha u'_\alpha$ , où  $\Lambda$  est la matrice diagonale dont le  $\alpha^{\text{ième}}$  élément est la valeur propre  $\lambda_\alpha$  et  $U$  le tableau des vecteurs propres unitaires correspondants. On a donc également :

$$(X'X)^{-1} = U\Lambda^{-1}U' = \sum_{\alpha=1}^p \frac{1}{\lambda_\alpha} u_\alpha u'_\alpha.$$

L'estimation de la matrice de covariances du vecteur  $a$  des coefficients vaut :

$$\text{Var}(a) = s^2 (X'X)^{-1} = s^2 \sum_{\alpha=1}^p \frac{1}{\lambda_\alpha} u_\alpha u'_\alpha$$

Sous cette forme on voit comment une ou plusieurs valeurs propres presque nulles rendent imprécis l'ajustement.



On a également évoqué le fait que la méthode des moindres carrés pouvait donner un poids excessif à des points éloignés (pouvant parfois être erronés ou aberrants).

On a vu d'autre part à la section 1.2 que l'analyse en composantes principales décrit la structure d'un tableau  $X$  en mettant en évidence les interrelations entre variables (colonnes de  $X$ ) ; elle permet également de visualiser les points-observations (points-lignes de  $X$ ) et donc d'aider à repérer d'éventuelles anomalies dans leur distribution. Enfin, on a vu que l'analyse fournit une base orthogonale hiérarchisée du sous-espace de  $\mathbb{R}^n$  appelé  $V_X$ .

Il est clair dans ces conditions qu'une analyse en composantes principales préalable permettra d'apprécier l'existence de colinéarités entre les variables explicatives, de détecter les redondances et compétitions entre prédicteurs ; de repérer les individus occupant des positions aberrantes ou simplement suspects. Il s'agit là d'une phase descriptive qui doit précéder la régression.

L'analyse peut également fournir des variables artificielles orthogonales (les coordonnées des points-observations sur les nouveaux axes) comme nouveaux prédicteurs : c'est la régression sur composantes principales, recommandée lorsque les variables explicatives sont nombreuses ou fortement corrélées entre elles. L'analyse factorielle joue donc un double rôle : un rôle d'exploration préalable et un rôle de régularisation<sup>1</sup>.

### a – Principe de la régression régularisée

Le principe revient à remplacer les  $p$  variables explicatives  $x_1, x_2, \dots, x_p$  par leurs  $p$  composantes principales qui engendrent le même sous-espace  $V_X$  à  $p$  dimensions. S'il existe  $r$  relations linéaires entre les variables explicatives, alors la transformation des  $p$  variables fournira  $q = p - r$  composantes principales. Il est possible ensuite d'exprimer les résultats de la régression en fonction des variables initiales. Nous nous plaçons dans  $\mathbb{R}^n$  où un point  $y$  est projeté sur le sous-espace  $V_X$  engendré par les vecteurs  $x_1, x_2, \dots, x_p$ .

Les  $p$  vecteurs propres  $u_k$  auxquels correspondent  $p$  composantes principales constituent une base orthonormée du sous-espace  $V_X$  sur lequel on veut projeter  $y$ .

On élimine le problème posé par la quasi-colinéarité si on supprime de cette base les  $p - r$  vecteurs  $u_k$  correspondant à des valeurs propres  $\lambda_k$  nulles ou très faibles.

---

<sup>1</sup> Les techniques de régularisation, largement utilisées en analyse discriminante, participent à la résolution de problèmes *mal posés* (ici : cas de colinéarité entraînant une singularité de la matrice  $X'X$ , et donc une impossibilité de calcul de  $a$ ) ou de problèmes *pauvrement posés* (ici : cas de quasi-colinéarité, entraînant une instabilité numérique de  $(X'X)^{-1}$  et du vecteur  $a$  des coefficients de régression). Pour une revue des traitements de la colinéarité dans le cas de la régression, cf. Palm et Lemma (1995).

Autrement dit on ne retient que les  $q$  premières composantes principales de variances non négligeables.

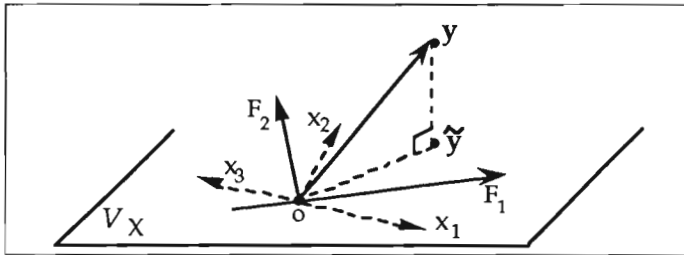


Figure 3.2 - 4  
Régression sur composantes principales

Les variables étant centrées, nous sommes dans le cas de l'analyse générale de la section 1.1. Le tableau  $X$  est reconstitué sur les  $q$  premiers axes factoriels (cf. [1.1 - 7] du § 1.1.5.b) par la formule ( $v_\alpha$  et  $u_\alpha$  sont unitaire) :

$$X^* = \sum_{\alpha=1}^q \sqrt{\lambda_\alpha} v_\alpha u'_\alpha = V_q \Lambda_q^{1/2} U'_q \quad (\text{avec } q < p)$$

où  $V_q$  et  $U_q$  sont les matrices respectivement d'ordre  $(n, q)$  et  $(p, q)$  contenant en colonne les vecteurs propres  $v_\alpha$  et  $u_\alpha$  et  $\Lambda_q$  la matrice diagonale  $(q, q)$  des valeurs propres.

On calcule<sup>1</sup> à partir de ce nouveau tableau le vecteur de coefficient  $a^*$  :

$$a^* = \sum_{\alpha=1}^q \frac{1}{\sqrt{\lambda_\alpha}} u_\alpha v'_\alpha y \quad [3.2 - 6]$$

Remarquons que  $a^*$  n'est plus unique, puisque tout vecteur de la forme  $a^* + c$  (avec  $c$  tel que  $U'c = 0$ ) satisfait aux équations [3.2 - 1].

Pour que la relation  $E(a^*) = \alpha$  soit vérifiée, il faut, dans le cas de l'estimation précédente, que le modèle théorique spécifie que  $\alpha$  soit de la forme  $U\beta$ ,  $\beta$  étant un vecteur quelconque à  $q$  composantes.

<sup>1</sup> Les équations [3.2 - 1] s'écrivent  $X'Xa = X'y$ , c'est-à-dire, en abandonnant provisoirement les indices  $q$  :

$$U\Lambda U'a = U\Lambda^{1/2}V'y$$

Le vecteur  $a$  n'ayant que  $q$  composantes indépendantes peut s'écrire sous la forme :  $a = Ub$

d'où puisque  $U'U = I$  (matrice unité  $(q, q)$ ) :

$$U\Lambda b = U\Lambda^{1/2}V'y$$

Prémultipliant les deux membres par  $U'$ , on obtient  $b$  :

$$b = \Lambda^{-1/2}V'y, \text{ donc } a = U\Lambda^{-1/2}V'y$$

Dans ces conditions, l'estimation de la matrice des covariances de  $\mathbf{a}^*$  (de rang  $q$ ) sera :

$$\text{Var}(\mathbf{a}^*) = s^2 \sum_{\alpha=1}^q \frac{1}{\lambda_{\alpha}} \mathbf{u}_{\alpha} \mathbf{u}'_{\alpha}$$

Notons que  $\mathbf{X} = \mathbf{X}^*$  s'il y a exactement  $q$  valeurs propres différentes de 0.

## b – Variables supplémentaires et régression

La procédure de mise en éléments supplémentaires dans une analyse en composantes principales constitue une variante descriptive de la régression multiple. D'un point de vue géométrique, les deux situations sont très similaires :

- les  $p$  variables explicatives engendrent un sous-espace  $V_X$  ayant au plus  $p$  dimensions sur lequel est projetée la variable à expliquer;
- les  $p$  variables actives de l'analyse engendrent aussi un sous-espace à au plus  $p$  dimensions que l'on réduit à  $q$  facteurs pour le visualiser et c'est sur ce sous-espace réduit à  $q$  dimensions que l'on projette les variables supplémentaires pour les situer par rapport aux variables actives.

La formule [3.2 - 6] précédente permet d'explicitier ce lien. Calculons à partir d'elle la nouvelle estimation  $\tilde{\mathbf{y}}^*$  de  $\mathbf{y}$  en utilisant la formule [1.1 - 4] du § 1.1.4 :

$$\tilde{\mathbf{y}}^* = \mathbf{X}^* \mathbf{a}^* = \sum_{\alpha=1}^q \mathbf{v}_{\alpha} \mathbf{v}'_{\alpha} \mathbf{y}$$

On a ainsi obtenu une expression de l'opérateur-projection  $\mathbf{P}_{X^*}$  sur l'espace des  $q$  premiers axes factoriels.

Le dernier membre rappelle clairement que la coordonnée  $(\mathbf{v}'_{\alpha} \mathbf{y})$  de  $\tilde{\mathbf{y}}^*$  sur l'axe unitaire  $\mathbf{v}_{\alpha}$  correspond au positionnement classique de  $\mathbf{y}$  en variable supplémentaire dans l'analyse dont les variables actives sont les colonnes de  $\mathbf{X}$ .

## c – Expression des coefficients dans la nouvelle base

Désignons par  $\mathbf{z}_{\alpha}$  le vecteur des nouvelles coordonnées des points sur l'axe  $\mathbf{u}_{\alpha}$ . Rappelons que l'on a les relations :

$$\mathbf{z}_{\alpha} = \mathbf{X}^* \mathbf{u}_{\alpha} = \mathbf{X} \mathbf{u}_{\alpha} = \sqrt{\lambda_{\alpha}} \mathbf{v}_{\alpha} \quad (\alpha = 1, 2, \dots, q)$$

L'ajustement sur la nouvelle base  $(\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_q)$  s'écrira :

$$\mathbf{y} = \mathbf{Zc} + \mathbf{e}$$

où  $Z$  est le tableau  $(n,q)$  des vecteurs orthogonaux  $z_\alpha$  et  $c$  le vecteur des  $q$  nouveaux coefficients de régression cherchés.

Puisque  $Z'Z = \Lambda$ , matrice diagonale dont les éléments diagonaux sont les valeurs propres, on a :

$$c = (Z'Z)^{-1}Z'y = \Lambda^{-1}Z'y$$

Cette situation idéale pour laquelle les variables explicatives sont orthogonales revient d'ailleurs à faire  $q$  régressions simples, car chacun des  $p$  coefficients peut être estimé séparément.

On a en effet :

$$c_\alpha = \frac{z'_\alpha y}{\lambda_\alpha} = \frac{\text{cov}(z_\alpha, y)}{\text{var}(z_\alpha)}$$

La matrice des covariances des coefficients  $c$  sera estimée par :

$$\text{Var}(c) = s^2(Z'Z)^{-1} = s^2\Lambda^{-1}$$

autrement dit ces coefficients sont non corrélés et ont pour variances les quantités :

$$\text{var}(c_\alpha) = \frac{s^2}{\lambda_\alpha}$$

### 3.2.6 Régression sur variables nominales : l'analyse de la variance

Lorsque les variables explicatives sont nominales, la régression multiple n'est autre que *l'analyse de la variance*, technique liée aux plans d'expériences et aux traitements statistiques des données expérimentales<sup>1</sup>. Il est courant d'opposer données d'observation et données expérimentales, en réservant les méthodes exploratoires pour les premières, et les méthodes inférentielles ou confirmatoires pour les secondes. La distinction n'est pas si nette en pratique : d'une part, nous l'avons vu, beaucoup de concepts et d'outils sont communs ; d'autre part, les champs d'application peuvent fréquemment se recouvrir, et une attitude méthodologique trop rigide pourrait être néfaste. D'où l'intérêt de connaître les principes et les possibilités des outils de l'analyse des données expérimentales.

---

<sup>1</sup> C'est R.A. Fisher qui est à l'origine de l'analyse de la variance et des plans d'expérience, dans une série d'articles datant des années vingt, repris dans l'ouvrage historique "The Design of Experiments" (Fisher, 1935). Citons également sur ce sujet les traités de Cochran et Cox (1957), de Cox (1958). Bailey (1981) et Steinberg et Hunter (1984) présentent des exposés synthétiques plus récents. En langue française, on pourra consulter les chapitres consacrés à ce thème dans les ouvrages de Dagnélie (1981) et Tomassone *et al.* (1993).

**a – Codage des variables nominales**

Supposons que l'on dispose sur une variable  $y$  de  $n$  observations classées selon  $p$  variables nominales  $x_1, \dots, x_l, \dots, x_p$  à respectivement  $m_1, \dots, m_l, \dots, m_p$  modalités.

Le tableau des variables explicatives  $X$  se présente maintenant sous la forme d'un tableau disjonctif complet  $[X_1, \dots, X_l, \dots, X_p]$ .

Cependant, pour chaque sous-tableau  $X_l$ , la somme des colonnes vaut 1. Il existe donc  $p$  relations linéaires entre les colonnes de  $X$ . Le tableau  $X$  n'est pas de plein rang et la matrice  $X'X$  n'est pas inversible.

Le problème peut être résolu par une régularisation de la régression (cf. § 3.2.5). Mais le fait que la nature des relations linéaires entre variables explicatives soit connue *a priori* (structure disjonctive complète du tableau) suggère d'autres possibilités de solutions.

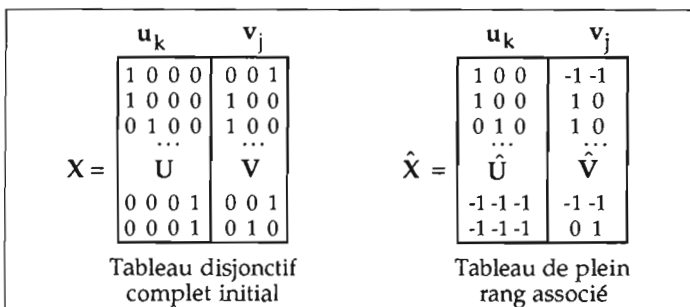
Pour éliminer la multicollinéarité, on peut ne retenir que  $m_l - 1$  modalités pour chaque variable  $x_l$  à  $m_l$  modalités. La modalité supprimée se recalcule évidemment à partir des autres. Une autre possibilité est également de supprimer une colonne de chaque sous-tableau mais après l'avoir retranchée aux colonnes restantes. Nous retiendrons ce deuxième codage mieux adapté au modèle linéaire avec interaction entre les variables explicatives.

Le tableau des variables explicatives ainsi recodé  $\hat{X}$  est de plein rang :

$$rang(\hat{X}) = \sum_{l=1}^p (m_l - 1)$$

Pour simplifier l'exposé, on se placera par la suite dans le cas où l'on dispose de deux variables nominales  $u$  et  $v$  ayant respectivement  $q$  et  $r$  modalités.

Notons  $u_k$  et  $v_j$ , les indicatrices des variables  $u$  et  $v$  avec  $1 < k < q$  et  $1 < j < r$ ,  $[U, V]$  le tableau disjonctif complet correspondant de dimension  $(n, q + r)$  et  $[\hat{U}, \hat{V}]$  le tableau disjonctif complet de plein rang et de dimension  $(n, q + r - 2)$  obtenu après recodage.



**Figure 3.2- 5**  
Tableaux des variables explicatives initial et recodé

La généralisation se fera sans difficulté.

**b – Modèle linéaire sans interaction**

On cherche à déterminer s'il existe un effet dû à la variable  $u$  et un effet dû à la variable  $v$ , autrement dit, si  $u$  et  $v$  ont une influence sur  $y$ .

Les variables sont ici considérées sans interaction et l'on dispose d'un modèle linéaire où les effets sont par conséquent additifs :

$$y_{ikj} = \mu + \alpha_i + \beta_j + \varepsilon_{ikj}$$

avec  $i = 1, \dots, n$  ;  $k = 1, \dots, q - 1$  et  $j = 1, \dots, r - 1$ . Ce modèle s'exprime sous forme matricielle par :

$$y = \mu \mathbf{1} + (\alpha_1 \mathbf{u}_1 \dots + \alpha_k \mathbf{u}_k \dots + \alpha_{q-1} \mathbf{u}_{q-1}) + (\beta_1 \mathbf{v}_1 \dots + \beta_j \mathbf{v}_j \dots + \beta_{r-1} \mathbf{v}_{r-1}) + \varepsilon$$

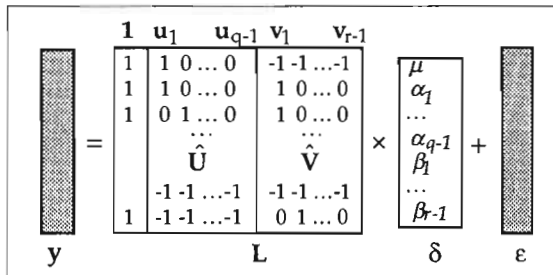
soit encore :

$$y = \mu \mathbf{1} + \hat{U}\alpha + \hat{V}\beta + \varepsilon$$

où  $\mathbf{1}$  est un vecteur de  $n$  composantes égales à 1 et  $\mu$  un coefficient scalaire.

Rassemblons dans un tableau  $L$  de dimension  $(n, q + r - 1)$  l'ensemble des variables explicatives artificielles et dans le vecteur  $\delta$  à  $(q + r - 1)$  composantes les coefficients  $\alpha_k, \beta_j$  et  $\mu$  du modèle. Il prend la forme matricielle :

$$y = L\delta + \varepsilon$$



**Figure 3.2 - 6**  
**Modèle de l'analyse de la variance :**  
**cas de deux variables  $u$  et  $v$  sans interaction**

Le problème est de tester si les  $\alpha_k$  (puis les  $\beta_j$ ) sont égaux entre eux, l'hypothèse alternative étant que l'un au moins des coefficients dans chaque groupe diffère des autres<sup>1</sup>.

On teste en d'autres termes les effets des variables  $u$  et  $v$ .

<sup>1</sup> La spécification du modèle est la même que lors de la régression multiple (résidus indépendants entre eux, de même variance). Pour procéder aux tests statistiques, il est nécessaire de supposer la normalité de la distribution des résidus.

On réalise alors le test de nullité simultanée des coefficients  $\alpha_k$ , ( $k = 1, \dots, q-1$ ) (cf. § 3.2.4.c).

Pour cela, on effectue successivement deux ajustements pour calculer d'une part  $S(\mu, \alpha, \beta)$  sur le modèle complet  $y = L\delta + \varepsilon$  et d'autre part  $S(\mu, \beta)$  sur le modèle réduit obtenu en supprimant dans  $L$  les  $q-1$  colonnes correspondant aux  $\alpha_k$ . La statistique du test sera d'après [3.2 - 5] :

$$F = \frac{(S(\mu, \beta) - S(\mu, \alpha, \beta))/(q-1)}{S(\mu, \alpha, \beta)/(n-q-r+1)}$$

On rejettera l'hypothèse nulle d'absence d'effet de la variable  $u$  si la probabilité de dépasser la valeur  $F$ , pour une variable de Fisher à  $(q-1)$  et  $(n-q-r+1)$  degrés de liberté, est jugée trop petite.

Pour tester l'existence d'un effet dû à la variable  $v$ , on procédera de façon analogue.

### c – Modèle linéaire avec interaction

Si l'on pense maintenant que l'effet de la modalité  $k$  de  $u$  peut être différent selon la modalité  $j$  de  $v$ , il faut ajouter au modèle l'effet d'interaction entre les deux variables  $u$  et  $v$ .

Cela peut se faire en juxtaposant au tableau disjonctif complet  $[\hat{U}, \hat{V}]$  le sous-tableau  $\hat{U} \times \hat{V}$  des interactions. On obtient  $\hat{U} \times \hat{V}$  en faisant le produit terme à terme des colonnes  $u_k$  par les colonnes  $v_j$ .

Puisque  $1 < k < q-1$  et  $1 < j < r-1$ , on engendre ainsi  $(q-1) \times (r-1)$  colonnes contenant les produits de deux indicatrices correspond à la conjonction des présences d'effet. On vérifie que le nouveau tableau ainsi construit  $[\hat{U}, \hat{V}, \hat{U} \times \hat{V}]$  est bien de plein rang  $q \times r$ . Le modèle s'exprime alors par :

$$y = \mu 1 + \hat{U}\alpha + \hat{V}\beta + (\hat{U} \times \hat{V})\gamma + \varepsilon$$

où  $\gamma$  est un vecteur à  $(q-1) \times (r-1)$  composantes.

#### - Test de l'effet de la variable $u$ et de l'effet de la variable $v$

Pour tester l'effet de la variable  $u$  on pose l'hypothèse nulle :

$$(H_0) \quad \alpha_k = 0 \quad (k = 1, \dots, q-1)$$

On effectue, comme pour le modèle sans interaction, le test de nullité simultanée des coefficients  $\alpha_k$ .

On calcule donc les sommes des carrés d'écart des ajustements sur le modèle complet et sur le modèle réduit, notées respectivement  $S(\mu, \alpha, \beta, \gamma)$  et  $S(\mu, \beta, \gamma)$ . On calcule ensuite la statistique  $F$  de Fisher à  $(q-1)$  et  $(n-pr)$  degrés de liberté d'après [3.2 - 5]. On agira de façon analogue pour tester l'effet de la variable  $v$ .

### - Test de l'interaction entre $u$ et $v$

Pour tester maintenant l'effet dû à l'interaction entre les deux variables  $u$  et  $v$ , on effectue le test de nullité simultanée des coefficients  $\gamma_j$  en calculant les quantités  $S(\mu, \alpha, \beta, \gamma)$  correspondant au modèle complet et  $S(\mu, \alpha, \beta)$  associé au modèle réduit où l'on a supprimé les  $(q-1) \times (r-1)$  colonnes correspondant aux  $\gamma_j$ . On calcule, toujours d'après [3.2 - 5], la statistique de Fisher à  $(q-1) \times (r-1)$  et  $(n - pr)$  degrés de liberté.

#### Remarques :

- 1) Il faut souligner que le choix du codage du tableau des variables explicatives pour avoir un tableau de plein rang est primordial ici alors qu'il était indifférent dans le modèle sans interaction.
- 2) La procédure développée dans le cas d'une interaction entre deux variables nominales peut être généralisée à des modèles comprenant plus de deux critères ( $u, v, w, \dots$ ), des interactions d'ordre 1 ( $uv, uw, vw, \dots$ ), des interactions d'ordre 2 ( $uvw, \dots$ ), etc. Cependant une certaine prudence s'impose pour plusieurs raisons. Tout d'abord, il est de plus en plus difficile d'apprécier et d'énoncer clairement la nature des hypothèses testées. D'autre part les interactions d'ordre élevé peuvent conduire à des tests "en chaîne" d'interprétation délicate ( $uv$  significatif,  $vw$  non significatif,  $uvw$  significatif, etc.). Enfin, on peut montrer qu'une interaction (surtout d'ordre élevé) peut n'être due qu'à la présence d'une observation légèrement aberrante (la procédure n'est pas robuste).

### 3.2.7 Régression sur variables mixtes : analyse de la covariance

Dans un modèle d'analyse de la variance, la valeur de la variable à expliquer est déterminée, à l'aléa  $\varepsilon$  près, par les classes dans lesquelles sont faites les mesures ou observations. On peut cependant imaginer un modèle où cette valeur est, à l'intérieur de chaque classe  $k$ , fonction également d'une ou plusieurs variables explicatives continues. On dira par exemple que la dépense individuelle en habillement est fonction du sexe  $u$  et pour chaque sexe fonction du revenu  $x$  de l'individu  $i$ .

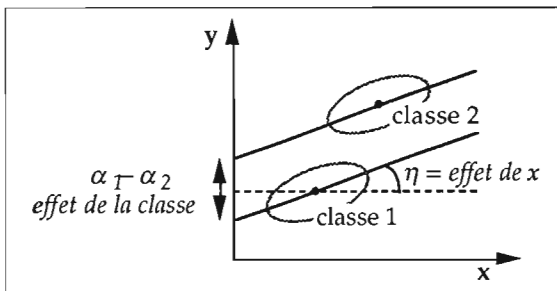


Figure 3.2 - 7  
Un modèle d'analyse de la covariance :  
variable nominale sans effet sur la pente de la régression



La figure 3.2 - 7 illustre un modèle où l'observation  $i$  dans la classe  $k$  serait déterminée par :

$$y_{ik} = \mu + \alpha_k + \eta x_{ik} + \varepsilon_{ik}$$

En donnant la même pente  $\eta$  aux deux droites passant par les centres de classe, on suppose ici que le revenu a le même effet quel que soit le sexe; la distance  $(\alpha_1 - \alpha_2)$  entre les deux droites mesure "l'effet du sexe". On aurait pu supposer un effet du revenu différencié suivant le sexe en traçant des droites non parallèles.

De tels modèles, où interviennent des variables nominales et des variables continues, sont appelés modèles d'analyse de la covariance. Ils vont se traduire par :

$$y = L^* \delta + \varepsilon$$

où  $L^*$  est le tableau de plein rang des variables explicatives.

**a – Modèles d'analyse de la covariance**

Plaçons-nous, pour simplifier l'exposé, dans le cas où le modèle contient une variable nominale  $u$  à  $q$  modalités et une variable continue  $x$ .

Le modèle le plus général correspondant au modèle complet suppose à la fois un effet dû à la variable nominale  $u$  et un effet  $x$  différencié pour chaque catégorie  $k$ ,  $1 < k < q - 1$ , ce qui s'exprime par :

$$y_{ik} = (\mu + \alpha_k) + (\eta + \beta_k) x_{ik} + \varepsilon_{ik} \tag{3.2 - 7}$$

Le tableau  $L$  est construit en deux parties : les  $q$  premières colonnes correspondent à l'analyse de la variance à un critère; les  $q - 1$  colonnes suivantes expriment de façon analogue l'effet différencié de  $x$  suivant la catégorie  $k$  de la variable  $u$ , mesuré autour de l'effet général représenté par la dernière colonne.

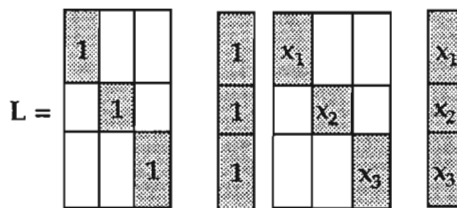


Figure 3.2 - 8  
Tableau des variables explicatives :  
cas d'une variable nominale  $u$  à 3 modalités et d'une variable continue  $x$

On remarquera que l'on obtient les  $q$  dernières colonnes comme une interaction entre la variable nominale  $u$  et la variable continue  $x$ , c'est-à-dire par multiplication terme à terme des  $q$  premières colonnes par  $x$ .

On notera  $S(\mu, \alpha, \eta, \beta)$  la somme de carrés d'écarts des ajustements sur le modèle complet [3.2 - 7].

**b – Test d'un effet différencié de x dans chaque classe k**

Pour tester l'existence d'un effet différencié de x dans chaque classe k, on effectuera un deuxième ajustement sur le modèle :

$$y_{ik} = (\mu + \alpha_k) + \eta x_{ik} + \varepsilon_{ik}$$

Ce modèle est la réduction du modèle complet [3.2 - 7], obtenu par introduction de l'hypothèse nulle :

$$(H_0) \quad \begin{cases} \beta_k = 0 & (k = 1, \dots, q-1) \\ \mu, \eta, \alpha_k & \text{quelconques} \end{cases}$$

La statistique du test s'obtient par application de la formule [3.2 - 5] :

$$F = \frac{(S(\mu, \alpha, \eta) - S(\mu, \alpha, \eta, \beta))/(q-1)}{S(\mu, \alpha, \eta, \beta)/(n-2q)}$$

On rejettera l'hypothèse nulle si la probabilité de dépasser la valeur lue dans la table de Fisher-Snedecor à  $(q-1)$  et  $(n-2q)$  degrés de liberté, est jugée trop petite.

**c – Test de l'effet de la variable u**

Pour tester l'existence de l'effet de la variable nominale u (tout en supposant cependant un effet différencié de x dans les classes), on calculera  $S(\mu, \eta, \beta)$  sur le modèle :

$$y_{ik} = \mu + (\eta + \beta_k) x_{ik} + \varepsilon_{ik}$$

pour le comparer à  $S(\mu, \alpha, \eta, \beta)$ . Ce modèle est la réduction du modèle complet [3.2 - 7] obtenu par introduction de l'hypothèse nulle :

$$(H_0) \quad \begin{cases} \alpha_k = 0 & (k = 1, \dots, q-1) \\ \mu, \eta, \beta_k & \text{quelconques} \end{cases}$$

La statistique du test fait référence à la formule [3.2 - 5] pour laquelle les degrés de liberté sont  $(q-1)$  et  $(n-2q)$ .

**d - Test d'un "effet classe global"**

On testera l'existence d'un "effet classe globale" à l'aide de  $S(\mu, \eta)$  calculé sur le modèle :

$$y_{ik} = \mu + \eta x_{ik} + \varepsilon_{ik}$$

pour le comparer à  $S(\mu, \alpha, \eta, \beta)$ . Ce modèle est la réduction du modèle complet [3.2 - 7], obtenu par introduction de l'hypothèse nulle :

$$(H_0) \quad \begin{cases} \alpha_k = 0 \text{ et } \beta_k = 0 & (k = 1, \dots, q-1) \\ \mu, \eta & \text{quelconques} \end{cases}$$

La statistique du test renvoie à la formule [3.2 - 5] ayant  $(2q-2)$  et  $(n-2q)$  degrés de liberté.

Elle permet de répondre à la question : est-ce que la valeur de  $y$  dépend de la classe, soit par des centres de classe distincts, soit par des pentes en  $x$  différentes ?

### e – Généralisation de l'analyse de la covariance

L'introduction de plusieurs variables continues ( $x_1, x_2, \dots$ ) ne présente aucune difficulté. Le déploiement de chacune d'elles se fait dans  $L$  comme le déploiement de la colonne  $x$  effectué précédemment. Les calculs de sommes de carrés d'écart et les constructions de tests s'effectuent selon les mêmes principes.

Il est plus délicat de généraliser la procédure au cas de plusieurs variables nominales. On rencontre en particulier les difficultés déjà évoquées en analyse de la variance lorsque l'on veut introduire un terme d'interaction entre les variables. Le problème est compliqué encore, dans la pratique, par la nécessité de choisir au départ le modèle *a priori* qui est censé représenter correctement le phénomène et qui servira de référence dans la construction des tests.

## 3.2.8 Choix des variables, généralisations du modèle

L'exposé qui précède ne fait que situer les principes de base du modèle linéaire par rapport aux méthodes descriptives de la première partie. Les méthodes présentées correspondent à une part notable des applications les plus courantes, mais à une part infime de la littérature théorique et technique sur le sujet, pour laquelle nous renvoyons le lecteur à la bibliographie citée au début du chapitre.

On évoquera brièvement deux points dans ce paragraphe de conclusion : le problème de la sélection des variables dans les modèles et celui de la généralisation du modèle.

### a – Sélection et choix des variables explicatives

La qualité de l'ajustement dépend également du choix des prédicteurs et il est souhaitable de retenir un nombre limité de variables, non redondantes et ayant un pouvoir prédictif.

Une technique souvent utilisée pour sélectionner les variables explicatives est la méthode pas-à-pas ou *stepwise*<sup>1</sup>. Elle consiste à effectuer une première régression simple sur une variable puis à ajouter successivement celles qui

---

<sup>1</sup> La méthode de Furnival et Wilson (Furnival, 1971 ; Furnival and Wilson, 1974) permet de calculer les meilleures régressions pour  $1, 2, \dots, p$  variables explicatives, par une exploration optimisée de toutes les possibilités. En pratique,  $p$  ne doit pas dépasser 40 pour que le volume de calcul reste raisonnable. Une telle procédure est recommandable car elle ne fait pas intervenir de critères externes (peu ou mal justifiés) pour inclure ou exclure des variables dans le modèle.

font augmenter le coefficient de corrélation multiple  $R^2$ , avec éventuellement remise en question des choix antérieurs. A chaque étape sont réalisés des tests sur les coefficients de régression ou sur des sous-ensembles afin de rejeter la variable ou d'éliminer éventuellement certaines variables introduites dans les étapes précédentes. Les critères d'Akaike (1973), de Mallows (1973), sont fréquemment utilisés pour sélectionner les modèles lors de ces procédures. Une revue des critères usuels se trouve dans Atkinson (1981). L'exploration des résidus est également très utilisée pour choisir ou compléter les variables du modèle, en général par des procédés graphiques (cf. Cook et Weisberg, 1982, 1994).

Les modèles graphiques (cf. par exemple : Whittaker, 1990 ; Wermuth et Cox, 1992 ; Fine, 1992) permettent, lorsque le nombre de variables explicatives n'est pas trop élevé, d'étudier les liaisons conditionnelles entre variables. Variables et liaisons sont représentées respectivement par les sommets et les arêtes de graphes de liaisons conditionnelles qui ont le mérite de conduire l'utilisateur à réfléchir sur la pertinence et les implications des modèles possibles.

Enfin on a vu qu'une analyse en composantes principales de tout ou partie des variables explicatives  $x_k$ , avec positionnement de la variable à expliquer  $y$  en élément supplémentaire, permet de positionner la ou les estimations  $\tilde{y}$  de  $y$  parmi les  $x_k$ . Il est également possible de positionner différents changements de variables, voire de nouvelles variables fonctions de plusieurs prédicteurs, et donc de porter une appréciation critique sur les redondances et complémentarités au sein du modèle et de ses extensions.

## b- Modèles linéaires généralisés

Ces modèles, présentés pour la première fois sous ce nom par Nelder et Wedderburn (1972), exposés de façon complète par McCullagh et Nelder (1989), généralisent le modèle linéaire de base sur deux points :

- 1- La combinaison linéaire notée  $\omega_i = a_0 x_{i0} + a_1 x_{i1} + \dots + a_p x_{ip}$  des variables explicatives n'est pas nécessairement l'espérance mathématique  $E(y_i)$  de la variable  $y_i$  mais peut être plus généralement une fonction  $g(\cdot)$  de  $E(y_i)$  (appelée *fonction lien*) et notée :

$$\omega_i = g[E(y_i)]$$

Pour le modèle linéaire classique :

$$\omega_i = E(y_i)$$

- 2- La loi des composantes de  $y$  appartient à la famille des lois exponentielles<sup>1</sup> (dont la loi normale est un cas particulier). Elle fait intervenir deux paramètres  $\theta$  et  $\varphi$ , et trois fonctions  $a(\cdot)$ ,  $b(\cdot)$ , et  $c(\cdot)$ .

<sup>1</sup> Cf. un exposé général dans : Dempster (1971) ; Berk (1972).

$$f_Y(y, \theta, \varphi) = e^{\left\{ \frac{y\theta - b(\theta)}{a(\varphi)} + c(y, \varphi) \right\}}$$

On voit que l'on obtient la fonction de densité de la loi normale :

$$f_Y(y; \theta, \varphi) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{\left\{ -\frac{(y-\mu)^2}{2\sigma^2} \right\}}$$

pour les spécifications suivantes des paramètres et des fonctions :

$$\theta = \mu; \quad \varphi = \sigma^2; \quad a(\varphi) = \varphi; \quad b(\theta) = \theta^2/2; \quad c(y, \varphi) = -1/2 \left\{ (y^2/\sigma^2) + \log(2\pi\sigma^2) \right\}$$

D'autres valeurs des paramètres et des fonctions conduisent aux lois binomiales, de Poisson, gamma.

L'ajustement du modèle se fait par la méthode du maximum de vraisemblance<sup>1</sup>, qui coïncide avec les moindres carrés dans le cas de la loi normale.

En faisant varier la loi de  $y$  et la *fonction lien*, le modèle linéaire généralisé inclut comme cas particulier une famille de modèles mettant en jeu des variables nominales, parmi lesquels les modèles log-linéaires (cf. section 3.4).

### 3.2.9 Modèles de variables latentes

Les modèles de variables latentes n'entrent pas dans le cadre du modèle linéaire général, mais ils sont apparentés à des modèles qui interviennent dans un cadre plus général, qui sont les *modèles à erreurs sur les variables* (exogènes)<sup>2</sup>.

Ces modèles ont été essentiellement développés en économétrie, où l'on distingue habituellement les modèles fonctionnels, ou à effet fixes (comme la régression multiple et le modèle linéaire dans son ensemble), et les modèles structurels ou à effet aléatoires (modèles de variables latentes). L'analyse factorielle en facteurs communs et spécifiques (*factor analysis*) est probablement le modèle le plus ancien<sup>3</sup>. Il est utilisé principalement par les psychologues et psychométriciens. Les développements auxquels il donne lieu sont complexes et diversifiés. On pourra consulter sur ce point les ouvrages de Harman (1967), Mulaik (1972).

<sup>1</sup> La méthode numérique de résolution est une méthode des moindres carrés pondérés itératifs très voisine de la méthode de Newton-Raphson.

<sup>2</sup> On trouvera un exposé des modèles et une note historique dans Malinvaud (1964).

<sup>3</sup> A l'origine des principes de la méthode se trouvent Spearman (1904) (analyse monofactorielle), puis Garnett (1919) et Thurstone (1947) (analyse multifactorielle).

### a – Le modèle

Cette méthode se propose de reconstituer, à partir d'un petit nombre  $q$  de *facteurs*, les corrélations existant entre  $p$  variables observées. On suppose l'existence d'un modèle *a priori* :

$$\underset{(p,1)}{x_i} = \underset{(p,q)}{\Gamma} \underset{(q,1)}{f_i} + \underset{(p,1)}{e_i} \quad [3.2-8]$$

Dans cette écriture  $x_i$  représente le  $i^{\text{ème}}$  vecteur observé des  $p$  variables;  $\Gamma$  est un tableau  $(p, q)$  de coefficients inconnus (avec  $q < p$ );  $f_i$  est la  $i^{\text{ème}}$  valeur du vecteur aléatoire et non observable de  $q$  *facteurs communs*; et  $e_i$  la  $i^{\text{ème}}$  valeur du vecteur non observable de résidus, lesquels représentent l'effet combiné de *facteurs spécifiques* et d'une perturbation aléatoire.

Ainsi par exemple, dans le cas des facteurs communs "f<sub>1</sub> = intelligence" et "f<sub>2</sub> = mémoire" que cherchaient les psychologues, le système [3.2 - 8] s'écrit pour le  $i^{\text{ème}}$  individu :

$$\begin{cases} x_{i1} = \gamma_{11}f_{i1} + \gamma_{12}f_{i2} + e_{i1} \\ x_{i2} = \gamma_{21}f_{i1} + \gamma_{22}f_{i2} + e_{i2} \\ \dots \\ x_{ip} = \gamma_{p1}f_{i1} + \gamma_{p2}f_{i2} + e_{ip} \end{cases}$$

Chaque observation de chaque variable est considérée comme une réalisation d'une variable aléatoire déterminée, par addition au résidu aléatoire spécifique, des deux variables aléatoires que sont les facteurs communs (avec des pondérations qui dépendent de chaque variable) <sup>1</sup>.

Désignons par  $X$  le tableau  $(n,p)$  dont la  $i^{\text{ème}}$  ligne est le vecteur transposé  $x_i'$  qui représente l'observation  $i$ . De même  $F$  désigne le tableau  $(n,q)$  non observable dont la  $i^{\text{ème}}$  ligne est  $f_i'$ ; et  $E$  le tableau  $(n,p)$  non observable dont la  $i^{\text{ème}}$  ligne est  $e_i'$ . Le modèle liant l'ensemble des observations aux facteurs hypothétiques s'écrit :

$$\underset{(n,p)}{X} = \underset{(n,q)}{F} \underset{(q,p)}{\Gamma'} + \underset{(n,p)}{E} \quad [3.2-9]$$

Dans cette écriture, seul  $X$  est observable, et le modèle est par conséquent indéterminé. Son identification et l'estimation des paramètres posent des problèmes complexes, sources d'une abondante littérature <sup>2</sup>. Une cascade d'hypothèses *a priori* supplémentaires va permettre d'écrire le problème sous une forme simplifiée, la seule que nous aborderons ici.

<sup>1</sup> Ainsi, on reconstitue approximativement les  $p$  notes d'un individu  $i$  dans  $p$  matières scolaires à partir de ses 2 notes factorielles, et de coefficients qui ne dépendent que des matières.

<sup>2</sup> Voir par exemple la synthèse et les références très complètes de Fine (1993). Il existe de nombreuses variantes de la méthode : axes obliques, rotations selon différents critères (varimax, quartimax, oblimax), recherches de structures simples, pour lesquelles on peut citer globalement l'ensemble des parutions de la revue *Psychometrika*.

Sans perte de généralité, nous supposons centrées les variables dont les observations sont les colonnes de  $\mathbf{X}$ , ainsi que les variables aléatoires que constituent les facteurs communs et les facteurs spécifiques. Nous utiliserons les notations suivantes:

- $\mathbf{W}$  matrice  $(p,p)$  des covariances théoriques entre variables;
- $\Phi$  matrice  $(q,q)$  des covariances théoriques entre facteurs communs;
- $\Delta$  matrice  $(p,p)$  des covariances théoriques entre facteurs spécifiques.

Appelons  $\mathbf{S}$  la matrice des covariances empiriques des observations  $\mathbf{X}$ , que nous supposons également centrées. Par définition et en vertu de [3.2 - 9], on a :

$$\mathbf{S} = \frac{1}{n} \mathbf{X}'\mathbf{X} = \frac{1}{n} (\mathbf{F}\Gamma' + \mathbf{E})'(\mathbf{F}\Gamma' + \mathbf{E})$$

c'est-à-dire :

$$\mathbf{S} = \frac{1}{n} \Gamma\mathbf{F}'\mathbf{F}\Gamma' + \frac{1}{n} \Gamma\mathbf{F}'\mathbf{E} + \frac{1}{n} \mathbf{E}'\mathbf{F}\Gamma' + \frac{1}{n} \mathbf{E}'\mathbf{E} \quad [3.2 - 10]$$

Aux hypothèses du modèle, nous ajouterons l'hypothèse *a priori* que les facteurs résiduels sont non corrélés aux facteurs communs; la matrice des covariances théoriques correspondantes étant nulle, nous considérerons comme négligeables dans [3.2 - 10] les matrices  $\frac{1}{n} \Gamma\mathbf{F}'\mathbf{E}$  et  $\frac{1}{n} \mathbf{E}'\mathbf{F}\Gamma'$  dont les espérances doivent être nulles. Ainsi la relation [3.2 - 10] prend la forme simplifiée :

$$\mathbf{S} = \frac{1}{n} \Gamma\mathbf{F}'\mathbf{F}\Gamma' + \frac{1}{n} \mathbf{E}'\mathbf{E}$$

correspondant à la relation théorique suivante pour le modèle :

$$\mathbf{W} = \Gamma\Phi\Gamma' + \Delta \quad [3.2 - 11]$$

Le problème d'estimation consiste à ajuster sur [3.2 - 11] une matrice  $\tilde{\mathbf{W}}$  qui, au regard d'un critère choisi par ailleurs, soit proche de la matrice des covariances empiriques  $\mathbf{S}$ . Mais afin d'obtenir une solution unique pour les paramètres de  $\Gamma$ ,  $\Phi$  et  $\Delta$ , il est nécessaire d'introduire des contraintes supplémentaires dans le modèle.

On suppose en général que les facteurs spécifiques sont non corrélés, c'est-à-dire que la matrice  $\Delta$  est diagonale. On impose de plus généralement que les facteurs communs soient orthogonaux et de variance unité, autrement dit la matrice  $\Phi$  est la matrice identité  $\mathbf{I}$  d'ordre  $q$ . La relation [3.2 - 11] du modèle s'écrit alors :

$$\mathbf{W} = \Gamma\Gamma' + \Delta$$

Sur cette relation le lien avec l'analyse en composantes principales apparaît clairement. Il s'agit dans ce cas de décomposer la matrice des covariances empiriques  $\mathbf{S}$  sous la forme:

$$\mathbf{S} = \mathbf{U}\Lambda\mathbf{U}'$$

où  $\Lambda$  est la matrice diagonale des valeurs propres (rangées) et  $\mathbf{U}$  le tableau des vecteurs propres unitaires correspondant. Cette relation s'écrit encore :

$$\mathbf{S} = (\mathbf{U}\Lambda^{1/2})(\mathbf{U}\Lambda^{1/2})' = \hat{\mathbf{U}}\hat{\mathbf{U}}'$$

où  $\hat{\mathbf{U}}$  est le tableau des vecteurs propres multipliés par les racines carrées des valeurs propres correspondantes.

Avec ce point de vue, l'analyse en facteurs communs et spécifiques suppose qu'en retranchant une matrice diagonale à éléments positifs ( $\tilde{\Delta}$  estimant  $\Delta$ ), on obtient une décomposition de la matrice des covariances empiriques sous la forme :

$$\mathbf{S} - \tilde{\Delta} = \Gamma\Gamma'$$

où  $\Gamma$  ne contient que  $q$  colonnes alors que dans  $\mathbf{S} = \hat{\mathbf{U}}\hat{\mathbf{U}}'$  le tableau  $\hat{\mathbf{U}}$  contenait  $p$  colonnes. On voit au passage qu'une analyse en composantes principales où les  $p - q$  dernières valeurs propres sont proches et voisines de 0, donnera des résultats très voisins de ceux d'une analyse à  $q$  facteurs communs orthogonaux.

### b- Estimation des paramètres inconnus

On n'insistera pas ici sur les problèmes posés par un tel modèle, qui font l'objet d'une abondante littérature. On donnera seulement quelques moyens pratiques de calcul.

Le problème essentiel est d'estimer  $\Delta$ , matrice diagonale des variances des résidus spécifiques. Une fois  $\Delta$  estimée par  $\tilde{\Delta}$ , il suffit de chercher les composantes principales (vecteurs propres) de  $(\mathbf{S} - \tilde{\Delta})$ ; on ne doit normalement trouver qu'un petit nombre de composantes différentes (statistiquement) de 0.

Nous allons examiner ici une spécification particulière du modèle, puis donner un algorithme de calcul dans le cas général.

#### - Cas de variances spécifiques égales

On suppose *a priori* que les facteurs spécifiques ont tous même variance théorique  $\sigma^2$ ; autrement dit par hypothèse  $\Delta = \sigma^2 \mathbf{I}$  :

$$\mathbf{W} = \Gamma\Gamma' + \sigma^2 \mathbf{I}$$

et, si on note  $s^2$  une estimation de  $\sigma^2$ , la relation [3.2 - 8] devient :

$$\mathbf{x}_i = \Gamma \mathbf{f}_i + \mathbf{se}_i$$

On obtiendrait une estimation de  $\Gamma$  en cherchant les composantes principales de la matrice  $(\mathbf{S} - s^2 \mathbf{I})$ . En effet, effectuant l'analyse de  $\mathbf{S}$ , on écrit :

$$\mathbf{S} = \mathbf{U}\Lambda\mathbf{U}'$$

et par conséquent :



$$\mathbf{S} - s^2\mathbf{I} = \mathbf{U}\mathbf{\Lambda}\mathbf{U}' - s^2\mathbf{U}\mathbf{U}' = \mathbf{U}(\mathbf{\Lambda} - s^2\mathbf{I})\mathbf{U}'$$

Les valeurs propres de  $(\mathbf{S} - s^2\mathbf{I})$  sont celles de  $\mathbf{S}$  diminuées de  $s^2$  (les vecteurs propres étant identiques). Puisque  $(\mathbf{S} - s^2\mathbf{I})$  doit être de rang  $q$ , il est nécessaire que  $s^2$  soit valeur propre multiple d'ordre  $p - q$  pour  $\mathbf{S}$ .

En particulier si, dans une analyse en composantes principales, les petites valeurs propres sont sensiblement égales, on peut considérer que les données sont engendrées par un modèle factoriel à variances spécifiques égales <sup>1</sup>.

### *-Une méthode de calcul dans le cas général*

La méthode que nous donnons ici est simple <sup>2</sup>. Elle procède de façon itérative, en posant au départ  $\tilde{\Delta} = 0$ . On calcule les vecteurs propres unitaires de  $\mathbf{S}$  rangés dans le tableau  $\mathbf{U}$  :

$$\mathbf{S} = \mathbf{U}\mathbf{\Lambda}\mathbf{U}' = \hat{\mathbf{U}}\hat{\mathbf{U}}'$$

Si l'on veut retenir  $q$  facteurs communs, on ne garde que les  $q$  premières colonnes de  $\hat{\mathbf{U}}$ , tableau que l'on notera  $\hat{\mathbf{U}}_1$ . On devrait pouvoir écrire :

$$\mathbf{S} = \hat{\mathbf{U}}_1\hat{\mathbf{U}}_1' + \tilde{\Delta}$$

On estimera donc provisoirement  $\tilde{\Delta}$  par les éléments diagonaux  $\tilde{\Delta}_1$  de  $(\mathbf{S} - \hat{\mathbf{U}}_1\hat{\mathbf{U}}_1')$ , et on calculera les  $q$  premiers vecteurs propres  $\hat{\mathbf{U}}_2$  de  $(\mathbf{S} - \tilde{\Delta}_1)$ .

A l'itération suivante on estime  $\tilde{\Delta}$  par les éléments diagonaux  $\tilde{\Delta}_2$  de  $(\mathbf{S} - \hat{\mathbf{U}}_2\hat{\mathbf{U}}_2')$  et l'on poursuit les opérations jusqu'à observer une convergence raisonnable du processus. On aura alors obtenu la décomposition cherchée :

$$\mathbf{S} = \Gamma\Gamma' + \tilde{\Delta}.$$

Mentionnons pour conclure ce bref aperçu les travaux d'Anderson et Rubin (1956) et de Lawley et Maxwell (1963) qui ont placé l'analyse factorielle en facteurs communs et spécifiques dans un cadre inférentiel classique.

<sup>1</sup> Ce modèle à variances spécifiques égales peut être justifié lorsque les  $p$  variables sont mesurées avec le même instrument (exemples : mensurations anthropométriques), et donc avec la même erreur.

<sup>2</sup> Cette procédure est parfois appelée analyse en facteurs principaux. Pour une première estimation de  $\Delta$ , on peut également prendre (Joreskog, 1963), lorsque  $\mathbf{S}$  est une matrice des corrélations,  $\delta_{jj} = 1 - R_j^2$ , où la quantité  $R_j^2$  est le coefficient de corrélation multiple de la variable  $j$  avec toutes les autres. Ainsi, une variable très peu corrélée avec les autres aura une variance spécifique forte. Une variable qui peut s'exprimer comme combinaison linéaire des autres aura une variance spécifique nulle.

Notons que  $1 - R_j^2$  est l'inverse du  $j^{\text{ème}}$  élément diagonal de  $\mathbf{S}^{-1}$ .

## Analyse factorielle discriminante

On désigne sous le nom d'*analyse discriminante* une famille de techniques destinées à classer (affecter à des classes préexistantes) des individus caractérisés par un certain nombre de variables numériques ou nominales.

L'origine de cette méthode remonte aux travaux de Fisher (1936) ou, de façon moins directe, à ceux de Mahalanobis (1936). Elle est une des techniques d'analyse multidimensionnelle les plus utilisées en pratique (Credit-scoring, diagnostic automatique, contrôle de qualité, prévision de risques, reconnaissance des formes).

L'*analyse factorielle discriminante* ou *analyse linéaire discriminante*, est une méthode à la fois descriptive et prédictive, qui donne lieu, comme les méthodes factorielles présentées au chapitre 1, à des calculs d'axes principaux. Elle peut être considérée comme une extension de la régression multiple dans le cas où la variable à expliquer est nominale et constitue la variable de partition. Ces deux techniques constituent d'ailleurs des cas particuliers de l'analyse canonique (cf. section 3.1).

Nous ne présenterons pas toutes les techniques d'analyse discriminante qui donnent lieu à une littérature presque aussi étendue que la régression et le modèle linéaire. Nous renvoyons le lecteur à des ouvrages spécifiques sur la question, notamment l'ouvrage de Tomassone et *al.* (1988) et les ouvrages édités par Celeux (1990) (discrimination à partir de variables continues) et Celeux et Nakache (1994) (discrimination à partir de variables qualitatives)<sup>1</sup>.

### 3.3.1 Formulation du problème et notations

On dispose de  $n$  individus ou observations décrits par un ensemble de  $p$  variables  $(x_1, x_2, \dots, x_p)$  et répartis en  $q$  classes définies a priori par la variable  $y$  nominale à  $q$  modalités<sup>2</sup>.

<sup>1</sup> Signalons dans la littérature de langue anglaise l'ouvrage de synthèse (riche de plus de 1200 références) de McLachlan (1992) et les articles, également de synthèse, de Lachenbruch et Goldstein (1979), de Gnanadesikan (1989) ; parmi les manuels classiques généralistes qui traitent de l'analyse discriminante, Anderson (1958, 2nd ed. 1984), Cacoullos (1973), Krishnaiah et Kanal (1982) ; parmi les manuels plus spécialisés, Goldstein et Dillon (1978), Hand (1981). Dans le domaine des méthodes statistiques de la reconnaissance des formes, outre l'ouvrage précité de McLachlan, les ouvrages de base sont Fukunaga (1972), Duda et Hart (1973), Devijver et Kittler (1982). Agrawala (1977) contient des réimpressions de références historiques.

<sup>2</sup> Dans ce chapitre, le vecteur  $y$  a des composantes entières donnant les numéros des classes, et  $Y$  désigne le tableau disjonctif d'ordre  $(n, q)$  correspondant.

L'analyse discriminante se propose dans un premier temps de séparer au mieux les  $q$  classes à l'aide des  $p$  variables explicatives. Dans un deuxième temps, elle cherche à résoudre le problème de l'affectation d'individus nouveaux, caractérisés par les  $p$  variables, à certaines classes déjà identifiées sur l'échantillon des  $n$  individus (appelé *échantillon d'apprentissage*).

On distingue par conséquent deux démarches successives, d'ordre descriptif puis décisionnel :

- chercher des fonctions linéaires discriminantes sur l'échantillon d'apprentissage de taille  $n$  qui sont les combinaisons linéaires des variables explicatives  $(x_1, x_2, \dots, x_p)$  dont les valeurs séparent au mieux les  $q$  classes.
- connaître la classe d'affectation de  $n'$  nouveaux individus décrits par les variables explicatives  $(x_1, x_2, \dots, x_p)$ . Il s'agit ici d'un problème de *classement* dans des classes préexistantes, par opposition au problème de *classification* (traité au chapitre 2) qui consiste à construire des classes les plus homogènes possibles dans un échantillon.

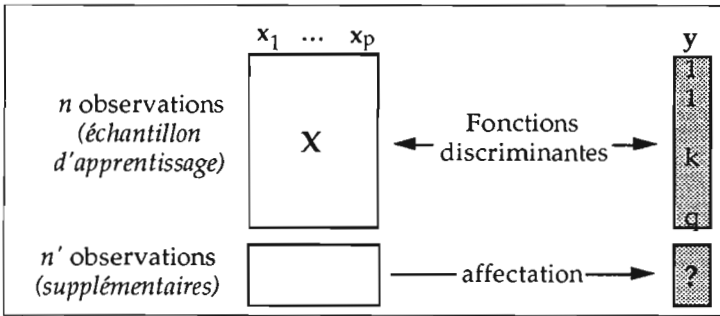


Figure 3.3 - 1  
Principe de l'analyse discriminante

Considérons pour fixer les idées le tableau de données  $(200, 30)$  qui contient, pour  $n = 200$  malades, les valeurs de  $p = 30$  variables issues d'analyses biologiques et d'exams cliniques. Il existe par ailleurs une partition de ces 200 malades selon  $q = 3$  catégories de diagnostics réalisés après des interventions beaucoup plus coûteuses que les 30 mesures précédentes. On se pose la question suivante : étant donné des patients supplémentaires (en nombre  $n'$ ) sur lequel on réalise les 30 analyses et exams, peut-on *prévoir* leurs catégories de diagnostic ? La question répond ici à un besoin pratique<sup>1</sup> :

<sup>1</sup> Les exemples les plus classiques d'analyse discriminante appartiennent sans doute au domaine médical (aide au diagnostic, aide à la décision en matière d'intervention) mais de nombreuses applications se développent dans le domaine du scoring bancaire (prévision de l'éventuelle défaillance d'un débiteur), du contrôle de qualité (prévision de qualité d'un produit en agro-industrie à partir de mesures externes) et surtout de la reconnaissance des formes (reconnaisances de caractères manuscrits ou d'images-radar, etc.).

est-ce-que des mesures nombreuses mais d'accès facile peuvent contenir une information sur un phénomène ou un état plus difficile à identifier ?

Soit le tableau des données  $X$  à  $n$  lignes (individus ou observations) et  $p$  colonnes (variables), de terme général  $x_{ij}$ . Les  $n$  individus sont partitionnés en  $q$  classes. Chaque classe  $k$  caractérise un sous-nuage  $I_k$  de  $n_k$  individus  $i$  avec :

$$\sum_{k=1}^q n_k = n$$

Par  $\bar{x}_{kj}$  on désigne la moyenne de la variable  $x_j$  dans la classe  $k$ . C'est la  $j^{\text{ème}}$  coordonnée du centre de gravité  $G_k$  du sous-nuage  $I_k$  :

$$\bar{x}_{kj} = \frac{1}{n_k} \sum_{i \in I_k} x_{ij} = G_{kj}$$

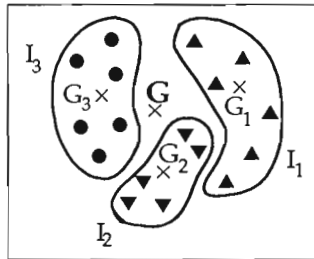


Figure 3.3 - 2  
Représentation du nuage des individus partitionnés

La moyenne de la variable  $x_j$  sur l'ensemble des individus qui correspond à la  $j^{\text{ème}}$  coordonnée du centre de gravité  $G$  du nuage des individus vaut :

$$\bar{x}_j = \frac{1}{n} \sum_{i=1}^n x_{ij} = \sum_{k=1}^q \frac{n_k}{n} \bar{x}_{kj} = G_j$$

### 3.3.2 Fonctions linéaires discriminantes

L'analyse factorielle discriminante consiste à rechercher les combinaisons linéaires de  $p$  variables explicatives ( $x_1, x_2, \dots, x_p$ ), généralement continues, qui permettent de séparer au mieux les  $q$  classes.

La première combinaison linéaire sera celle dont la variance entre les classes (inter-classes) est maximale, afin d'exalter les différences entre les classes, et dont la variance à l'intérieur des classes (intra-classes) minimale pour que l'étendue dans les classes soit délimitée. Puis, parmi les combinaisons linéaires non corrélées à la première, on recherchera celle qui discrimine le mieux les classes, etc.

Ces combinaisons linéaires seront les *fonctions linéaires discriminantes*.

Désignons par  $a(i)$  la valeur, pour l'individu  $i$ , d'une combinaison linéaire à des  $p$  variables préalablement centrées :

$$a(i) = \sum_{j=1}^p a_j (x_{ij} - \bar{x}_j)$$

La variance  $var(\mathbf{a})$  de la nouvelle variable synthétique  $a(i)$  vaut, puisque  $a(i)$  est centrée :

$$var(\mathbf{a}) = \frac{1}{n} \sum_{i=1}^n a^2(i) = \frac{1}{n} \sum_{i=1}^n \left[ \sum_{j=1}^p a_j (x_{ij} - \bar{x}_j) \right]^2$$

$$var(\mathbf{a}) = \frac{1}{n} \sum_{i=1}^n \sum_{j=1}^p \sum_{j'=1}^p a_j a_{j'} (x_{ij} - \bar{x}_j)(x_{ij'} - \bar{x}_{j'})$$

En intervertissant les sommations et en posant :

$$t_{jj'} = \frac{1}{n} \sum_{i=1}^n (x_{ij} - \bar{x}_j)(x_{ij'} - \bar{x}_{j'}) = cov(x_j, x_{j'})$$

la variance de la combinaison des variables  $\mathbf{a}$  peut s'écrire :

$$var(\mathbf{a}) = \sum_{j=1}^p \sum_{j'=1}^p a_j a_{j'} cov(x_j, x_{j'}) = \mathbf{a}' \mathbf{T} \mathbf{a}$$

où  $\mathbf{a}$  désigne le vecteur dont les  $p$  composantes sont  $a_1, \dots, a_p$  et  $\mathbf{T}$  désigne la matrice des covariances des  $p$  variables, de terme général  $t_{jj'}$ .

Nous allons montrer que la variance de  $\mathbf{a}$  se décompose en variance intra-classes et en variance inter-classes, ce qui correspond à une décomposition analogue de la matrice des covariances  $\mathbf{T}$ .

### **a – Décomposition de la matrice de covariance**

La covariance totale entre deux variables  $x_j$  et  $x_{j'}$  s'écrit :

$$cov(x_j, x_{j'}) = \frac{1}{n} \sum_{k=1}^q \left[ \sum_{i \in I_k} (x_{ij} - \bar{x}_j)(x_{ij'} - \bar{x}_{j'}) \right] = t_{jj'}$$

Comme en *analyse de la variance*, nous allons décomposer  $cov(x_j, x_{j'})$  en somme de covariances *intra-classes* (à l'intérieur des classes) et covariances *inter-classes* (entre les classes).

Pour cela nous partirons de l'identité, pour  $i, j, k$  :

$$(x_{ij} - \bar{x}_j) = (x_{ij} - \bar{x}_{kj}) + (\bar{x}_{kj} - \bar{x}_j)$$

La somme entre crochets dans la formule de la covariance se décompose alors en quatre termes, dont deux sont nuls.

En effet, par définition de  $\bar{x}_{kj}$  :

$$\sum_{i \in I_k} (x_{ij} - \bar{x}_{kj})(x_{kj'} - \bar{x}_{j'}) = (x_{kj'} - \bar{x}_{j'}) \sum_{i \in I_k} (x_{ij} - \bar{x}_{kj}) = 0$$

de façon analogue, les sommes ci-dessous s'annulent :

$$\sum_{i \in I_k} (x_{kj} - \bar{x}_j)(x_{ij'} - \bar{x}_{kj'}) = 0$$

Il reste la formule dite *formule de décomposition de Huyghens* (ou équation d'analyse de la variance) :

$$t_{jj'} = d_{jj'} + e_{jj'}$$

avec :

$$d_{jj'} = \frac{1}{n} \sum_{k=1}^q \sum_{i \in I_k} (x_{ij} - \bar{x}_{kj})(x_{ij'} - \bar{x}_{kj'})$$

$$e_{jj'} = \sum_{k=1}^q \frac{n_k}{n} (\bar{x}_{kj} - \bar{x}_j)(\bar{x}_{kj'} - \bar{x}_{j'})$$

Ces  $p^2$  relations se notent sous forme matricielle<sup>1</sup> :

$$\mathbf{T} = \mathbf{D} + \mathbf{E} \quad [3.3-1]$$

Ainsi, la variance d'une combinaison linéaire  $\mathbf{a}$  des variables se décompose d'après la relation [3.3-1] en variance interne et variance externe :

$$\mathbf{a}'\mathbf{T}\mathbf{a} = \mathbf{a}'\mathbf{D}\mathbf{a} + \mathbf{a}'\mathbf{E}\mathbf{a} \quad [3.3-2]$$

Rappelons que, parmi toutes les combinaisons linéaires des variables, on cherche celles qui ont une variance intra-classes minimale et une variance inter-classes maximale. En projection sur l'axe discriminant  $\mathbf{a}$ , chaque sous-nuage doit être, dans la mesure du possible, à la fois bien regroupé et bien séparé des autres sous-nuages.

Il s'agit donc de chercher  $\mathbf{a}$  tel que le quotient  $\mathbf{a}'\mathbf{E}\mathbf{a}/\mathbf{a}'\mathbf{D}\mathbf{a}$  soit maximal (ou  $\mathbf{a}'\mathbf{D}\mathbf{a}/\mathbf{a}'\mathbf{E}\mathbf{a}$  minimal).

D'après la relation [3.3-2] il est équivalent de minimiser  $\mathbf{a}'\mathbf{T}\mathbf{a}/\mathbf{a}'\mathbf{E}\mathbf{a}$  ou de rendre maximal  $f(\mathbf{a})$  tel que :

$$f(\mathbf{a}) = \frac{\mathbf{a}'\mathbf{E}\mathbf{a}}{\mathbf{a}'\mathbf{T}\mathbf{a}}$$

### **b – Calcul des fonctions linéaires discriminantes**

La fonction  $f(\mathbf{a})$  à maximiser est le rapport de la variance inter-classes à la variance totale. Cette fonction étant homogène de degré 0 en  $\mathbf{a}$  (invariante si

<sup>1</sup> La matrice des covariances Totale  $\mathbf{T}$  se décompose en une matrice d'inertie intra-classes  $\mathbf{D}$  (*Dans* les classes) et une matrice d'inter-classes  $\mathbf{E}$  (*Entre* les classes).

$\mathbf{a}$  est changé en  $\xi \mathbf{a}$ ,  $\xi$  étant un scalaire quelconque), il est équivalent de chercher le maximum de la forme quadratique  $\mathbf{a}'\mathbf{E}\mathbf{a}$  sous la *contrainte* quadratique  $\mathbf{a}'\mathbf{T}\mathbf{a} = 1$ .

Ceci conduit à la relation<sup>1</sup> :

$$\mathbf{E}\mathbf{a} = \lambda\mathbf{T}\mathbf{a} \quad [3.3 - 3]$$

Lorsque la matrice des covariances  $\mathbf{T}$  est inversible, on obtient :

$$\mathbf{T}^{-1}\mathbf{E}\mathbf{a} = \lambda\mathbf{a}$$

$\mathbf{a}$  est vecteur propre de  $\mathbf{T}^{-1}\mathbf{E}$  relatif à la plus grande valeur propre  $\lambda$ .

En prémultipliant les deux membres de [3.3 - 3] par le vecteur  $\mathbf{a}'$  on constate que  $\mathbf{a}'\mathbf{E}\mathbf{a}$ , le maximum cherché, n'est autre que  $\lambda$ .

La plus grande valeur propre  $\lambda$ , quotient de la variance *externe* de la fonction discriminante par la variance *totale*, est inférieure à 1 d'après la relation [3.3 - 1]. On l'appelle quelquefois *pouvoir discriminant* de la fonction  $\mathbf{a}$ .

### Remarque

En rendant maximum le quotient  $\mathbf{b}'\mathbf{E}\mathbf{b}/\mathbf{b}'\mathbf{D}\mathbf{b}$  les combinaisons linéaires discriminantes  $\mathbf{b}$  seraient alors les vecteurs propres de la matrice  $\mathbf{D}^{-1}\mathbf{E}$  où la matrice  $\mathbf{D}^{-1}$  définit la *métrique de Mahalanobis*. La valeur propre  $\mu$  correspondant, solution de  $\mathbf{D}^{-1}\mathbf{E}\mathbf{b} = \mu\mathbf{b}$  est reliée à  $\lambda$  par la formule :

$$\mu = \frac{\lambda}{1 - \lambda}$$

On a évidemment  $\mu \geq \lambda$ , puisque la variance interne est toujours inférieure à la variance totale.

Le vecteur  $\mathbf{b}$  est comme  $\mathbf{a}$  solution de l'équation [3.3 - 3] mais doit respecter la contrainte  $\mathbf{b}'\mathbf{D}\mathbf{b} = 1$ .

Les vecteurs  $\mathbf{a}$  et  $\mathbf{b}$  sont liés par la relation<sup>2</sup> :

$$\mathbf{a} = (\sqrt{1 - \lambda}) \mathbf{b}$$

## c – Diagonalisation d'une matrice symétrique

La matrice  $\mathbf{T}^{-1}\mathbf{E}$  n'est pas symétrique. Mais il est possible de se ramener à la diagonalisation d'une matrice  $(q, q)$  symétrique. (Rappelons que  $p$  est le nombre de variables et  $q$  le nombre de classes avec dans la plupart des applications  $q < p$ ).

<sup>1</sup> Comme en analyse générale (section 1.1) ou en analyse canonique (section 3.1), nous sommes conduits à annuler le vecteur des dérivées partielles du *lagrangien*  $\mathcal{L} = \mathbf{a}'\mathbf{E}\mathbf{a} - \lambda(\mathbf{a}'\mathbf{T}\mathbf{a} - 1)$  par rapport à  $\mathbf{a}$ , ce qui donne la relation :  $2\mathbf{E}\mathbf{a} - 2\lambda\mathbf{T}\mathbf{a} = 0$ , d'où finalement  $\mathbf{E}\mathbf{a} = \lambda\mathbf{T}\mathbf{a}$ .

<sup>2</sup> Posant  $\mathbf{a} = \xi \mathbf{b}$ , les deux relations  $\mathbf{a}'\mathbf{E}\mathbf{a} = \lambda$  et  $\mathbf{b}'\mathbf{E}\mathbf{b} = \mu$  conduisent à la relation  $\xi^2 \mathbf{b}'\mathbf{E}\mathbf{b} = \lambda$ , d'où :  $\xi^2 \mu = \lambda$  et  $\xi = \sqrt{1 - \lambda}$

En effet la matrice  $E$ , de terme général :

$$e_{jj'} = \sum_{k=1}^q \frac{n_k}{n} (\bar{x}_{kj} - \bar{x}_j)(\bar{x}_{kj'} - \bar{x}_{j'})$$

est le produit d'une matrice  $C$  à  $p$  lignes et  $q$  colonnes par sa transposée; cette matrice  $C$  a pour terme général :

$$c_{jk} = \sqrt{\frac{n_k}{n}} (\bar{x}_{kj} - \bar{x}_j) \quad [3.3 - 4]$$

Avec la décomposition  $E = CC'$ , la relation [3.3 - 3] s'écrit :

$$CC'a = \lambda Ta$$

Posons :

$$a = T^{-1}Cw \quad [3.3 - 5]$$

cette relation s'écrit alors :

$$CC'T^{-1}Cw = \lambda Cw \quad [3.3 - 6]$$

Il est clair que tout vecteur propre  $w$  relatif à une valeur propre  $\lambda$  (différente de 0) de la matrice symétrique  $C'T^{-1}C$  d'ordre  $(q, q)$  vérifie également [3.3-6]. Le vecteur  $a$  et le scalaire  $\lambda$  vérifient alors la relation [3.3 - 3]. Il suffit en pratique d'effectuer la diagonalisation de cette matrice symétrique<sup>1</sup>, puis d'en déduire  $a$  par la transformation [3.3 - 5].

### 3.3.3 Cas de deux classes : équivalence avec la régression multiple

Lorsque la variable  $y$  ne prend que deux valeurs, chacune caractérisant une classe, des simplifications apparaissent. L'analyse discriminante est alors un cas particulier de la régression multiple.

On repérera les deux classes par les indices 1 et 2. La matrice des covariances  $E$  entre classes a pour terme général :

$$e_{jj'} = \frac{n_1}{n} (\bar{x}_{1j} - \bar{x}_j)(\bar{x}_{1j'} - \bar{x}_{j'}) + \frac{n_2}{n} (\bar{x}_{2j} - \bar{x}_j)(\bar{x}_{2j'} - \bar{x}_{j'})$$

avec :

$$\bar{x}_j = \frac{n_1}{n} \bar{x}_{1j} + \frac{n_2}{n} \bar{x}_{2j}$$

En remplaçant  $\bar{x}_j$  par sa valeur et en tenant compte du fait que  $n_1 + n_2 = n$ , on trouve :

<sup>1</sup>De plus cette matrice symétrique d'ordre  $(q, q)$  sera en général notablement plus petite que la matrice non-symétrique  $T^{-1}E$  d'ordre  $(p, p)$ .



$$e_{jj'} = \frac{n_1 n_2}{n^2} (\bar{x}_{1j} - \bar{x}_{2j})(\bar{x}_{1j'} - \bar{x}_{2j'})$$

La matrice symétrique  $\mathbf{E}$  d'ordre  $(p, p)$  et de rang 1, peut être considérée comme le produit d'une matrice colonne  $\mathbf{c}$  par sa transposée :

$$\mathbf{E} = \mathbf{c}\mathbf{c}'$$

avec :

$$c_j = \frac{\sqrt{n_1 n_2}}{n} (\bar{x}_{1j} - \bar{x}_{2j})$$

La relation [3.3 - 3] s'écrit alors :

$$\mathbf{T}^{-1}\mathbf{c}\mathbf{c}'\mathbf{a} = \lambda\mathbf{a}$$

Prémultiplions les deux membres par  $\mathbf{c}'$  :

$$[\mathbf{c}'\mathbf{T}^{-1}\mathbf{c}]\mathbf{c}'\mathbf{a} = \lambda\mathbf{c}'\mathbf{a}$$

La quantité entre crochets est un scalaire, égal par conséquent à  $\lambda$  qui est ici une valeur propre unique car  $\mathbf{E}$  est de rang 1.

Cette valeur propre vaut donc :  $\lambda = \mathbf{c}'\mathbf{T}^{-1}\mathbf{c}$

$\lambda$  est appelée *distance généralisée* entre les deux classes ou encore "*Distance de Mahalanobis*". Le vecteur propre correspondant :

$$\mathbf{a} = \mathbf{T}^{-1}\mathbf{c}$$

est l'unique fonction discriminante.

Considérons un vecteur  $\mathbf{w}$  à  $n$  composantes, défini par :

$$w_i = \begin{cases} \sqrt{n_1/n_2} & \text{si le } i^{\text{ème}} \text{ individu appartient à la classe 1} \\ -\sqrt{n_2/n_1} & \text{s'il appartient à la classe 2} \end{cases}$$

La régression multiple expliquant  $\mathbf{w}$  par les colonnes de  $\mathbf{X}$  conduit au vecteur de coefficients noté ici  $\mathbf{b}$  :

$$\mathbf{b} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{w}, \quad \text{avec : } \frac{1}{n}\mathbf{X}'\mathbf{X} = \mathbf{T}$$

On vérifie que :

$$\frac{1}{n}\mathbf{X}'\mathbf{w} = \mathbf{c}$$

d'où :

$$\mathbf{b} = \mathbf{T}^{-1}\mathbf{c}$$

Le vecteur des *coefficients de régression*  $\mathbf{b}$  coïncide par conséquent avec le vecteur des composantes de la *fonction discriminante*  $\mathbf{a}$  calculé précédemment <sup>1</sup>.

<sup>1</sup> Notons cependant que les tests et autres procédures statistiques seront d'une autre nature.

### 3.3.4 Lien avec d'autres méthodes

L'analyse factorielle discriminante est un cas particulier de l'analyse canonique lorsque l'un des deux ensembles de variables est formé par les indicatrices d'une partition. Lorsque les deux ensembles sont formés de variables indicatrices, on retrouve l'analyse des correspondances, qui est une double analyse discriminante (cf. aussi § 3.1.3). On peut également présenter la méthode comme une analyse en axes principaux du nuage des points moyens dans une métrique particulière.

#### a – L'analyse canonique

Comme en analyse des correspondances multiples, la variable nominale à  $q$  classes sera représentée par un codage disjonctif complet. On construit ainsi une matrice  $Y$  à  $n$  lignes et  $q$  colonnes de terme général  $y_{ik}$  valant 1 si l'individu  $i$  appartient à la classe  $k$  ou 0 sinon. Autrement dit, nous ajoutons aux variables initiales  $X$  des variables *artificielles*  $Y$  qui indiquent l'appartenance aux diverses classes.

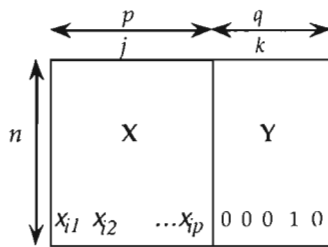


Figure 3.3 - 3  
Tableau de données  $[X, Y]$

Les  $p$  colonnes des variables observées du sous-tableau  $X$  seront centrées et notées  $\hat{X}$ . Nous poserons :

$$\hat{x}_{ij} = x_{ij} - \bar{x}_j$$

Notons qu'à la différence de l'analyse canonique, les colonnes de  $Y$  ne sont pas centrées : la somme des éléments de la  $k^{\text{ème}}$  colonne vaut  $n_k$ .

L'analyse canonique du tableau  $[\hat{X}, Y]$  conduit à chercher le vecteur propre  $a$  de la matrice  $N$  (formule [3.1 - 4] du § 3.1.2.a) :

$$N = (\hat{X}'\hat{X})^{-1}\hat{X}'Y(Y'Y)^{-1}Y'\hat{X}$$

Explicitons les différents éléments de la matrice  $N$  en tenant compte de la nature particulière des colonnes de  $Y$  :

- la matrice  $\frac{1}{n}\hat{X}'\hat{X}$  n'est autre que la matrice des covariances empiriques désignée précédemment par  $T$ .

- la matrice  $\mathbf{D} = \mathbf{Y}'\mathbf{Y}$  est diagonale et son  $k^{\text{ème}}$  élément diagonal vaut  $n_k$ , effectif de la  $k^{\text{ème}}$  classe<sup>1</sup>.

- la matrice à  $p$  lignes et  $q$  colonnes  $\mathbf{H} = \hat{\mathbf{X}}'\mathbf{Y}$  a pour terme général :

$$h_{jk} = \sum_{i=1}^n \hat{x}_{ij} y_{ik} = \sum_{i=1}^n (x_{ij} - \bar{x}_j) y_{ik} = \sum_{i \in I_k} (x_{ij} - \bar{x}_j) = n_k (\bar{x}_{kj} - \bar{x}_j)$$

En vertu de la relation [3.3 - 4], on peut écrire :

$$h_{jk} = \sqrt{nn_k} c_{jk}$$

soit :

$$\mathbf{H} = \hat{\mathbf{X}}'\mathbf{Y} = \sqrt{n} \mathbf{C} (\mathbf{Y}'\mathbf{Y})^{1/2}$$

Ces dernières remarques nous permettent d'écrire :

$$\hat{\mathbf{X}}'\mathbf{Y}(\mathbf{Y}'\mathbf{Y})^{-1}\mathbf{Y}'\hat{\mathbf{X}} = n \mathbf{C}\mathbf{C}' = n \mathbf{E}$$

puisque :

$$(\hat{\mathbf{X}}'\hat{\mathbf{X}})^{-1} = \frac{1}{n} \mathbf{T}^{-1}$$

la matrice  $\mathbf{N}$  devient finalement :  $\mathbf{N} = \mathbf{T}^{-1}\mathbf{E}$

et le vecteur  $\mathbf{a}$  cherché vérifie bien la relation [3.3 - 3] :

$$\mathbf{E}\mathbf{a} = \lambda \mathbf{T}\mathbf{a}$$

Nous pouvons également noter que l'on a, pour les deux types d'analyse, la même contrainte de normalisation :

$$\mathbf{a}'\mathbf{T}\mathbf{a} = 1$$

Il y a donc coïncidence entre variable canonique et fonction discriminante. L'analyse discriminante apparaît ainsi comme un cas particulier de l'analyse canonique (sans centrage préalable des variables indicatrices) lorsque l'un des deux ensembles est constitué de vecteurs booléens décrivant la partition de l'ensemble des individus.

## b – L'analyse des correspondances

Lorsque le sous-tableau  $\mathbf{X}$  décrit lui aussi une partition en  $p$  classes, les résultats du paragraphe précédent montrent immédiatement que l'analyse des correspondances est un cas particulier de l'analyse factorielle discriminante.

---

<sup>1</sup> En effet, on a la relation  $\sum_{i=1}^n y_{ik} y_{ik}' = \delta_{kk'} n_k$  car l'individu  $i$  appartient soit à la classe  $k$ , soit à la classe  $k'$ ;  $\delta_{kk'} = 1$  si  $k=k'$  et vaut 0 sinon. Pour  $k=k'$ , il y aura autant de termes non nuls dans la somme que d'individus dans la classe  $k$ .

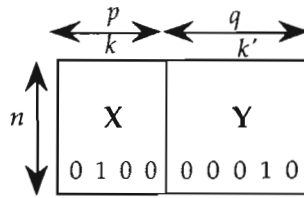


Figure 3.3 - 4  
Tableau de données [X,Y]

Les deux sous-tableaux X d'ordre  $(n,p)$  et Y d'ordre  $(n,q)$  de la matrice des données [X,Y] sont formés de variables indicatrices et jouent maintenant des rôles analogues. Dans ce cas, les matrices  $X'X$  et  $Y'Y$  sont diagonales et ont pour  $k^{\text{ème}}$  élément les effectifs de la classe  $k$  de chacune des partitions ; la matrice  $X'Y$  n'est autre que le tableau de contingence d'ordre  $(p,q)$  croisant les deux partitions  $P_X$  et  $P_Y$ .

Conformément aux conventions adoptées en analyse des correspondances, on notera<sup>1</sup> :

-  $f_{k.}$ , le  $k^{\text{ème}}$  élément diagonal de la matrice  $\frac{1}{n}X'X (=D_p)$ , ( $k \leq p$ )

-  $f_{.k'}$ , le  $k'^{\text{ème}}$  élément diagonal de la matrice  $\frac{1}{n}Y'Y (=D_q)$ , ( $k' \leq q$ )

-  $f_{kk'}$ , l'élément générique de la matrice  $\frac{1}{n}X'Y (=F)$ , d'ordre  $(p,q)$

Rappelons les formules établies au paragraphe 3.1.2 reliant les variables canoniques :

$$\mathbf{a} = \frac{1}{\lambda}(X'X)^{-1}X'Y\mathbf{b} \quad \text{et} \quad \mathbf{b} = \frac{1}{\lambda}(Y'Y)^{-1}Y'X\mathbf{a}$$

Leurs composantes s'écrivent :

$$a_k = \frac{1}{\lambda} \sum_{k'=1}^q \frac{f_{kk'}}{f_{k.}} b_{k'} \quad \text{et} \quad b_{k'} = \frac{1}{\lambda} \sum_{k=1}^p \frac{f_{kk'}}{f_{.k'}} a_k$$

On reconnaît, sous cette forme, les relations barycentriques de l'analyse des correspondances [1.3 - 12] et [1.3 - 13] reliant les coordonnées des deux nuages sur un même axe factoriel.

Cette identité suffit à établir qu'une analyse des correspondances est une analyse canonique particulière où les tableaux X et Y contiennent les variables indicatrices de deux partitions<sup>2</sup>.

<sup>1</sup>  $n$  est ici l'effectif global alors qu'il était désigné par  $k$  à la section 1.3.

<sup>2</sup> La première racine canonique  $\lambda^2$  est l'homologue de la première valeur propre, notée  $\lambda$  précédemment pour l'analyse des correspondances.

Les sous-espaces  $V_X$  et  $V_Y$  ont maintenant en commun la première bissectrice<sup>1</sup> de  $\mathbb{R}^n$ ; leur plus petit angle est donc nul.

Son cosinus (=1) est la valeur propre triviale déjà rencontrée en analyse des correspondances lorsque l'analyse est faite par rapport à l'origine et non par rapport au centre de gravité.

On a alors  $\lambda = 1$ ,  $a_i = 1$  et  $b_j = 1$ , pour tout  $i$  et tout  $j$  dans les relations écrites ci-dessus. Le fait de centrer le tableau  $X$  revient à projeter les points-colonnes sur le sous-espace orthogonal à la première bissectrice.

Cette opération ne modifie donc pas les variables canoniques non triviales.

L'analyse des correspondances apparaît comme une *double analyse discriminante* car chacun des blocs dans  $[X, Y]$  décrit une partition et aucun d'entre eux n'est privilégié. Les fonctions linéaires discriminantes coïncident avec les facteurs de l'analyse des correspondances<sup>2</sup> du tableau de contingence d'ordre  $(p, q)$  croisant les deux partitions.

### c – Une analyse en axes principaux avec une métrique particulière

L'analyse factorielle discriminante peut être considérée comme une analyse générale du nuage des  $q$  centres de gravité des classes  $k$  munis des masses  $n_k/n$  et avec la métrique  $T^{-1}$  ou la métrique  $D^{-1}$  dite de *Mahalanobis*.

Le nombre d'axes discriminants est égal à  $q - 1$  dans le cas où  $n > p > q$ .

Il suffit en effet de se reporter au paragraphe 3.3.2.c précédent où est intervenu pour la première fois le tableau  $C$  des moyennes centrées.

L'analyse générale de ce tableau  $C$  avec la métrique  $T^{-1}$ , selon les résultats du paragraphe 1.1.6.a du chapitre 1 (analyse générale avec une métrique quelconque : ici,  $X = C$ ,  $M = T^{-1}$  et  $N = I$ ) conduit, pour trouver l'axe factoriel  $u$ , à la relation :

$$C'CT^{-1}u = \lambda u$$

Posant  $T^{-1}u = a$ , où  $a$  est le facteur (opérateur projection) correspondant à l'axe factoriel  $u$  :

$$C'Ca = \lambda Ta$$

De la même façon, avec la métrique  $D^{-1}$ , on obtient :

$$C'Ca = \lambda Da$$

<sup>1</sup> La somme des colonnes de  $X$  et la somme des colonnes de  $Y$  constituent le vecteur dont toutes les composantes valent 1.

<sup>2</sup> Cette présentation permet de montrer directement que les valeurs propres de l'analyse des correspondances, étant des *coefficients de corrélation canonique* (ou des pouvoirs discriminants) sont inférieures ou égales à 1. De plus on pourra interpréter les valeurs propres de l'analyse des correspondances en terme de *pouvoir discriminant* des facteurs (axes factoriels) vis-à-vis des partitions étudiées.

Choisir la métrique  $D^{-1}$  pour analyser le nuage des points-moyens, c'est considérer comme équidistantes du centre  $j$  (par exemple) des zones équiprobables (au sens des ellipsoïdes de densité) d'équation :

$$(x - \bar{x}_j)' D^{-1} (x - \bar{x}_j) = \text{constante}$$

Grâce à cette métrique, la distance est interprétée en terme de "vraisemblance d'appartenance".

Ainsi, sur la figure 3.3 - 5, où sont représentées trois classes ayant mêmes ellipsoïdes de densité (équation ci-dessus,  $D$  étant la matrice des covariances interne commune à chaque groupe), les points A et B sont équidistants (selon la métrique  $D^{-1}$ ) du centre de classe  $G_1$ .

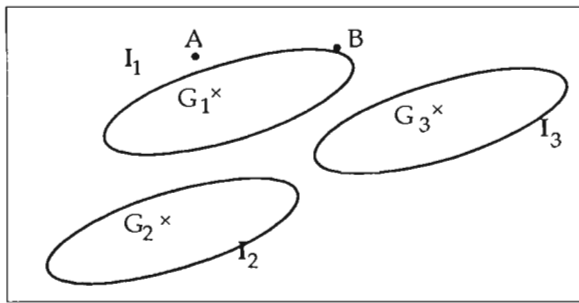


Figure 3.3 - 5  
Illustration de la métrique  $D^{-1}$

Avec la métrique euclidienne usuelle, B serait affecté plutôt à la classe 3 qu'à la classe 1. On voit donc l'intérêt de faire intervenir cette métrique dans l'analyse des centres<sup>1</sup>. Nous reviendrons sur cette question au paragraphe suivant dévolu aux règles d'affectation.

### 3.3.5 Principes des règles d'affectation (ou de classement)

Une fois trouvées les fonctions discriminantes qui séparent au mieux les individus répartis en  $q$  classes, on veut trouver la classe d'affectation d'un nouvel individu, pour lequel on connaît les valeurs des variables  $(x_1, x_2, \dots, x_p)$ .

Une règle simple et géométrique d'affectation est de choisir la classe dont le centre de gravité est le plus proche du point-individu. La métrique

<sup>1</sup> Il est clair que cette métrique prend en compte une certaine anisotropie (orientation préférentielle) de la densité. Elle n'a cependant de sens que si les ellipsoïdes de densité sont les mêmes à l'intérieur de chaque classe. C'est précisément ce qui caractérise l'analyse discriminante linéaire, par opposition à l'analyse discriminante quadratique, qui autorise des densités de formes différentes, et donc des métriques différentes pour chaque classe.

généralement utilisée dans les applications les plus courantes est celle de *Mahalanobis globale* ( $\mathbf{D}^{-1}$ ), ou *locale* ( $\mathbf{D}_k^{-1}$ , où  $\mathbf{D}_k$  est la matrice des covariances internes au groupe  $I_k$ ).

Cette approche purement géométrique ne prend cependant pas en compte les probabilités *a priori* des différentes classes, qui peuvent être très inégales dans certaines applications (prévision de défaillance par exemple, ou diagnostic d'un événement rare). Le modèle bayésien d'affectation permet d'enrichir ce point de vue.

### a – Le modèle bayésien d'affectation

Au moment de l'apprentissage, nous savons que l'individu  $i$  appartient au groupe  $I_k$  (appartenance codée par la valeur :  $y_i = k$ ) et nous calculons une estimation de la probabilité  $P(x_i | I_k)$ , c'est-à-dire la probabilité de  $x_i$  sachant que  $I_k$  est réalisé.

Au moment de l'affectation d'un individu nouveau noté  $x$ , on peut calculer les différents  $P(x | I_k)$  pour  $k = 1, 2, \dots, q$ . Il paraît raisonnable d'affecter  $x$  à la classe  $I_k$  pour laquelle  $P(x | I_k)$  est maximale.

Cependant, ce ne sont pas les probabilités  $P(x | I_k)$  qu'il faudrait connaître mais les probabilités  $P(I_k | x)$ , c'est-à-dire la probabilité du groupe  $I_k$  sachant que  $x$  est réalisé.

Le théorème de Bayes<sup>1</sup> permet de procéder à cette *inversion des probabilités*.

Il exprime  $P(I_k | x)$  en fonction de  $P(x | I_k)$ ,  $P(I_k)$  et  $P(x)$  :

$$P(I_k | x) = \frac{P(x | I_k)P(I_k)}{P(x)}$$

$P(I_k)$  est la probabilité *a priori* du groupe  $k$ .  $P(x)$  s'exprime en fonction de  $P(x | I_k)$  et de  $P(I_k)$  ; d'où la formulation classique du théorème de Bayes :

$$P(I_k | x) = \frac{P(x | I_k)P(I_k)}{\sum_{k=1}^q P(x | I_k)P(I_k)}$$

Le dénominateur est le même pour toutes les classes. La classe d'affectation de  $x$  sera celle pour laquelle le produit  $P(x | I_k) \times P(I_k)$  est maximal. Si les probabilités *a priori*  $P(I_k)$  des classes sont égales pour toutes les valeurs de  $k$ , les classements selon  $P(I_k | x)$  et  $P(x | I_k)$  sont identiques.

<sup>1</sup> Pour un exposé de l'approche bayésienne qui donne un cadre conceptuel spécifique à la théorie de l'estimation et de la décision statistique, voir Robert (1992).

Pour tester l'efficacité des règles d'affectation, on mesure les erreurs de classement par des méthodes de rééchantillonnage, notamment la validation croisée ou le bootstrap (cf. § 4.2.2). Comme dans le cas du modèle linéaire, le choix des variables explicatives est une opération délicate. L'étude de la stabilité des fonctions discriminantes est difficile. Les règles d'affectation ainsi que l'estimation des taux d'erreur de classement dépendent souvent de la taille de l'échantillon d'apprentissage.

### b – Le modèle bayésien dans le cas normal

Notons  $f_k(\mathbf{x})$  la densité de probabilité de  $\mathbf{x}$  connaissant  $I_k$  dans le cas multinormal,  $\mu_k$  et  $\Sigma_k$  désignant respectivement la moyenne et la matrice des covariances théoriques à l'intérieur du groupe  $I_k$  :

$$f_k(\mathbf{x}) = \frac{1}{(2\pi)^{p/2} |\Sigma_k|^{1/2}} e^{-\frac{1}{2}(\mathbf{x} - \mu_k)' \Sigma_k^{-1} (\mathbf{x} - \mu_k)}$$

que l'on préférera écrire :

$$f_k(\mathbf{x}) = (2\pi)^{-p/2} |\Sigma_k|^{-1/2} \exp\left\{-\frac{1}{2}(\mathbf{x} - \mu_k)' \Sigma_k^{-1} (\mathbf{x} - \mu_k)\right\}$$

L'affectation se fera selon la règle :

$$\text{choisir } \hat{k} \text{ tel que } f_{\hat{k}}(\mathbf{x})P(I_{\hat{k}}) = \max_{k \leq q} \{f_k(\mathbf{x})P(I_k)\}$$

ce qui est équivalent à trouver le minimum sur  $k$  de la fonction  $sc_k(\mathbf{x})$  appelée *score discriminant* :

$$sc_k(\mathbf{x}) = (\mathbf{x} - \mu_k)' \Sigma_k^{-1} (\mathbf{x} - \mu_k) + \text{Log} |\Sigma_k| - 2 \text{Log} P(I_k) \quad [3.3 - 7]$$

Dans le cas où les distributions dans chaque classe ont même matrice des covariances (cas illustré par la figure 3.3 - 5), la densité s'écrit :

$$f_k(\mathbf{x}) = (2\pi)^{-p/2} |\Sigma|^{-1/2} \exp\left\{-\frac{1}{2}(\mathbf{x} - \mu_k)' \Sigma^{-1} (\mathbf{x} - \mu_k)\right\}$$

Il suffit alors de prendre pour score discriminant :

$$sc_k(\mathbf{x}) = (\mathbf{x} - \mu_k)' \Sigma^{-1} (\mathbf{x} - \mu_k) - 2 \text{Log} P(I_k) \quad [3.3 - 8]$$

Si de plus les probabilités *a priori*  $P(I_k)$  sont égales, le score discriminant coïncide avec la distance de Mahalanobis :

$$sc_k(\mathbf{x}) = (\mathbf{x} - \mu_k)' \Sigma^{-1} (\mathbf{x} - \mu_k) \quad [3.3 - 9]$$

et la règle bayésienne d'affectation devient la recherche du centre le plus proche selon cette distance.

Le score discriminant donné par la formule [3.3 - 7] correspond à l'*analyse discriminante quadratique*. Les cloisons interclasses données par l'équation  $sc_k(\mathbf{x}) = sc_{k'}(\mathbf{x})$ , ( $k \neq k'$ ), sont en effet des hyperquadriques.



Les scores discriminants donnés par les formules [3.3 - 8] ou [3.3 - 9] correspondent à l'analyse discriminante linéaire. Dans l'équation  $sc_k(\mathbf{x}) = sc_{k'}(\mathbf{x})$ , ( $k \neq k'$ ), les termes du second degré en  $\mathbf{x}$  disparaissent et les cloisons interclasses sont dans ce cas des hyperplans.

Ces hyperplans ont une équation de la forme :

$$\mathbf{x}' \Sigma^{-1}(\mu_{k'} - \mu_k) = \text{constante}$$

Notons que le calcul suppose connus les paramètres théoriques  $\mu_k$  et  $\Sigma_k$ .

Ils suggèrent de substituer en pratique les estimations empiriques aux paramètres théoriques. Cette substitution est également encouragée par l'approche descriptive développée au début de cette section, dans laquelle les distances de Mahalanobis sont apparues de façon naturelle, en cherchant à maximiser le rapport variance externe sur variance interne, sans recours à l'hypothèse de normalité.

Les scores discriminants utilisés en pratique<sup>1</sup>, lorsque l'hypothèse de normalité est plausible, sont donc ceux présentés ici avec utilisation des estimations empiriques des paramètres.

### c – Autres règles d'affectation

Il existe d'autres méthodes de discrimination que celles apparentées à l'analyse factorielle discriminante ou au modèle multinormal. Elles impliquent d'autres règles d'affectations.

Citons, parmi les méthodes les plus utilisées<sup>2</sup> : les méthodes d'estimation non-paramétriques de la densité, connues également sous le nom de méthodes des *noyaux* (de Rosenblatt ou de Parzen), et les méthodes d'affectation (également non-paramétriques) utilisant les *m plus proches voisins*.

#### - Estimation de la densité par noyaux

Une méthode simple de discrimination consisterait à diviser l'espace multidimensionnel de l'échantillon d'apprentissage en cellules de volumes comparables  $v_r$  puis de compter, à l'intérieur de chaque classe  $k$ , ( $k \leq q$ ), les  $n_{rk}$  observations contenues dans chaque cellule  $r$ .

La fréquence  $n_{rk}/n_k$  est une estimation de la probabilité qu'une observation de la catégorie  $k$  appartienne à la cellule  $v_r$ . La règle de Bayes permet alors d'affecter une observation supplémentaire  $\mathbf{x}$  à une catégorie  $k$ , après avoir déterminé la cellule  $v_r$  qui la contient.

<sup>1</sup> Il n'est cependant pas aisé de démontrer l'optimalité de cette démarche intuitive, sauf dans des contextes asymptotiques assez particuliers (cf. Anderson, 1958 ; Friedman, 1989).

<sup>2</sup> D'autres techniques de discrimination seront évoquées plus loin (méthodes neuronales, régression logistique).

Cette méthode est malheureusement impossible à mettre en oeuvre car le nombre de cellules devient vite prohibitif dans un espace à  $p$  dimensions et les échantillons n'ont pas une taille suffisante pour permettre une estimation de fréquence à l'intérieur de chaque cellule.

On peut, pour la classe  $k$ , entourer d'une cellule chaque point observé, de façon à décrire la densité dans l'espace  $\mathbb{R}^p$ . Si le point à affecter  $\mathbf{x}$  tombe à l'intersection de trois cellules de la classe  $k$  par exemple et en dehors des cellules relatives aux autres classes, cela signifiera qu'il est dans une zone de forte densité pour la classe  $k$  et donc qu'il a plus de chance d'appartenir à cette classe qu'aux autres. Cette idée, présentée ici de façon intuitive, est celle des noyaux de Rosenblatt (1956).

Au lieu d'entourer les points de cellules de volumes fixes, on peut les entourer d'une sorte de halo, une zone de densité qui décroît lorsqu'on s'éloigne du point, de façon à procéder à un lissage de cette densité dans l'espace multidimensionnel. C'est la méthode d'estimation directe de la densité par noyaux à laquelle on attache le nom de Parzen (1962).

La méthode des noyaux consiste à estimer la densité de probabilité à l'intérieur de la classe  $k$  dans l'espace  $\mathbb{R}^p$  par une formule du type :

$$f_k(\mathbf{x}) = \frac{1}{h^p n_k} \sum_{i=1}^{n_k} K\left(\frac{\mathbf{x} - \mathbf{x}_i}{h}\right) \quad [3.3 - 10]$$

La fonction  $K(\mathbf{z})$  doit vérifier les relations  $K(\mathbf{z}) \geq 0$ , et  $\int K(\mathbf{z}) d\mathbf{z} = 1$ . Elle pourra être choisie parmi les densités de probabilité usuelles. On note que l'on a bien dans ces conditions :

$$\int f_k(\mathbf{x}) d\mathbf{x} = 1.$$

On utilise souvent la densité de la loi normale sphérique :

$$K(\mathbf{z}) = (2\pi)^{-p/2} \exp\left\{-\frac{1}{2}\mathbf{z}'\mathbf{z}\right\}$$

Le paramètre  $h$  qui intervient dans la formule [3.3 - 10] est la *dimension de la fenêtre*. Dans le cas des noyaux normaux sphériques, il correspond à l'écart-type de la densité locale autour de chaque point. Si  $h$  est petit, le lissage risque d'être mauvais; si  $h$  est trop grand, il risque d'être excessif. Le choix de la dimension de la fenêtre est une des difficultés de ces méthodes d'estimation directe de densité<sup>1</sup>.

- *Règle des  $m$  plus proches voisins* (Fix et Hodges, 1951)

Cette technique, utilisée surtout en reconnaissance des formes, résout d'une autre façon le problème des cellules à densité trop faible: on étend le voisinage autour du point  $\mathbf{x}$  jusqu'à ce qu'il contienne  $m$  points de

<sup>1</sup> Cf. Tomassone et al. (1988), Silverman (1986), Delecroix (1983), Hand (1982). Le paramètre  $h$ , supposé unique dans la formule [3.3 - 10] peut avoir, dans des modèles plus généraux, une valeur différente pour chacune des  $p$  variables et également pour chacune des  $q$  classes.

l'échantillon d'apprentissage. On affecte  $x$  à la classe la plus représentée dans ce voisinage.

Cette méthode est particulièrement simple à mettre en œuvre, surtout dans un processus d'apprentissage progressif, car il n'y a pas de fonctions complexes à recalculer pour prendre en compte les nouveaux individus qui enrichissent l'échantillon d'apprentissage.

Elle nécessite cependant des effectifs importants, des calculs d'affectation coûteux (pour les exigences de la reconnaissance des formes, où le classement s'effectue souvent en temps réel) si les paramètres  $m$  ou  $p$  sont grands<sup>1</sup>.

#### d – Qualité des règles de classement

Il existe un cadre inférentiel paramétrique, apparenté à *l'analyse multidimensionnelle de la variance*, qui permet de tester l'hétérogénéité des classes (test de l'égalité des moyennes  $\mu_k$ , test de l'égalité des matrices de covariances internes  $D_k$ ). Ces tests (mentionnés dans la plupart des manuels de référence cités au début de la section 3.3) dont la robustesse est difficile à établir, sont moins utilisés depuis l'avènement des méthodes non-paramétriques de rééchantillonnage qui seront évoquées à la section 4.2 du chapitre 4.

On esquissera ici, pour les besoins des développements qui suivront, la méthode dite de *validation croisée*.

##### *La validation croisée*

La mesure de la qualité d'une discrimination se fait à partir des pourcentages de bien classés (ou de mal classés) dans chaque classe, et du pourcentage global de bien classés. Cette mesure peut également, dans certaines applications, faire intervenir des coûts de mauvais classement.

On peut calculer un pourcentage de bien classés sur l'échantillon d'apprentissage, ce qui donnera une idée optimiste de la qualité de la discrimination. Ce pourcentage de bien classés augmente avec le nombre de paramètres du modèle, et peut être excellent si le nombre de paramètres est considérable, sans pour cela assurer que le modèle permet de réaliser une prévision correcte. Le pourcentage de mal classés dans ces conditions est appelé le *taux d'erreur apparent* ou encore le *taux d'erreur par resubstitution*.

<sup>1</sup> Il existe des ponts théoriques entre la méthode des  $m$  plus proches voisins et l'estimation directe de densité dans le cas de variables binaires (cf. Fix et Hodges, 1951; Aitchison et Aitken, 1976). Il est également possible, on l'a évoqué, de travailler avec des noyaux adaptatifs, en faisant varier la dimension de la fenêtre  $h$  ou en tenant compte des distances des  $m$  plus proches voisins. Pour une discussion de ces diverses variantes, voir McLachlan (1992). Sur les divers algorithmes de  $m$  plus proches voisins utilisés en reconnaissance des formes, cf. Dubuisson (1990). Sur les problèmes posés par des probabilités *a priori* inégales, cf. Chateau (1994).

La méthode des échantillons-tests<sup>1</sup> recommande d'effectuer la discrimination sur une partie seulement de l'échantillon d'apprentissage (disons 80%) et de tester les règles de discrimination sur les 20% non utilisés.

On peut améliorer le calcul du taux d'erreur en divisant l'échantillon d'apprentissage en  $m$  parties égales, en calculant la règle sur un échantillon partiel formé de  $m-1$  parties, et le taux d'erreur sur la partie restante, ce qui peut être fait de  $m$  façons différentes. Ceci permet donc de calculer un taux d'erreur moyen sur un échantillon aussi important que l'échantillon d'apprentissage.

Plus  $m$  est proche de  $n$ , plus on se rapproche de la situation réelle de classement. La validation croisée<sup>2</sup> correspond au cas  $m = n$ , autrement dit, au cas pour lequel on effectue  $n$  discriminations en excluant à chaque fois une observation. Cette méthode est évidemment coûteuse en calcul mais on peut parfois mettre en œuvre des algorithmes évitant des recalculs complets des fonctions discriminantes<sup>3</sup>.

La minimisation du taux d'erreur par validation croisée peut être utilisée comme critère pour calculer les paramètres de certains modèles de discrimination.

### 3.3.6 Régularisation en analyse discriminante

Comme la régression multiple (dont elle est un cas particulier dans le cas où la variable nominale à prédire n'a que deux catégories, cf. § 3.3.3), l'analyse factorielle discriminante nécessite l'inversion d'une matrice des covariances des prédicteurs (la matrice totale  $\mathbf{T}$  ou la matrice intraclasses  $\mathbf{D}$ ).

Dans le cas de l'analyse discriminante quadratique, le calcul des *distances de Mahalanobis locales* demande d'inverser les matrices de covariances internes à chaque classe  $\mathbf{D}_k$  (dont  $\mathbf{D}$  est une moyenne pondérée).

Ces matrices  $\mathbf{D}$  ou  $\mathbf{T}$ , et surtout les matrices  $\mathbf{D}_k$ , calculées sur un effectif  $n_k$  plus petit que  $n$ , peuvent être mal conditionnées ou même singulières.

C'est systématiquement le cas en analyse discriminante qualitative lorsque les prédicteurs sont des variables nominales codées sous forme disjonctive comme en analyse des correspondances multiples ou en analyse de la variance (cf. § 3.3.7).

<sup>1</sup> On peut faire remonter cette pratique à Highleyman (1962), mais elle a probablement dû être utilisée antérieurement, tant son principe relève du bon sens. Elle a été prônée notamment par Romeder (1973).

<sup>2</sup> Attribuée à Lachenbruch et Mickey, 1968, cette méthode (*cross-validation*) aurait été utilisée dès 1964 par des chercheurs russes, selon Toussaint (1974). Ses propriétés ont été étudiées par Stone (1974) et Geisser (1975). Une revue est faite par Hand (1986).

<sup>3</sup> Cf. par exemple, Celeux (1990) pour le cas des fonctions linéaires discriminantes.

On présentera brièvement ci-dessous une méthode de régularisation proposée par Friedman et la méthode de régularisation par axes principaux déjà proposée pour la régression (§ 3.2.5). Cette méthode a l'avantage de fournir une description préalable de l'espace des prédicteurs et des possibilités ultérieures de filtrage et de sélection de l'information.

### a – Analyse régularisée de Friedman (1989)

Dans cette méthode de régularisation, une nouvelle estimation  $\mathbf{D}_k(\lambda, \gamma)$  est calculée pour chaque matrice des covariances locales  $\mathbf{D}_k$ , qui devient une moyenne pondérée des matrices des covariances globales et locales (rôle du poids  $\lambda$ ) et de la matrice unité (rôle du poids  $\gamma$ ) :

$$\mathbf{D}_k(\lambda, \gamma) = (1 - \gamma)\mathbf{D}_k(\lambda) + \frac{\gamma}{p} \text{tr}[\mathbf{D}_k(\lambda)]\mathbf{I}$$

avec :

$$\mathbf{D}_k(\lambda) = \frac{(1 - \lambda)\mathbf{D}_k + \lambda\mathbf{D}}{(1 - \lambda)n_k + \lambda n}$$

Le scalaire  $\text{tr}[\mathbf{D}_k(\lambda)]$  est la trace de la matrice  $\mathbf{D}_k(\lambda)$ .

La détermination des paramètres  $\lambda$  et  $\gamma$  se fait en optimisant les pourcentages de bien classés obtenus par validation croisée.

Ces techniques donnent des résultats intéressants dans le cas de tableaux de données petits ou moyens, lorsque le problème initial est *mal posé* ( $n \leq p$ ) ou *pauvrement posé* ( $n > p$ , mais encore comparable à  $p$ )<sup>1</sup>.

Dans le cas de grandes matrices clairsemées cependant, l'échelle du phénomène crée de nouveaux problèmes. Il est alors nécessaire de comprendre ce qui se passe dans les espaces de dimension élevée.

Est-il vraiment nécessaire de garder tous les axes principaux ? Est-il possible de filtrer l'information de base caractérisée parfois par un haut niveau de bruit ? L'analyse par axes principaux répond à ces préoccupations.

### b – Analyse régularisée par axes principaux

Du point de vue numérique, la diagonalisation est une opération plus sûre que l'inversion des matrices. La théorie de la perturbation<sup>2</sup> nous apprend que la stabilité des vecteurs propres est une fonction croissante des différences entre valeurs propres consécutives. Dans ce contexte, s'il est nécessaire d'éliminer les dimensions correspondant à des valeurs propres nulles, il peut être aussi avantageux d'éliminer les dimensions

<sup>1</sup> Voir aussi Callant (1991) pour une technique d'estimation des paramètres  $\lambda$  et  $\gamma$ .

<sup>2</sup> Cf. par exemple : Wilkinson (1965); Kato (1966) et les travaux de Escofier et Leroux (1972) utilisant les résultats de ces théories en analyse factorielle.

correspondant aux petites valeurs propres, qui sont très sensibles aux perturbations du tableau de données<sup>1</sup>.

#### - Axes principaux de l'échantillon total

La technique de réduction qui sera utilisée durant la première étape dépend de la nature et des propriétés statistiques des données de base<sup>2</sup>. Une simple décomposition aux valeurs singulières suffit pour une régularisation numérique, si l'on ne désire pas de description de l'espace des prédicteurs.

Les nouvelles coordonnées de l'individu  $i$  sur l'axe principal  $r$  issu de l'analyse de l'échantillon total sont désignées par  $z_{ri}$ ,

$$z_{ri} = u'_r(x_i - \bar{x})$$

où ici  $u_r$  est le  $r^{\text{ème}}$  vecteur propre normalisé de  $T$  matrice des covariances totales correspondant à la valeur propre  $\alpha_r$ ;  $u_r$  est aussi la  $r^{\text{ème}}$  colonne de la matrice  $U$  d'ordre  $(p, r_{max})$  (où  $r_{max}$  est le nombre de valeurs propres retenues).|

La distance euclidienne usuelle dans  $\mathbb{R}^p$  de tout point  $i$  au point-moyen  $G_k$  de la classe  $k$  (le point  $i$  peut ne pas appartenir à la classe  $k$  ni à l'échantillon d'apprentissage) peut s'écrire :

$$d^2(i, G_k) = \sum_{j=1}^p (x_{ij} - \bar{x}_{kj})^2 \quad [3.3 - 11]$$

si  $r_{max} = p'$  ( $p'$  désignant le rang de la matrice de données  $X$ ), cette même distance s'écrit, pour la nouvelle base :

$$d^2(i, G_k) = \sum_{r=1}^{r_{max}} (z_{ir} - \bar{z}_{kr})^2 \quad [3.3 - 12]$$

avec  $\bar{z}_{kr} = u'_r(\bar{x}_k - \bar{x})$ .

La distance de tout point  $i$  au centre  $G_k$  de la classe  $k$  dans la métrique  $T^{-1}$  (intervenant en analyse discriminante linéaire, cf. 3.3.4.c) est telle que :

$$D^2(i, G_k) = \sum_{r=1}^{r_{max}} \frac{(z_{ir} - \bar{z}_{kr})^2}{\alpha_r} \quad [3.3 - 13]$$

On a toujours  $r_{max} \leq p'$ . La distance  $D^2(i, G_k)$  est dite *régularisée* si  $r_{max} < p'$  ou si  $r_{max} = p'$  avec  $p' < \text{Min}(n, p)$ .

<sup>1</sup> Cf. Les travaux de Wold (1976). Benzécri (1977 a) recommande que les analyses discriminantes soient réalisées sur les axes d'une analyse factorielle préalable.

<sup>2</sup> Analyse en composantes principales dans le cas où les prédicteurs sont des variables continues, situation retenue au cours des développements qui précèdent ; mais cette réduction pourra aussi être une analyse des correspondances dans le cas de fréquences ou des correspondances multiples dans le cas de variables nominales.

### - Axes principaux de l'échantillon projeté

Si l'on substitue à la matrice de données  $\mathbf{X}$ , de terme général  $x_{ij}$ , la matrice  $\hat{\mathbf{X}}$  de terme général  $\hat{x}_{ij} = x_{ij} - \bar{x}_{kj}$  où  $k$  est l'indice de la classe  $I_k$  à laquelle appartient l'observation  $i$  et où  $\bar{x}_{kj}$  désigne la moyenne de la variable  $j$  dans cette classe<sup>1</sup>, on est conduit à diagonaliser la matrice  $\mathbf{D}$  (au lieu de  $\mathbf{T}$ ). Les valeurs propres de  $\mathbf{D}$  sont notées  $\hat{\alpha}_r$  et les coordonnées des observations sur les nouveaux axes principaux  $\hat{\mathbf{u}}_r$  sont notées  $\hat{z}_{ir}$ .

La distance de tout point  $i$  au centre  $G_k$  de la classe  $k$  dans la métrique  $\mathbf{D}^{-1}$  (distance de Mahalanobis globale) est telle que :

$$\hat{D}^2(i, G_k) = \sum_{r=1}^{r_{\max}} \frac{(\hat{z}_{ir} - \bar{\hat{z}}_{kr})^2}{\hat{\alpha}_r} \quad [3.3 - 14]$$

$\hat{D}^2(i, G_k)$  est régularisée si  $r_{\max} = p''$  (où  $p''$  désigne le rang de la matrice transformée  $\hat{\mathbf{X}}$ ) quand  $p'' < \text{Min}(n, p)$  ou si  $r_{\max} < p''$ .

### - Axes principaux dans les groupes

Pour chaque classe  $I_k$ , les matrices de covariances d'ordre  $(r_{\max}, r_{\max})$  sont calculées séparément. On les exprimera ici à partir des coordonnées de l'analyse globale précédente.

Les nouvelles coordonnées de l'individu  $i$  sur l'axe principal  $s$  de l'analyse réalisée à l'intérieur de la classe  $I_k$  (il s'agit donc dans ce cas d'une simple analyse en composantes principales non normée) sont<sup>2</sup> :

$$w_{ski} = \mathbf{v}'_{sk}(\mathbf{z}_i - \bar{\mathbf{z}}_k)$$

où  $\mathbf{v}_{sk}$  est le  $s^{\text{ème}}$  vecteur propre normalisé de  $\mathbf{U}'\mathbf{D}_k\mathbf{U}$  correspondant à la valeur propre  $\beta_{sk}$  ( $\beta_{sk}$  est également valeur propre de  $\mathbf{D}_k$ ).

Avec ces coordonnées, on peut évidemment retrouver les distances usuelles, calculées cette fois dans chacune des  $q$  nouvelles bases (pour tout point  $i$  et tout point-moyen  $G_k$ ), lorsque le nombre  $s_{\max}(k)$  d'axes retenus à ce stade pour la classe  $k$ , vérifie :  $s_{\max}(k) = r_{\max}$ .

$$d^2(i, G_k) = \sum_{s=1}^{s_{\max}(k)} (w_{ski} - \bar{w}_{ks})^2 \quad [3.3 - 15]$$

avec :

$$\bar{w}_{ks} = \mathbf{v}'_{ks}(\bar{\mathbf{z}}_k - \bar{\mathbf{z}})$$

<sup>1</sup> Comme l'opération de centrage global, cette opération correspond à une projection  $\mathbf{P}$ . Si  $\mathbf{Y}$  désigne le tableau disjonctif complet d'ordre  $(n, q)$  décrivant la partition à prédire, l'opérateur projection s'écrit :  $\mathbf{P} = \mathbf{I} - \mathbf{Y}(\mathbf{Y}'\mathbf{Y})^{-1}\mathbf{Y}'$ . On peut parler dans ces conditions d'analyse interne ou conditionnelle : comme en analyse de la variance, on a éliminé la dispersion due aux classes en supposant que celles-ci avaient un effet additif.

<sup>2</sup> Cette formule de projection sur l'axe  $t$  est évidemment valable pour des points n'appartenant pas à la catégorie  $k$  (points supplémentaires ou illustratifs).

La distance de Mahalanobis locale (intervenant en analyse discriminante quadratique) peut s'écrire :

$$\mathcal{D}^2(i, G_k) = \sum_{s=1}^{s_{\max}(k)} \frac{(w_{ski} - \bar{w}_{ks})^2}{\beta_{sk}} \quad [3.3 - 16]$$

Une telle distance peut être "régularisée" à deux niveaux :

- une première fois si  $r_{\max} < p'$  ( $p'$  désigne le rang du tableau de donnée) ;
- de nouveau si  $s_{\max}(k) < r_{\max}$ .

On a noté que, si  $s_{\max}(k) = r_{\max} = p$ , les distances données par les formules [3.3 - 11], [3.3 - 12] et par les  $q$  formules [3.3 - 15] (il y a  $q$  bases orthonormées différentes donc  $q$  formules différentes) sont toutes égales.

#### - Exemple numérique d'application

L'exemple qui suit concerne les effets de la dimension des sous-espaces sur les pourcentages de bien-classés, à la fois dans les échantillons d'apprentissage et dans les échantillons-tests.

Le jeu de données utilisé est un tableau binaire clairsemé de dimensions (634, 83) contenant 4039 cases non-nulles<sup>1</sup>.

L'ensemble des 634 lignes (répondants) peut être réparti en  $q = 3$  classes d'âge. Le problème est de savoir dans quelle mesure ces classes âge peuvent être prédites à partir des réponses. Notre critère d'évaluation de la discrimination est le pourcentage de succès (bien classés), qui sera calculé systématiquement à la fois pour l'échantillon d'apprentissage et pour un échantillon-test qui comprend le tiers (211 individus) de l'échantillon global.

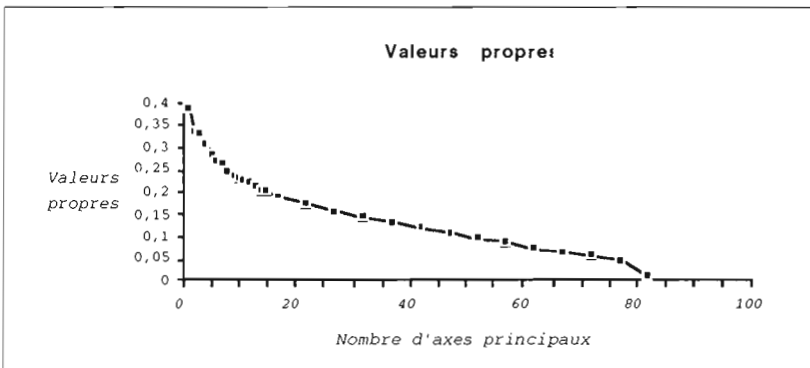


Figure 3.3 - 6  
Séquence des valeurs propres de la première analyse.

<sup>1</sup> Il s'agit pour cet exemple de 4039 occurrences de  $p = 83$  mots utilisés dans  $n = 634$  réponses à une question ouverte dans une enquête (cf. Lebart, 1992).



La première étape est un changement d'axes par analyse des correspondances. La séquence des valeurs propres, visible sur la figure précédente (cf. figure 3.3 - 6), est assez typique des tableaux clairsemés : la décroissance des valeurs propres est très lente, presque linéaire après l'axe 15. Les 15 premières valeurs propres correspondent à 37% de la trace. Chacun des axes restant correspond approximativement à 1% de la trace.

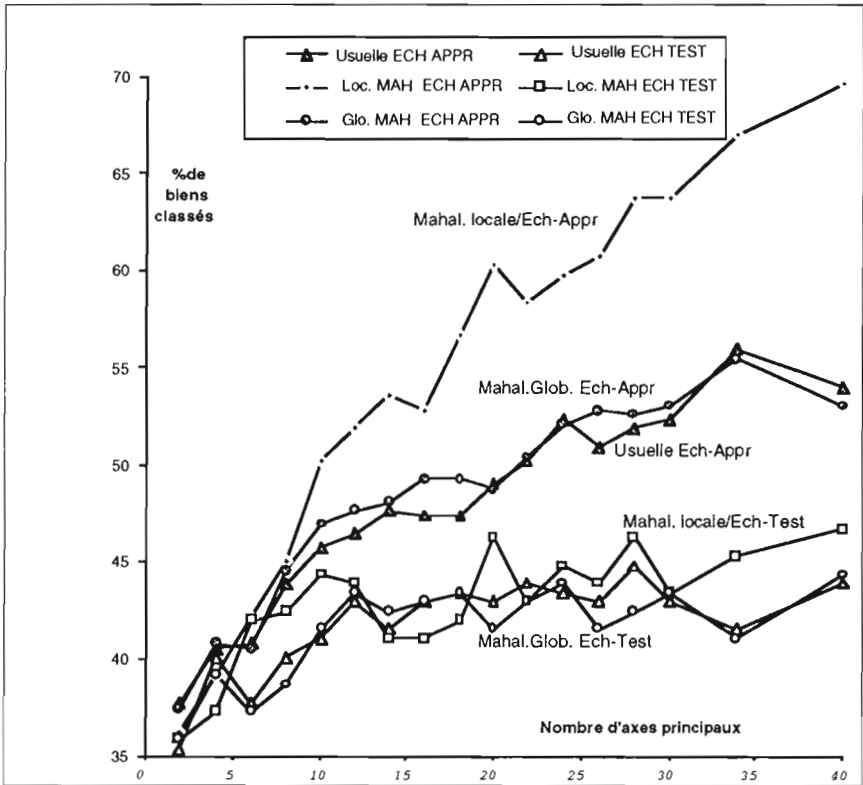


Figure 3.3 - 7

Trajectoires des pourcentages de bien classés en fonction du nombre d'axes principaux (axe des abscisses) selon trois distances et selon le type d'échantillon (test ou apprentissage)

La figure 3.3 - 7 montre les trajectoires des pourcentages de succès obtenus pour chacune des trois distances précédentes : *Distance euclidienne usuelle* (formule [3.3 - 12]), *distance de Mahalanobis globale* (formule [3.3 - 14]), *distance de Mahalanobis locale* (formule [3.3 - 16]).

On note que les taux correspondant aux échantillons d'apprentissage augmentent continûment avec le nombre d'axes alors que les taux correspondant aux échantillons-tests sont pratiquement stabilisés dès l'axe 15 (avec les notations ci-dessus, on peut choisir sans inconvénient  $s_{max}(k) = r_{max} = 15$  alors que  $p = 83$  et  $p' = 82$ ).

Parmi les trajectoires des échantillons d'apprentissage, le pourcentage de bien classés correspondant à la *distance de Mahalanobis locale* croît fortement et atteint un niveau de 70% de succès pour 40 axes. Une telle distance dépendant d'un nombre de paramètres plus important que les deux autres, s'adapte souplement aux données d'apprentissage<sup>1</sup>, sans fournir d'amélioration notable sur les échantillons-tests.

Pour cet exemple, la *distance de Mahalanobis globale* a des performances très voisines de la distance euclidienne usuelle. Les performances sont légèrement supérieures pour l'échantillon d'apprentissage.

Cet exemple met bien en évidence la puissance du filtrage réalisé par l'analyse factorielle préalable. La plupart des traits structuraux susceptibles de donner lieu à une prévision sont retenus dans l'espace à 15 dimensions des premiers axes.

### 3.3.7 Discrimination sur variables nominales

L'analyse factorielle discriminante que nous venons de présenter s'applique à un ensemble de  $n$  individus répartis en  $q$  classes définies *a priori* par la variable nominale  $y$  et décrits par  $p$  variables  $(x_1, x_2, \dots, x_p)$  continues. Lorsque les  $p$  variables explicatives sont nominales, le calcul des fonctions linéaires discriminantes ne peut plus être appliqué, en raison des singularités de la matrice  $X$ , mais la régularisation par axes principaux permettra de lever cette difficulté.

#### a – Analyse factorielle discriminante qualitative

Comme pour tout traitement de variables nominales, on procède au codage disjonctif complet des  $p$  variables explicatives. L'analyse factorielle discriminante qualitative consiste alors en une analyse factorielle discriminante classique sur les indicatrices des variables explicatives.

La matrice des variables explicatives  $X = [X_1, X_2, \dots, X_p]$  n'est pas inversible puisqu'il existe  $p$  relations linéaires entre les colonnes du tableau disjonctif complet. On peut alors, comme pour l'analyse de la variance, supprimer une modalité de chaque variable nominale ce qui ne modifie pas le sous-espace des variables explicatives  $V_X$ . Ceci ne suffit pas à assurer que la matrice réduite est bien conditionnée.

La régularisation par axes principaux revient dans ce cas à réaliser une analyse discriminante classique sur les facteurs de l'analyse des correspondances multiples<sup>2</sup>.

<sup>1</sup> Plus le nombre de paramètres augmente, plus l'apprentissage se rapproche de ce que l'on appelle en intelligence artificielle *l'apprentissage par coeur*, c'est-à-dire une adaptation trompeuse à une situation donnée, sans possibilité de généralisation.

<sup>2</sup> Enchaînement connu en particulier sous le nom de méthode DISQUAL (Saporta, 1977).

On procède alors en effectuant :

- une analyse des correspondances sur le tableau disjonctif complet ; les  $p$  variables nominales sont donc remplacées par  $h$  variables continues qui sont les  $h$  facteurs de l'analyse des correspondances multiples.
- une analyse factorielle discriminante sur les  $h$  variables continues dont les valeurs sont les coordonnées sur les axes factoriels de l'analyse des correspondances multiples.

Compte tenu du nombre généralement important de facteurs de l'analyse des correspondances multiples, on retiendra les facteurs les plus discriminants et qui ne figurent pas toujours parmi les premiers<sup>1</sup>.

### b – Analyse discriminante barycentrique

L'analyse discriminante barycentrique revient simplement à faire l'analyse des correspondances du tableau croisant la variable à expliquer  $y$  avec les variables explicatives ( $x_1, x_2, \dots, x_p$ ) (empilement de tables de contingences) : les lignes sont constituées par les modalités de  $y$  et les colonnes par la juxtaposition des modalités de ( $x_1, x_2, \dots, x_p$ ).

Il s'agit en fait d'une bande du tableau de Burt qui permet de décrire les liaisons existant entre la variable à expliquer et l'ensemble des variables explicatives (cf. §1.4.7.b; Saporta, 1975 a ; Leclerc, 1976).

En plaçant en éléments supplémentaires de nouveaux individus caractérisés par les variables explicatives, on réalise une réaffectation similaire à celle l'analyse discriminante (cf. Nakache *et al.*, 1977).

Dans le cas où les variables explicatives sont indépendantes deux à deux, l'analyse discriminante barycentrique est équivalente à l'analyse factorielle discriminante qualitative (puisque l'analyse d'une bande du tableau de Burt est alors équivalente à l'analyse du tableau complet). Dans le cas général, elle est, en théorie, moins performante puisque, comme nous l'avons vu dans §1.4.7.b, elle ne tient pas compte des liaisons entre les variables explicatives. Elle est cependant largement utilisée en raison de sa simplicité et sa robustesse (cf. Carlier, in : Celeux et Nakache, 1994).

### c – Note sur le "scoring"

Fréquemment utilisée par les organismes bancaires cherchant à prévoir la défaillance éventuelle d'un client (individu ou entreprise), la méthode dite de "scoring" permet une mise en forme simple des résultats d'une analyse discriminante généralement à deux groupes. Elle n'est pas à proprement

---

<sup>1</sup> Que ce soit pour l'analyse factorielle discriminante qualitative et, nous allons le voir, pour l'analyse discriminante barycentrique, il est conseillé de procéder au préalable à une *première sélection des variables nominales explicatives* en croisant par exemple chacune d'entre elles avec la partition à expliquer  $y$ , en calculant les  $\chi^2$  correspondants, et gardant celles qui correspondent aux  $\chi^2$  les plus significatifs.

parler une méthode de discrimination sur variables nominales ; mais elle utilise les résultats d'analyses discriminantes sur variables nominales ou continues pour construire une *fonction de score*<sup>1</sup>. On dispose ainsi d'un instrument de décision accessible pour affecter un individu dans un groupe.

Dans le cas de deux groupes, on obtient une seule fonction discriminante : la combinaison linéaire des variables qui sépare au mieux les deux groupes d'individus. Un individu est affecté à l'un des groupes si la fonction prend pour lui une valeur supérieure à un certain seuil.

Cette fonction discriminante est ensuite transformée en un système équivalent de coefficients attribués aux modalités des variables nominales ou aux éventuelles variables continues (en général après une sélection sévère). Cette transformation fournit la fonction score dont les coefficients constituent des notes attachées aux modalités ou aux variables.

Pour chaque individu, on calcule le score<sup>2</sup> c'est-à-dire la somme des notes associées aux prédicteurs. On affectera alors cet individu à un groupe si son score est supérieur à un certain seuil. L'introduction d'une tolérance d'erreur de classement permet en fait de définir trois zones de décisions sur la fonction score : la zone des scores élevés, celle des scores faibles et une zone d'indécision pour laquelle un individu n'est pas automatiquement classé.

### 3.3.8 Discrimination et réseaux de neurones

Ce paragraphe ne constitue qu'une brève note bibliographique destinée à orienter le lecteur statisticien désireux d'aborder les techniques neuronales de discrimination.

Développées au milieu des années quatre-vingt, les méthodes neuronales (ou réseaux neuronaux ou encore réseaux neuro-mimétiques) ont renouvelé et stimulé la discipline connue sous le nom de *reconnaissance de formes* qui recouvre beaucoup d'applications industrielles (notamment des applications en temps réel) des méthodes de discrimination.

Fondées au départ sur des analogies biologiques et sur un effort de modélisation des mécanismes de perception visuelle et auditive, ces méthodes ont acquis depuis une certaine autonomie. Les relations avec la statistique ont été frileuses en raison de différences d'approches et de vocabulaires<sup>3</sup>. Mais des ponts ont été jetés et les années récentes ont vu la

---

<sup>1</sup> Cf. dans le cas d'analyses appliquées à la détection de défaillances d'entreprises (à partir de sélection de variables continues) : Bardos (1984, 1989).

<sup>2</sup> Les enchaînements de calculs de l'analyse discriminante qualitative, la fonction score ainsi que l'analyse barycentrique (construction d'une bande du tableau de Burt) sont prévus dans le logiciel SPAD.N.

<sup>3</sup> Ce sont des informaticiens en milieu industriel qui sont à l'origine de ces méthodes.

parution d'une série d'articles de revue ou de synthèse<sup>1</sup> qui ont prouvé la complémentarité des points de vue et l'enrichissement mutuel à attendre des contacts et échanges entre statisticiens et neuromiméticiens.

Schématiquement, disons que les statisticiens peuvent compléter la panoplie des modèles qui leur sont familiers avec les modèles essentiellement *non-linéaires* et à seuils qui sont attachés aux réseaux de neurones. La structure de ces réseaux permet d'autre part des calculs parallèles indispensables pour une implémentation matérielle directe de ces méthodes et des utilisations en temps réel, domaine peu abordé par les statisticiens. Inversement, l'essentiel de ce qui concerne l'inférence ou la validation des démarches et des résultats est à mettre au crédit des approches statistiques. Ces aspects sont reconnus comme indispensables dès qu'il s'agit de comparer des modèles, d'évaluer des risques, de calculer des taux d'erreurs, préoccupations caractéristiques d'une discipline arrivée à maturité.

On évoquera seulement dans cette note bibliographique le modèle neuronal le plus répandu dans le cadre de la discrimination qui est le *perceptron multi-couche*, puis on dira quelques mots des méthodes non-supervisées.

### a – Schéma et modèle du perceptron multi-couches

Le contexte est le même que celui qui a été défini au début de cette section. On dispose d'une variable qualitative  $y$  à  $q$  modalités (ou catégories) que l'on doit prédire à partir de  $p$  variables ( $x_1, x_2, \dots, x_p$ ) prédictrices. On dispose par ailleurs de  $n$  individus ou observations (échantillon d'apprentissage) décrits par les  $p$  variables ( $x_1, x_2, \dots, x_p$ ) et pour lesquels on connaît la classe d'affectation notée ici  $y_k$  ( $k \leq q$ ).

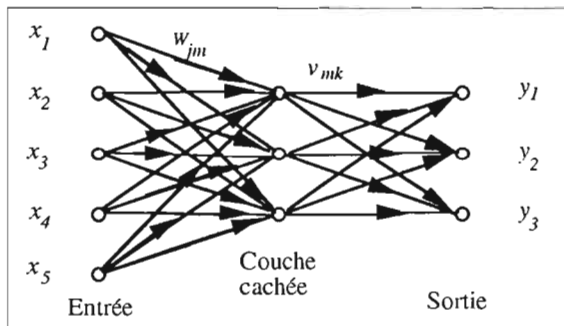


Figure 3.3 - 8  
Perceptron à une couche cachée

La figure 3.3 - 8 se commente de la façon suivante en utilisant le vocabulaire et les concepts de l'approche neuronale : la couche d'entrée est formée de

<sup>1</sup> Citons en particulier les articles de synthèse de Ripley (1993, 1994) et de Cheng et Titterington (1994).

$p = 5$  entrées, auxquelles seront appliquées des coefficients appelés les *poids synaptiques*  $w_{jm}$ . La *couche cachée* comprend  $c = 3$  neurones qui seront chacun *activés* par une intégration (en général fonction monotone de la somme) des  $p$  signaux en provenance de la couche d'entrée. La même opération a lieu pour les  $q = 3$  éléments de la couche de sortie mettant en jeu des poids synaptiques  $v_{mk}$ .

En termes de modèle analytique, on écrira :

$$y_k = \Phi_0 \left\{ a_k + \sum_{m=1}^c v_{mk} \Phi \left( a_m + \sum_{j=1}^p w_{jm} x_j \right) \right\} \quad [3.3 - 17]$$

Dans cette formule, la fonction  $\Phi$  est dans la plupart des applications la fonction logistique qui sera abordée à la section 3.4. Elle s'écrit :

$$\Phi(z) = \frac{\exp\{z\}}{1 + \exp\{z\}}$$

La fonction  $\Phi_0$  peut être selon les cas linéaire, logistique, ou à seuil (par exemple :  $\Phi_0(z) = 0$  si  $z \leq 0$  et  $\Phi_0(z) = 1$  si  $z > 0$ ).

On voit que la figure 3.3 - 8 est utile pour visualiser l'enchaînement de fonctions correspondant aux étapes du traitement. La lecture de droite à gauche de la figure correspond bien sûr à une lecture de gauche à droite de la formule [3.3 - 17]. Il y a  $\{c(p+1) + q(c+1)\}$  paramètres à estimer.

L'équation [3.3 - 17] correspond à une observation  $(i)$ . On a en réalité  $n$  équations de ce type, faisant chacune intervenir  $q$  valeurs  $y_k^{(i)}$  (valeurs 0 ou 1 s'il s'agit d'appartenance à une classe d'une partition en  $q$  classes) et  $p$  valeurs  $x_j^{(i)}$ .

L'estimation des paramètres se fait en minimisant une fonction de perte, qui peut simplement être la somme des carrés des écarts entre les valeurs calculées  $\tilde{y}_k^{(i)}$  et les valeurs observées  $y_k^{(i)}$  dans l'échantillon d'apprentissage<sup>1</sup>.

Remarquons que pour une sortie binaire (deux classes possibles pour  $y$  qui peut alors être un scalaire prenant les valeurs 0 ou 1) et un perceptron sans couche cachée, on se trouve dans le cadre du modèle de la régression logistique évoqué en section 3.4.4.

La formule [3.3 - 17] s'écrit alors :

$$y = \Phi_0 \left\{ \Phi \left( a_m + \sum_{j=1}^p w_{jm} x_j \right) \right\} \quad [3.3 - 18]$$

<sup>1</sup> L'estimation numérique se fait par une méthode de descente de gradient dite de *back-propagation*. (cf. Werbos, 1974, 1990; Rumelhart et al., 1986). Pour un programme de calcul, cf. Proriot (1991), ou la procédure NEURO du logiciel SPAD.N.

Ici, la fonction  $\Phi_0$  peut être une fonction à seuil, qui convertit la probabilité donnée par le modèle logistique proprement dit (à l'intérieur des accolades) en l'une des deux valeurs 0 ou 1.

Si l'on réduit les deux fonctions  $\Phi_0$  et  $\Phi$  à la fonction identique  $\Phi(x) = x$ , on retrouve la régression multiple (cf. section 3.2) et l'analyse discriminante à deux groupes (cf. paragraphe 3.3.3) qui en sont des cas particuliers.

Cet exemple très simple du perceptron multi-couches montre donc que les généralisations les plus évidentes par rapport aux modèles explicatifs usuels de la statistique concernent la présence éventuelle des fonctions  $\Phi_0$  et  $\Phi$  et l'existence d'une ou plusieurs couches cachées qui autorisent des interventions non-linéaires des paramètres<sup>1</sup>.

## b – Modèles non-supervisés ou auto-organisés

Alors que les modèles supervisés (pour lesquels on dispose d'un échantillon d'apprentissage permettant d'estimer les paramètres) correspondent tout à fait à la démarche de la régression et de l'analyse discriminante, les modèles non-supervisés ou auto-organisés sont le pendant des méthodes purement exploratoires.

Reprenons l'exemple du perceptron multicouche, pour lequel nous supposerons les fonctions  $\Phi_0$  et  $\Phi$  linéaires ou (sans perte de généralité dans ce cas) égales à la fonction identique. Nous supposerons de plus que les variables sont des variables numériques centrées, et que les termes constants sont nuls.

La formule [3.3 - 17] s'écrit :

$$y_k = \left\{ \sum_{m=1}^c v_{mk} \left( \sum_{j=1}^p w_{jm} x_j \right) \right\} = \sum_{j=1}^p \left( \sum_{m=1}^c v_{mk} w_{jm} \right) x_j \quad [3.3 - 19]$$

que l'on peut écrire sous la forme :

$$y = \mathbf{VWx}, \text{ soit } y^{(i)} = \mathbf{VWx}^{(i)} \text{ pour chaque observation } i.$$

On peut poser  $\mathbf{A} = \mathbf{VW}$ . La matrice  $\mathbf{A}$  est d'ordre  $(q, p)$ , si la taille  $c$  de la couche cachée n'introduit pas de restriction sur le rang de  $\mathbf{A}$ , qui est au plus le plus petit des trois nombres  $q, c, p$ .

En l'absence de contrainte sur  $\mathbf{A}$ , on est dans le cadre de la régression multiple simultanée comportant plusieurs variables endogènes, qui revient à faire  $q$  régressions multiples (cf. § 3.6.2 b).

<sup>1</sup> Notons que dans un modèle général comme celui de la formule [3.3 - 17], il n'est pas nécessaire de retenir toutes les flèches entre deux couches consécutives (certain poids synaptiques peuvent être nuls *a priori*, d'autres peuvent avoir une valeur fixe, et ainsi réduire le nombre de paramètres à estimer).

La solution s'obtient ici en rendant minimale la somme sur les  $n$  observations :

$$S = \sum_i \left( \mathbf{y}^{(i)} - \mathbf{Ax}^{(i)} \right)' \left( \mathbf{y}^{(i)} - \mathbf{Ax}^{(i)} \right)$$

Dans les modèles non-supervisés dits d'auto-association, on ne connaît pas  $\mathbf{y}$ , (il n'y a pas de "professeur") et on utilise l'artifice qui consiste à remplacer  $\mathbf{y}$  par  $\mathbf{x}$  (cf. Baldi et Hornik, 1989).

Ceci semble une trivialité, et est effectivement une trivialité si la couche cachée possède autant d'éléments que  $\mathbf{x}$  ( $c = p$ ) et s'il n'y a pas de contraintes sur  $\mathbf{A}$  (auquel cas on a la solution  $\mathbf{A} = \mathbf{I}$ ).

Mais si la couche cachée est notablement plus réduite que les couches d'entrée et de sortie, ( $c \ll p$ ), elle forme un étranglement et le réseau réalise une *compression* du signal d'entrée.

On veillera à rendre minimale la quantité  $S_1$  :

$$S_1 = \sum_i \left( \mathbf{x}^{(i)} - \mathbf{VWx}^{(i)} \right)' \left( \mathbf{x}^{(i)} - \mathbf{VWx}^{(i)} \right)$$

On s'efforce donc de réduire le plus possible la déformation moyenne de  $\mathbf{x}$  après intervention du réseau, qui n'est autre ici qu'une projection sur un sous-espace de dimension  $c$  inférieure à  $p$ . La solution est fournie par l'analyse en composantes principales du tableau  $\mathbf{X}$  (qui est aussi une décomposition aux valeurs singulières, puisque nous avons supposé les variables centrées) dont les  $n$  lignes sont les vecteurs  $\mathbf{x}_i'$ .

Ainsi, par exemple, avec un seul neurone dans la couche cachée, la matrice  $\mathbf{VW}$  est de rang 1, ce qui conduira au premier axe de l'analyse en composantes principales de  $\mathbf{X}$ . Une démonstration complète, incluant le cas supervisé (régressions multiples simultanées avec contrainte de rang pour le tableau de coefficients) se trouve dans Baldi et Hornik (1989, *op. cit.*).

L'*auto-organisation*, notion étudiée et formalisée par Kohonen (1989), qui est un des pionniers de l'approche neuronale, est donc rendue possible par la structure interne du réseau.

D'autres travaux sont relatifs aux algorithmes à lecture directe, comme l'algorithme de diagonalisation par approximation stochastique proposé par Benzécri (1969 b), antérieurement aux approches neuronales<sup>1</sup>.

Ces algorithmes peuvent en effet être interprétés en terme d'apprentissage et d'auto-organisation. Un algorithme identique à une normalisation près a été proposé indépendamment par Oja et Karhunen (1981), puis amélioré par la suite par ces auteurs et d'autres neuromiméticiens. Ce domaine, qui a des applications potentielles importantes en compression d'image, a depuis été

---

<sup>1</sup> On trouvera une étude plus numérique de la convergence de l'algorithme dans Lebart (1974), et le programme correspondant dans Lebart *et al.* (1977).



très développé. Sur les liens entre réseaux neuronaux et analyse en composantes principales, cf. Oja (1982), Bourlard et Kamp (1988), Sirat (1991), Oja (1992).

Enfin, une autre approche non-supervisée, plus proche des méthodes de classification, est celle des cartes auto-organisées (*self organizing maps*) de Kohonen (Kohonen, 1989; Cottrell et Fort, 1987). L'algorithme est assez similaire à celui de la méthode d'agrégation autour de centres mobiles (k-means) (démarrage aléatoire, affectations aux centres de distances minimales, obtention de minima locaux) mais conduit à une représentation plane (cf. Ritter *et al.*, 1992).

### c – Statistique et méthodes neuronales

On complétera cet aperçu par un résumé de l'intervention de Tibshirani lors d'une discussion sur la synthèse de Cheng et Titterington (1994, *op.cit.*). Cette intervention commence par une remarque générale sur la statistique et les réseaux de neurones :

"Les statisticiens ont tendance à travailler avec des modèles plus interprétables car, pour eux, mesurer l'effet des variables est plus important que la prédiction".

Tibshirani répond ensuite à deux questions :

#### - *Que peut apprendre un statisticien d'un neuro-miméticien ?*

- 1 "On devrait moins se soucier de l'optimalité statistique que de trouver des méthodes qui fonctionnent, spécialement sur les grands ensembles de données.
- 2 On devrait plus attaquer les problèmes réels auxquels se consacrent les neuro-miméticiens : reconnaissance de l'écriture et de la parole, prédiction des structures de l'ADN. Comme le dit John Tukey : il vaut mieux avoir une solution approchée d'un problème réel que la solution exacte d'un problème trop simplifié.
- 3 Les modèles à très nombreux paramètres peuvent être utiles pour la prédiction, spécialement pour les grands tableaux de données et les données bruitées.
- 4 Modéliser des combinaisons linéaires des variables d'entrées est très utile, car cela prend en compte des traits structurels et réduit la dimension.
- 5 Des algorithmes itératifs comme la descente de gradient (avec taux d'erreurs) peuvent éviter des ajustements trop complaisants.
- 6 Nous (statisticiens) devrions mieux nous vendre..."

#### - *Que peut apprendre un neuromiméticien d'un statisticien ?*

- 1 "Il devrait plus s'intéresser à l'optimalité statistique, ou au moins, aux propriétés statistiques des méthodes.

- 2 Il devrait faire plus d'efforts pour comparer ses méthodes à des méthodes statistiques plus simples. Il serait alors surpris de voir que la régression fait souvent aussi bien qu'un perceptron multi-couches. Il ne devrait jamais utiliser un modèle compliqué alors qu'un modèle simple suffit."

Ces remarques n'épargnent pas les statisticiens, qui ont devant eux une profusion d'idées nouvelles et un vaste chantier ouvert. Ceux d'entre eux qui se consacrent à l'analyse exploratoire des grands tableaux se sentent cependant moins concernés par les deux premières critiques de Tibshirani.

#### *Autres références*

Outre les trois articles de synthèse précités, on mentionnera, toujours pour un lectorat de statisticien : l'ouvrage de base de Hertz *et al.* (1991), l'article plus théorique de Amari (1990), sur les fondements mathématiques des méthodes. Mentionnons également l'article de Hornik (1994) décrivant, à l'intention des statisticiens, le perceptron multicouche et les algorithmes d'analyses en composantes principales par apprentissage comme deux intersections importantes entre les deux disciplines. En Français, on consultera les ouvrages généralistes de Bourret *et al.* (1991) et de Milgram (1993). Pour des exposés faisant le lien avec l'approche "analyse des données", Gallinari *et al.* (1988), Lelu (1991), Chabanon et Dubuisson (1991).

# Modèles log-linéaires

Les modèles log-linéaires permettent d'étudier et de modéliser les liaisons entre plusieurs variables nominales en tenant compte de leurs éventuelles interactions.

On peut considérer l'analyse des tables de contingences multiples par ces modèles comme une analyse descriptive car aucune variable ne joue le rôle privilégié de variable à prévoir. Mais ces modèles s'apparentent aussi, par leur démarche, à l'analyse de la variance (sélection de modèles sur la base de tests statistiques), ce qui justifie leur présentation dans cette partie. Ils nous amènent d'ailleurs à la régression logistique qui peut être considérée comme l'analogue de la régression multiple sur variables nominales. La régression logistique se propose en effet de prévoir une variable dichotomique à l'aide d'une ou de plusieurs variables (de nature quelconque) en prenant en compte l'effet propre de chaque variable et l'effet éventuel des interactions<sup>1</sup>.

### 3.4.1 Formulation du problème et principes de base

Présentons le problème à partir d'un exemple médical. Considérons un échantillon d'individus ayant été irradiés accidentellement. Ces individus sont caractérisés par un état (être *décédés* ou *non* à la suite de leucémie : variable nominale à 2 modalités), par la dose de radiations reçue mesurée en *Rad* (variable continue ordonnée en 6 modalités) et par l'âge au moment des accidents (variable continue regroupée en 5 modalités).

Ces données se présentent sous forme d'un tableau de contingence  $K$  croisant ces trois variables de terme général  $k_{ijl}$ .

On s'intéresse aux relations existant entre ces trois variables : sont-elles indépendantes ou non dans leur ensemble ou une variable est-elle indépendante conditionnellement à une ou aux deux autres ? Autrement dit, on cherche à connaître la structure des liaisons entre ces données en tenant compte des interactions entre les 3 variables.

---

<sup>1</sup> Les modèles log-linéaires et logistiques donnent lieu à des publications nombreuses. Après les premiers travaux de Birch (1963) et Goodman (1970), il faut mentionner les ouvrages de base de Haberman (1974), Bishop, Fienberg, Holland (1975), Fienberg (1980). Plus récemment, Dobson (1983), Agresti (1990), Christensen (1990) rédigent des synthèses enrichies de contributions personnelles. Goodman (1986, 1991) fait des rapprochements avec certains aspects de l'analyse des correspondances. Anderson (1982) réalise une revue très complète du modèle logistique. L'ouvrage collectif édité par Celeux et Nakache (1994) présente les contributions des modèles log-linéaires et logistiques à la discrimination.

D'une manière générale,  $p$  variables nominales  $x_1, x_2, \dots, x_p$  ayant respectivement  $m_1, m_2, \dots, m_p$  modalités, constituent un tableau de contingence multidimensionnel à  $p$  entrées comprenant  $m_1 \times m_2 \times \dots \times m_p$  cases. Le terme général  $k_{ij\dots p}$  de cet hypercube de contingence indique le nombre d'individus ayant répondu simultanément aux modalités  $i, j, \dots, p$  de  $x_1, x_2, \dots, x_p$  avec  $1 < i < m_1, 1 < j < m_2, \dots, 1 < p < m_p$ .

L'effectif total d'individus observés est noté  $k$  avec :

$$k = \sum_{i,j,\dots,p} k_{ij\dots p}$$

Les hypothèses que nous formulons sur les liaisons entre ces  $p$  variables nous amènent à construire des tableaux de fréquences théoriques espérées  $T$  de terme général  $t_{ij\dots p}$ . La confrontation des fréquences observées  $k_{ij\dots p}$  et des fréquences théoriques  $t_{ij\dots p}$  va permettre de tester ces hypothèses.

On construira par conséquent autant de tableaux  $T$  (et donc de modèles log-linéaires) qu'il y a d'hypothèses à tester.

Dans le cas d'un tableau de contingence à deux dimensions, on construit, sous l'hypothèse d'indépendance entre les deux variables, le tableau  $T$  tel que  $t_{ij} = t_{i.} t_{.j}$ . Le test du  $\chi^2$  permet de rejeter ou non cette hypothèse en confrontant le tableau théorique  $T$  au tableau des fréquences observées  $K$ .

Ainsi les modèles log-linéaires peuvent être considérés comme une généralisation du test du  $\chi^2$  à un ensemble de  $p$  variables nominales ( $p > 2$ ), la difficulté résidant alors dans le choix des modèles, c'est-à-dire des hypothèses concernant les liaisons entre les variables.

### 3.4.2 Ajustement d'un modèle log-linéaire

On suppose que la fréquence observée  $k_{ij\dots p}$  est la réalisation d'une variable aléatoire  $x_{ij\dots p}$  d'espérance mathématique inconnue  $t_{ij\dots p}$ .

$$E(x_{ij\dots p}) = t_{ij\dots p}$$

Nous envisagerons successivement le cas du tableau de contingence à deux dimensions et celui à  $p$  entrées. Les notations étant lourdes dans le cas général, nous nous bornerons à  $p = 3$  pour simplifier l'exposé.

#### a – Tableau de contingence à deux entrées

Intéressons-nous d'abord à la relation entre deux variables nominales, la *risque de décès* et la *dose de radiation reçue*, par exemple. Dans ce cas, deux hypothèses peuvent être formulées : y a-t-il indépendance ou non entre les deux variables ?

En supposant  $t_{ij}$  non nul, le modèle log-linéaire le plus complet décompose le logarithme népérien de l'espérance  $t_{ij}$  sous la forme :

$$\log(t_{ij}) = \alpha_0 + \alpha_1(i) + \alpha_2(j) + \alpha_{12}(ij)$$

Par analogie avec l'analyse de la variance,  $\log(t_{ij})$  se décompose en une somme de coefficients  $\alpha$  décrivant plusieurs effets :

- $\alpha_0$ , l'effet global;
- $\alpha_1(i)$ , l'effet dû à la variable  $x_1$ ,
- $\alpha_2(j)$ , l'effet dû à la variable  $x_2$ ,
- $\alpha_{12}(ij)$ , l'effet dû à l'interaction entre les variables  $x_1$  et  $x_2$ .

Afin d'avoir une solution unique, on impose les contraintes suivantes :

$$\sum_i \alpha_1(i) = \sum_j \alpha_2(j) = \sum_i \alpha_{12}(ij) = \sum_j \alpha_{12}(ij) = 0$$

Sous l'hypothèse d'indépendance des deux variables, la fréquence espérée s'exprime par  $t_{ij} = t_i \cdot t_j$ . Dans ce cas, tous les coefficients d'interaction  $\alpha_{12}(ij)$  sont nuls. Le modèle log-linéaire correspondant à cette hypothèse s'écrit :

$$\log(t_{ij}) = \alpha_0 + \alpha_1(i) + \alpha_2(j)$$

La nullité des interactions traduit l'hypothèse d'indépendance entre les deux variables. A partir des coefficients  $\alpha_0$ ,  $\alpha_1(i)$  et  $\alpha_2(j)$ , on calcule le tableau des fréquences théoriques espérées noté **T**.

### **b – Tableau de contingence à $p$ entrées**

On généralise ces modèles au cas de plus de deux variables. Pour trois variables par exemple, le modèle qui prend en compte toutes les liaisons entre les variables est le suivant :

$$\begin{aligned} \log(t_{ijl}) = & \alpha_0 + \alpha_1(i) + \alpha_2(j) + \alpha_3(l) \\ & + \alpha_{12}(ij) + \alpha_{13}(il) + \alpha_{23}(jl) + \alpha_{123}(ijl) \end{aligned} \quad [3.4 - 1]$$

Ce modèle est appelé *modèle saturé*. Il contient tous les effets et toutes les interactions qu'il est possible de définir avec les variables disponibles.

Les coefficients  $\alpha_0, \alpha_1(i), \dots, \alpha_{123}(ijl)$  traduisent des effets différents :

- $\alpha_0$ , l'effet global;
- $\alpha_1(i), \alpha_2(j), \alpha_3(l)$ , les effets principaux;
- $\alpha_{12}(ij), \alpha_{13}(ik), \alpha_{23}(jl)$ , les effets dus aux interactions deux à deux des variables;
- $\alpha_{123}(ijl)$ , l'effet dû à l'interaction à trois variables;

On impose la nullité de la somme des coefficients du modèle faisant intervenir une modalité d'une variable sur l'ensemble des modalités de cette même variable.

Par exemple pour la variable  $x_1$  et pour tout  $1 < i < m_1$ , on a :

$$\sum_i a_1(i) = \sum_i a_{12}(ij) = \sum_i a_{13}(il) = \sum_i a_{123}(ijl) = 0$$

et il en est de même pour les autres variables.

Le modèle [3.4 - 1], comme tous les modèles saturés, permet de reconstituer exactement le tableau de fréquence  $\mathbf{K}$ . Celui-ci présentant souvent un trop grand nombre de coefficients, on va rechercher un ou des modèles ayant moins de coefficients mais devant reconstituer le mieux possible le tableau  $\mathbf{K}$  (principe de parcimonie). Ceci est réalisé en annulant certains termes du modèle saturé.

Si on arrive à une reconstitution correcte du tableau  $\mathbf{K}$ , l'hypothèse de nullité des coefficients supprimés ne peut pas être rejetée. Ces modèles non saturés mettent alors en évidence les liaisons les plus significatives entre les variables.

Dans le cas de deux variables, l'hypothèse de nullité du terme d'interaction s'interprète en terme d'indépendance. Si cette hypothèse est rejetée, on incriminera une dépendance entre les deux variables. Lorsque l'on s'intéresse à plus de deux variables, l'interprétation est plus complexe :

- pour exprimer l'*indépendance mutuelle* entre toutes les variables  $x_1, x_2, x_3$ , on annule tous les termes d'interactions. Cela nous conduit au modèle :

$$\log(t_{ijl}) = \alpha_0 + \alpha_1(i) + \alpha_2(j) + \alpha_3(l)$$

- pour exprimer l'*indépendance conditionnelle* de deux variables  $x_1$  et  $x_2$  par rapport à  $x_3$ , on annule tous les termes d'interaction contenant les indices relatifs aux variables  $x_1$  et  $x_2$  c'est-à-dire :

$$\alpha_{12}(ij) = \alpha_{123}(ijl) = 0$$

on en déduit le modèle suivant :

$$\log(t_{ijl}) = \alpha_0 + \alpha_1(i) + \alpha_2(j) + \alpha_3(l) + \alpha_{13}(il) + \alpha_{23}(jl)$$

Chaque modèle log-linéaire met ainsi en évidence une liaison particulière entre les variables : la dépendance ou l'indépendance mutuelle des variables dans leur ensemble ou l'indépendance de certaines variables conditionnellement à une ou plusieurs autres.

Pour des modèles à plus de trois variables, on trouvera des compléments sur les interactions, dans par exemple, Agresti (1990).

### c – modèles hiérarchiques

Un modèle log-linéaire est dit hiérarchique si la condition suivante est vérifiée : quand un coefficient d'interaction est présent dans le modèle, les coefficients des variables mises en jeu et toutes les interactions d'ordre inférieur sont aussi dans le modèle.

Par exemple, si dans un modèle à 5 variables on trouve l'interaction  $x_{135}$ , alors le modèle, pour être hiérarchique, doit contenir au moins  $x_1$ ,  $x_3$  et  $x_5$  ainsi que les interactions d'ordre inférieur  $x_{13}$ ,  $x_{15}$  et  $x_{35}$ .

Parmi les modèles log-linéaires possibles dans le cas d'un tableau de contingence à deux variables, certains modèles sont hiérarchiques :

- $\log(t_{ij}) = \alpha_0 + \alpha_1(i) + \alpha_2(j) + \alpha_{12}(ij)$
- $\log(t_{ij}) = \alpha_0 + \alpha_1(i) + \alpha_2(j)$

et d'autres ne le sont pas :

- $\log(t_{ij}) = \alpha_0 + \alpha_1(i) + \alpha_{12}(ij)$ ;
- $\log(t_{ij}) = \alpha_0 + \alpha_2(j) + \alpha_{12}(ij)$ ;
- $\log(t_{ij}) = \alpha_0 + \alpha_{12}(ij)$

Traditionnellement et pour des raisons de simplicité d'interprétation, on se limite aux modèles hiérarchiques.

### 3.4.3 Estimation et tests d'ajustement du modèle

On se donne un modèle traduisant une hypothèse exprimée par la nullité de certains coefficients  $\alpha$ . On cherche ainsi à estimer les fréquences théoriques pour construire puis confronter le tableau  $\hat{T}$  des estimations au tableau  $K$  des fréquences observées. Cette confrontation est réalisée par des tests d'ajustement. Ils permettent de rejeter ou non l'hypothèse sur les liaisons exprimée par le modèle.

#### a – Estimation des paramètres

Les fréquences théoriques espérées  $t_{ijl}$  sont en général estimées par la méthode du maximum de vraisemblance. Elle consiste à rechercher les paramètres qui maximisent la fonction de vraisemblance  $\mathcal{L}(k_{ijl}, t_{ijl})$ .

Pour cela, on suppose que les variables aléatoires  $x_{ijl}$  suivent soit une loi de Poisson, soit une loi multinomiale<sup>1</sup>.

On montre alors (cf. par exemple Haberman, 1974) que maximiser  $\mathcal{L}(k_{ijl}, t_{ijl})$  revient à maximiser :

$$\sum_{i,j,l} k_{ijl} \log(t_{ijl})$$

<sup>1</sup> Ce sont des hypothèses assez naturelles dans le cas des tables de contingence multidimensionnelles. Brièvement dit, la loi de Poisson correspond au cas où l'effectif total  $k$  n'est pas fixé ou borné a priori.

On calcule les estimations  $\hat{t}_{ijl}$  des fréquences espérées  $t_{ijl}$  données par le modèle. On peut utiliser la méthode de régression pondérée de Grizzle et *al.* (1969) ou celle des algorithmes itératifs (méthode de Newton-Raphson ou méthode des moindres carrés itératifs) qui est la méthode la plus répandue, utilisée pour tous les modèles linéaires généralisés, dont les modèles log-linéaires sont des cas particuliers<sup>1</sup>.

### b – Tests d'ajustement

Pour comparer le tableau des fréquences estimées  $\hat{T}$  avec le tableau des fréquences observées  $K$ , deux tests (voisins) sont généralement utilisés :

- le test du  $\chi^2$  de Karl Pearson :

$$\chi^2 = \sum_{i,j,l} \frac{(k_{ijl} - \hat{t}_{ijl})^2}{\hat{t}_{ijl}}$$

- le test du rapport de vraisemblance<sup>2</sup> :

$$G^2 = -2 \sum_{i,j,l} k_{ijl} \log \frac{\hat{t}_{ijl}}{k_{ijl}}$$

Les statistiques  $\chi^2$  et  $G^2$  suivent une distribution du  $\chi^2$  à  $m$  degrés de liberté où  $m$  est le nombre de cases du tableau auquel on soustrait le nombre de coefficients estimés. Pour l'une et l'autre de ces statistiques, les valeurs augmentent avec le nombre de variables introduites dans le modèle.

Plus ces statistiques sont voisines de zéro, meilleur est l'ajustement. Elles sont nulles pour le modèle saturé. On recherche le modèle le plus simple (peu de paramètres) et qui reste acceptable (bon ajustement).

### c – Choix du modèle

Le choix du modèle log-linéaire est d'autant plus difficile que le nombre de variables est élevé. La méthode dite "combinatoire" est une des méthodes possibles pour obtenir un "bon" modèle. A partir du modèle saturé, on construit des modèles plus simples en retirant un à un les termes d'interaction. La statistique  $G^2$  croît progressivement et l'on peut arrêter la procédure lorsqu'elle augmente plus rapidement. On retiendra alors le

<sup>1</sup> Cf. Haberman (1974), Nelder et Wedderburn, (1972), McCullagh et Nelder (1989), Christensen, (1990).

<sup>2</sup>  $G^2$  est aussi une mesure de proximité entre les distributions de fréquence  $\hat{T}$  et  $K$  selon la théorie de l'information développée en particulier par Kullback et Leibler (1951), Kullback (1959). En fait la première formule ( $\chi^2$ ) correspond au premier terme

non nul du développement limité de  $G^2$ , en écrivant :  $G^2 = -2 \sum_{i,j,l} k_{ijl} \log \left( 1 + \frac{\hat{t}_{ijl} - k_{ijl}}{k_{ijl}} \right)$ .



modèle correspondant et l'on en déduira les liaisons importantes entre les variables<sup>1</sup>.

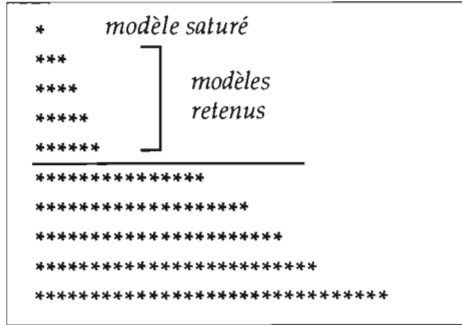


Figure 3.4 - 1  
Histogramme de  $G^2$  et recherche du palier de croissance

Cette méthode combinatoire est applicable aux modèles mettant en jeu un petit nombre de variables. Pour fixer les idées, avec 4 variables, il y a 167 modèles hiérarchiques possibles. Il existe un nombre considérable de travaux sur ce problème de sélection de modèles (problème qui se pose également dans le cas de la régression, mais de façon moins complexe). La multiplication des tests pose des problèmes de *comparaisons multiples* spécifiques (Gabriel, 1969; Aitkin, 1979).

On peut restreindre la recherche aux *modèles graphiques* (sous-ensemble des modèles hiérarchiques) et à l'intérieur de ceux-ci aux *modèles décomposables*. Whittaker (1990) fait une présentation générale des modèles graphiques et une revue des problèmes de sélection des modèles log-linéaires graphiques<sup>2</sup>.

### 3.4.4 La régression logistique

Dans le cadre d'une régression logistique, la problématique est différente mais le modèle utilisé est étroitement lié au modèle log-linéaire.

La régression logistique, comme l'analyse discriminante, cherche à décrire la liaison entre une variable nominale  $y$  (variable à expliquer) et un ensemble de  $p$  variables ( $x_1, x_2, \dots, x_p$ ). On veut également connaître l'effet d'une

<sup>1</sup> On note que l'estimation du critère d'Akaike (1973), fonction de la statistique  $G^2$ , est souvent utilisé pour sélectionner un modèle et mesurer sa qualité. Elle offre l'avantage d'être obtenue sans étudier l'ensemble des modèles possibles (cf. Morineau et al., 1995). Ce critère équivaut asymptotiquement à la validation croisée (Stone, 1977).

<sup>2</sup> Les références de base sur les modèles graphiques sont Wermuth (1976), et Darroch et al. (1980). Pour une synthèse récente, voir Wermuth et Cox (1992). On pourra consulter Fine (in : Droesbeke et al., 1992), de Falguerolles et Jmel (1993).

variable sur la variable à expliquer en tenant compte des liaisons qu'elle entretient avec les autres variables du modèle.

Le plus souvent la variable à expliquer est dichotomique et les variables explicatives sont nominales. Les  $n$  individus caractérisés par l'ensemble des  $p$  variables sont partitionnés en deux groupes définis par les modalités de la variable  $y$ . C'est dans ce cadre que nous nous plaçons.

Pour reprendre l'exemple du paragraphe 3.4.1, on désire étudier par exemple l'influence de la dose de radiation reçue et de l'âge des individus au moment des accidents sur le risque de décès par leucémie.

### a – Le modèle logistique

On suppose que la probabilité qu'un individu a d'appartenir au premier groupe  $I_1$  ( $y = 1$ ) dépend des valeurs des variables explicatives ( $x_1, x_2, \dots, x_p$ ) observées sur cet individu.

On note  $\mathbf{x}$  le vecteur dont les  $p$  composantes sont les valeurs des variables explicatives.

Le modèle logistique se propose de fournir une estimation de cette probabilité notée  $\pi(\mathbf{x})$  :

$$\pi(\mathbf{x}) = P(I_1 | \mathbf{x}) = P(y = 1 | \mathbf{x}).$$

Le théorème de Bayes (§ 3.3.5.a) nous permet d'écrire dans le cas de deux groupes  $I_1$  et  $I_2$  :

$$P(I_1 | \mathbf{x}) = \frac{P(\mathbf{x} | I_1)P(I_1)}{P(\mathbf{x} | I_1)P(I_1) + P(\mathbf{x} | I_2)P(I_2)}$$

qui s'écrit encore :

$$P(I_1 | \mathbf{x}) = \frac{\frac{P(\mathbf{x} | I_1)P(I_1)}{P(\mathbf{x} | I_2)P(I_2)}}{1 + \frac{P(\mathbf{x} | I_1)P(I_1)}{P(\mathbf{x} | I_2)P(I_2)}} \quad [3.4 - 2]$$

Cette formule ne fait intervenir que les quotients des deux probabilités conditionnelles de l'observation  $\mathbf{x}$ .

Dans le cas multinormal avec matrices des covariances  $\Sigma$  égales dans les deux groupes, chacune des deux probabilités conditionnelles s'écrit, pour  $k = 1, 2$  :

$$P(I_k | \mathbf{x}) = (2\pi)^{-p/2} |\Sigma|^{-1/2} \exp\left\{-\frac{1}{2}(\mathbf{x} - \mu_k)' \Sigma^{-1}(\mathbf{x} - \mu_k)\right\}$$

Le quotient des probabilités pondérées fait disparaître les termes du second degré en  $\mathbf{x}$  et s'écrit comme l'exponentielle d'une forme linéaire en  $\mathbf{x}$  avec terme constant (fonction affine de  $\mathbf{x}$ ) :

$$\frac{P(\mathbf{x} | I_1)P(I_1)}{P(\mathbf{x} | I_2)P(I_2)} = \exp\{\beta' \mathbf{x} + b\}$$

Pour alléger les notations, le vecteur  $\mathbf{x}$  désignera désormais un vecteur à  $p+1$  composantes (avec  $x_0 = 1$  et les autres composantes égales à celles de l'ancien  $\mathbf{x}$ ) et le nouveau vecteur de coefficients sera désigné par  $\alpha$ , de sorte que  $\beta' \mathbf{x} + b$  s'écrit maintenant  $\alpha' \mathbf{x}$ .

Ceci permet de réécrire la formule [3.4 - 2] et conduit à l'expression du modèle logistique :

$$\pi(\mathbf{x}) = \frac{\exp\{\alpha' \mathbf{x}\}}{1 + \exp\{\alpha' \mathbf{x}\}} = \frac{\exp\left\{\sum_{j=0}^p \alpha_j x_j\right\}}{1 + \exp\left\{\sum_{j=0}^p \alpha_j x_j\right\}}, \quad [3.4 - 3]$$

où les  $\alpha_j$ , composantes du vecteur  $\alpha$ , sont les coefficients inconnus du modèle. Il s'agit d'un modèle qui ne fait pas intervenir de termes d'interaction entre les variables explicatives<sup>1</sup>.

On peut écrire [3.4 - 3] sous la forme :

$$\frac{\pi(\mathbf{x})}{1 - \pi(\mathbf{x})} = \exp\{\alpha' \mathbf{x}\},$$

ou encore :

$$\log \frac{\pi(\mathbf{x})}{1 - \pi(\mathbf{x})} = \alpha' \mathbf{x} = \sum_{j=0}^p \alpha_j x_j \quad [3.4 - 4]$$

La fonction :

$$F(\pi(\mathbf{x})) = \log \frac{\pi(\mathbf{x})}{1 - \pi(\mathbf{x})}$$

est appelée fonction *Logit*.

### Remarques

1) Les modalités de la variable nominale seront codées 0 ou 1. Comme pour l'analyse de la variance, on élimine, pour chaque variable nominale, une de ses modalités. Le coefficient associé est égal à 0 et cette modalité est appelée traditionnellement "situation de référence" : on mesure en fait les différences avec la ou les autres modalités de la même variable.

2) Le modèle logistique, ou de régression logistique, ou de discrimination logistique, s'applique à une famille de distributions de  $\mathbf{x}$  plus générale que la loi multinormale

<sup>1</sup> Le modèle a été proposé originellement par Cornfield (1962). Étudié notamment par Cox (1972), il a été situé dans le cadre du modèle linéaire généralisé (cf. section 3.3) par Nelder et Wedderburn (1972). Une revue de ses applications en analyse discriminante est faite par Anderson (1982). Cf. également Hosmer et Lemeshow (1989), Devaud (1985).

avec matrices de covariances égales qui nous a servi à l'introduire. Il suffit, on l'a vu, que le quotient des probabilités conditionnelles s'exprime comme l'exponentielle d'une fonction affine de  $\mathbf{x}$ . Ceci est le cas de la plupart des distributions de la famille exponentielle (cf. § 3.2.8.b) dans certaines conditions (Anderson, 1982).

### b – Estimation et tests des coefficients

Pour estimer les coefficients  $\alpha_j$  du modèle, on utilise le plus souvent la méthode du maximum de vraisemblance.

Les  $n$  observations  $(y_i, \mathbf{x}_i)$  [où  $\mathbf{x}_i = (x_{i1}, x_{i2}, \dots, x_{ip})'$ ] sont indépendantes et les  $y_i$  sont des variables de Bernoulli.

La vraisemblance  $\mathcal{L}(\alpha, y_i)$  pour une observation s'écrit :

$$\mathcal{L}(\alpha, y_i) = \pi(\mathbf{x}_i)^{y_i} [1 - \pi(\mathbf{x}_i)]^{1-y_i}$$

et pour l'ensemble des observations, on a :

$$\mathcal{L}(\alpha, \mathbf{y}) = \prod_{i=1}^n \mathcal{L}(\alpha, y_i) = \prod_{i=1}^n \pi(\mathbf{x}_i)^{y_i} [1 - \pi(\mathbf{x}_i)]^{1-y_i}$$

La procédure d'estimation revient à rechercher la valeur  $\hat{\alpha}$  de  $\alpha$  qui maximise le logarithme de la vraisemblance :

$$\log[\mathcal{L}(\alpha, \mathbf{y})] = \sum_i \left[ y_i \log \frac{\pi(\mathbf{x}_i)}{1 - \pi(\mathbf{x}_i)} + \log[1 - \pi(\mathbf{x}_i)] \right]$$

soit encore en exprimant  $\pi(\mathbf{x}_i)$  par la fonction Logit (cf. formule [3.4 - 4]:

$$\log[\mathcal{L}(\alpha, \mathbf{y})] = \sum_i y_i \alpha' \mathbf{x}_i + \sum_i \log[1 + \exp(\alpha' \mathbf{x}_i)]$$

Pour apprécier l'éventuelle *non-influence* d'une variable ou d'une modalité  $x_j$  sur la variable  $y$ , on teste l'hypothèse nulle  $H_0$  :

$$(H_0) : \quad \alpha_j = 0$$

On considère alors la statistique de Student :

$$t = \frac{\hat{\alpha}_j}{\sqrt{\text{Var}(\hat{\alpha}_j)}}$$

où  $\hat{\alpha}_j$  est la  $j^{\text{ème}}$  composante de l'estimateur  $\hat{\alpha}$  et  $\text{Var}(\hat{\alpha}_j)$  est la variance estimée associée à cette composante<sup>1</sup>.

Pour tester l'influence d'une variable nominale à  $q$  modalités, on procède à un test de nullité des  $q$  coefficients  $\alpha_j$  affectés à ses modalités. D'une manière générale, l'hypothèse  $H_0$  stipulant une éventuelle non-influence d'un

<sup>1</sup> On peut également tester la significativité du coefficient  $\alpha_j$  à partir de la *statistique de Wald* qui est le carré de celle de Student.

ensemble de  $q$  variables  $(x_1, x_2, \dots, x_p)$  sur  $y$ , s'exprime par la nullité des  $q$  coefficients associés :

$$(H_0) : \quad \alpha_1 = \alpha_2 = \dots = \alpha_q = 0$$

Notons  $\hat{\alpha}_0$  l'estimateur des  $\alpha_j$  sous l'hypothèse  $H_0$  et  $\hat{\alpha}$  l'estimateur des coefficients du modèle alternatif.

On teste l'hypothèse nulle en calculant la statistique du rapport de vraisemblance :

$$\Lambda = 2(\ell(\hat{\alpha}, \mathbf{y}) - \ell(\hat{\alpha}_0, \mathbf{y}))$$

On démontre qu'elle suit une distribution du  $\chi^2$  à  $q$  degrés de liberté sous des hypothèses de travail convenables. Si l'hypothèse nulle est rejetée, on en déduit qu'au moins une des  $q$  variables (ou une modalité de la variable nominale) influe sur la variable  $y$ .

### c – Comparaison de deux modèles

Considérons deux modèles emboîtés : le modèle 1 à  $p$  variables explicatives et le modèle 2 à  $p + q$  variables explicatives comportant entre autres celles du modèle 1. Choisir le modèle 1, c'est supposer nuls les  $q$  coefficients existant dans le modèle 2 et non dans le modèle 1.

En référence au test de nullité d'un ensemble de coefficients, on retiendra le modèle 1 si l'hypothèse de nullité des  $q$  coefficients n'est pas rejetée, c'est-à-dire si la statistique du rapport de vraisemblance  $\Lambda$  est inférieure à la valeur critique du  $\chi^2$  à  $q$  degrés de liberté<sup>1</sup>.

### d – Modèle avec interaction

Un fois établi le modèle logistique réduit, certains utilisateurs proposent, pour affiner les résultats, d'introduire des termes d'interaction entre les prédicteurs. Pour cela, on ajoute certains produits des  $x_j$ .

Par exemple pour un modèle à deux variables explicatives, le modèle s'écrira :

$$F(\pi(x)) = \alpha_0 + \alpha_1 x_1 + \alpha_2 x_2 + \alpha_{12} x_1 x_2$$

La notion d'interaction d'ordre élevé est complexe. Un terme d'interaction d'ordre 2 en modèle logistique correspond au terme d'interaction d'ordre 3 en modèle log-linéaire.

<sup>1</sup> En pratique, le choix du modèle logistique repose sur la comparaison de modèles emboîtés. On adopte une procédure pas à pas en commençant par prendre en compte le modèle comportant le plus de variables explicatives que l'on compare à un modèle restreint comprenant un sous-ensemble des prédicteurs. On procédera généralement par élimination progressive des variables ne modifiant pas de manière significative la vraisemblance jusqu'à avoir un modèle ne pouvant plus être réduit. Cette procédure n'assure cependant qu'un optimum local.

### 3.4.5 Modèles log-linéaire et analyse des correspondances

Le modèle log-linéaire et l'analyse des correspondances ne répondent pas aux mêmes préoccupations et ne fournissent pas des résultats de même nature. Ce sont en fait des techniques complémentaires.

Le modèle logistique, qui est un véritable modèle explicatif, est plutôt à comparer aux techniques d'analyse discriminante de la section précédente. Comme elles, il peut gagner à être mis en œuvre dans un environnement de méthodes exploratoires, mais il n'est pas en concurrence directe avec ces dernières méthodes.

D'assez nombreux travaux ont porté sur la comparaison des différentes approches dans des contextes d'application divers, parfois sensiblement éloignés des contextes réels<sup>1</sup>.

#### a – Des champs d'application différents

Bien que s'appliquant aux mêmes types de variables, les variables nominales, ces deux méthodes ont des problématiques et des champs d'application différents.

Le *modèle log-linéaire* s'applique avec profit lorsque l'on dispose de peu de variables (rarement plus de cinq variables surtout si elles ont beaucoup de modalités) avec cependant beaucoup d'individus, pour que les cellules de l'hypertable de contingence obtenue en croisant les variables ne soient pas vides. Le nombre des sous-modèles explicitant les liaisons entre les variables augmente beaucoup plus vite que le nombre de variables. On augmente alors le nombre de coefficients à tester et donc les chances de trouver des effectifs nuls, ce qui rend les résultats plus instables. De ce fait, le modèle log-linéaire est bien adapté lorsque le problème posé permet de

---

<sup>1</sup> On ne mentionnera ici qu'un petit nombre de publications sur ce thème en suivant un ordre chronologique : Daudin et Trécourt (1980) sont parmi les premiers à faire une comparaison sur une table de contingence à 6 entrées ( $21 \times 2 \times 2 \times 2 \times 2$ ) entre une des analyses des correspondances possibles et le modèle log-linéaire. Escoufier (1982), Lauro et Decarli (1982) proposent également des rapprochements entre utilisations des méthodes. Leclerc et al. (1985) comparent sur un exemple approfondi l'analyse des correspondances et la régression logistique. Van der Heijden et de Leeuw (1985), Van der Heijden (1987), puis Van der Heijden et al. (1989) proposent une méthodologie de l'utilisation simultanée de l'analyse des correspondances et des modèles log-linéaires en préconisant de décrire par des analyses des correspondances les résidus des modèles log-linéaires. D'autres comparaisons et applications se trouvent dans Worsley (1987) et plus généralement dans le numéro spécial 35 -3 (1987) de la Revue de Statistique Appliquée, animé par le L.S.P. de l'Université Paul Sabatier. Cf. également Hudon (1990), Tenenhaus et al. (1993). Gilula (1986), Gilula et Ritov (1990), Goodman (1986, 1991) étudient les performances de l'analyse des correspondances et des modèles log-linéaires dans le contexte d'utilisation des modèles qu'ils ont eux-mêmes développés pour les tables de contingences multiples ou à modalités ordonnées (approche confirmatoire pour des tables de dimensions très réduites).

procéder à une sélection préalable des variables et de formuler les hypothèses nulles.

*L'analyse des correspondances binaires* (sur vraies tables de contingence, que l'on appelle parfois tables de contingence binaire ou à double entrées) s'applique avec profit lorsque les deux partitions mises en correspondances (colonnes et lignes actives) sont relativement importantes : par exemple, tables de contingence croisant 95 départements métropolitains et 12 causes de décès, tables croisant 373 communes de la région parisienne et 29 catégories socio-professionnelles. Pour des petites tables de contingence, la fonction de l'analyse des correspondances est surtout iconographique, illustrative<sup>1</sup>.

*L'analyse des correspondances multiples* (sur tableaux disjonctifs complets) est utile pour les tableaux de type "sous-fichiers d'enquête" : en général une à plusieurs dizaines de variables nominales, de 200 à 20 000 individus. Il n'est pas rare que l'hypertable de contingence soit à 99% vide<sup>2</sup>.

Qu'il s'agisse de correspondances binaires ou multiples, la dichotomie entre *variables actives et illustratives* est fondamentale. C'est elle qui permet de confronter une information illimitée au sous-espace des variables actives, dont la description ne constitue qu'une phase préliminaire.

Enfin, ces méthodes ne permettent que de décrire des tableaux. Et une table de contingence multiple permet de construire différents types de tableaux. Si l'on s'intéresse aux interactions d'ordre élevé entre certaines variables, on construira de nouvelles variables en croisant ces variables et en considérant selon les cas la nouvelle variable comme active ou supplémentaire.

C'est le problème sous-jacent qui permet de guider la démarche : choix des tableaux à décrire dans un cas, choix des modèles à sélectionner et à éprouver dans l'autre. Rappelons également que l'usage simultané de la classification et des analyses en axes principaux fait partie intégrante de la démarche exploratoire.

Le tableau 3.4 - 1 résume ces différences d'objectifs et d'applications dans le cas de l'analyse des correspondance multiples.

Certains travaux de confrontation entre méthodes perdent de leur portée en raison de la méconnaissance des vocations (essentiellement attestée par une expérience pratique) de chacune des approches. Il est vrai que le *paradoxe pédagogique* inhérent à l'analyse des données - *comment prouver sur un modèle réduit l'efficacité de méthodes qui ne sont utiles et profitables que sur de grands tableaux* - ne facilite pas la tâche d'explication de la vocation réelle de ces méthodes.

<sup>1</sup> Nous reviendrons au chapitre suivant sur la validité des résultats et la méthodologie générale d'emploi des méthodes.

<sup>2</sup> Ainsi, pour une petite batterie de 10 questions à 4 modalités posées à 1000 répondants, l'hypertable présente  $4^{10}$  ( $\approx 10^9$ ) cases; moins d'une case sur 1000 sera non-vide.

**Tableau 3.4 - 1**  
**Vocations spécifiques des deux approches**

Modèle log-linéaire	Correspondances multiples
<ul style="list-style-type: none"> <li>- Description des interactions entre plus de deux variables dans un cadre inférentiel.</li> <li>- Des hypothèses sur les liaisons doivent être formulées au préalable.</li> <li>- Est limité à peu de variables (en pratique moins de 5).</li> <li>- Met en jeu toutes les cases d'un hypercube de contingence :</li> </ul>	<ul style="list-style-type: none"> <li>- Description des liaisons entre les variables prises deux à deux sous forme essentiellement graphique.</li> <li>- N'impose aucune hypothèse sur les liaisons, mais impose une certaine homogénéité de l'ensemble des variables actives.</li> <li>- N'est pas limitée dans le nombre de variables</li> <li>- Met seulement en jeu les faces de l'hypercube représentées par le tableau de Burt :</li> </ul>
<ul style="list-style-type: none"> <li>- Méthode par essence confirmatoire, utilisée pour explorer l'univers des modèles. On cherche celui ou ceux qui s'adaptent le mieux aux observations.</li> <li>- Les individus n'apparaissent pas.</li> <li>- La notion de variable supplémentaire n'est pas directement pertinente.</li> </ul>	<ul style="list-style-type: none"> <li>- Méthode descriptive et exploratoire de la structure intrinsèque des données.</li> <li>- Les individus peuvent jouer un rôle central. L'analyse sert souvent à produire des typologies d'individus.</li> <li>- La notion de variable supplémentaire est fondamentale.</li> </ul>

Il faut reconnaître cependant que si l'analyse des correspondances est bien utile dans le cas des grandes tables de contingences à deux entrées et dans le cas des grands tableaux disjonctifs complets, elle est beaucoup plus délicate à utiliser dans le cas intermédiaire des petites tables de contingence multidimensionnelles.

Pour ce type de tableau aux facettes peu nombreuses, l'intérieur de la table de contingence (croisements de plus de deux variables), s'il contient des effectifs suffisants, est intéressant à décrire de façon détaillée. Une analyse des correspondances multiples sur un tableau comportant trois ou quatre



variables nominales donne des résultats assez grossiers, d'une stabilité douteuse.

Il existe en la matière des savoir-faire, sans qu'une méthodologie rigoureuse se soit imposée définitivement : on peut juxtaposer des tranches en ligne ou en colonnes (cf. par exemple van der Heijden (1987) pour le cas des données longitudinales) ; juxtaposer des tableaux obtenus par croisements des variables initiales ; positionner en éléments supplémentaires les croisements de variables deux à deux dans les plans factoriels d'une analyse des correspondances multiples ; dans certains cas, réaliser une analyse factorielle multiple (cf. § 3.8.3). D'autres approches seront évoquées plus loin. C'est à propos de ce type d'applications que l'on pourra parler de complémentarité entre les méthodes.

### **b – Liens théoriques entre l'analyse des correspondances et les modèles log-linéaires**

L'analyse des correspondances analyse l'écart entre un tableau de fréquence  $f_{ij}$  et un tableau modèle  $f_i f_j$  correspondant à l'hypothèse d'indépendance. Lorsque cet écart est significatif<sup>1</sup>, elle décrit de façon suggestive les associations privilégiées entre lignes et colonnes responsables de cet écart.

Ce principe d'analyse est manifestement insuffisant pour les tables de contingence à plus de deux entrées. Certes, l'analyse des correspondances multiples constitue *une* généralisation possible de cette démarche, réalisant une sorte de compromis entre tous les croisements des variables prises deux à deux. Cette généralisation est opératoire lorsque le nombre et la nature des variables nominales exclut une étude méthodique de leurs interactions : on a alors à traiter un tableau (individus  $\times$  variables), comme en analyse en composantes principales.

Mais il n'existe pas d'analogue du théorème d'Eckart et Young dans le cas des tableaux tridimensionnels<sup>2</sup>. Il ne peut donc exister dans ce cas de démarche exploratoire aussi bien assise que dans le cas des tableaux à double entrée.

La démarche proposée par van der Heijden et de Leeuw (1985) puis développée par van der Heijden (1987), qui s'apparente aux analyses partielles évoquées à la section 3.6, va effectivement dans le sens d'une

<sup>1</sup> Le classique  $\chi^2$  permet d'alerter l'utilisateur sur la signification de cet écart, mais les premières valeurs propres de l'analyse des correspondances, ainsi que les taux d'inertie correspondants, peuvent également mesurer des écarts que le  $\chi^2$  ne décèle pas; cf. § 1.3.4.a.

<sup>2</sup> Ce que l'on peut exprimer dans les termes suivants : il existe une décomposition hiérarchique unique d'un élément du produit tensoriel de deux espaces euclidiens en une somme de produits tensoriels de vecteurs appartenant à chacun des deux espaces. Mais une telle décomposition n'est pas unique dans le cas de d'un élément du produit tensoriel de plus de deux espaces euclidiens (cf. Benzécri, 1973; Tome 2B, n°6 [RED.TENS.]).

utilisation synergique des deux approches : utiliser le modèle log-linéaire pour éliminer l'effet complexe de certaines variables et utiliser l'analyse des correspondances pour décrire les résidus que le modèle log-linéaire ne permet pas d'expliquer.

Elle rejoint une généralisation de l'analyse des correspondances introduite par Escofier (1984) qui permet d'introduire des modèles moins restrictifs. L'analyse factorielle des correspondances se généralise à un modèle différent du modèle d'indépendance en supposant que les marges du tableau de référence sont distinctes de celles du tableau étudié.

Les liens théoriques entre l'analyse des correspondances et les modèles log-linéaires sont très ténus, même dans des contextes relativement simples. Après Escoufier (1982), Worsley (1987), van der Heijden *et al.* (1989), écrivons ce que pourrait être un modèle de l'analyse des correspondances dans le cas d'une approximation bi-dimensionnelle de la loi  $f_{ij}$ .

La formule de reconstitution des données en analyse des correspondances (cf. section 1.3.3.h) peut s'écrire, en retenant deux axes :

$$f_{ij} \approx f_{i.} \cdot f_{.j} \left\{ 1 + \sum_{h=1}^2 \sqrt{\lambda_h} \varphi_h(i) \psi_h(j) \right\}$$

ce qui suggère un modèle de la forme :

$$f_{ij} \approx e_{ij} = c p_i q_j (1 + r_{1i} s_{1j} + r_{2i} s_{2j})$$

où les coefficients inconnus, assujettis aux mêmes contraintes que leurs homologues de la formule de reconstitution, sont déterminés de façon à rendre minimale une distance entre  $f_{ij}$  et  $e_{ij}$ <sup>1</sup>.

Le modèle peut s'écrire, si les valeurs propres  $\lambda_1$  et  $\lambda_2$  sont petites par rapport à 1, ce qui est le cas au voisinage de l'indépendance :

$$\log f_{ij} \approx \log e_{ij} = a_0 + a_i + b_j + r_{1i} s_{1j} + r_{2i} s_{2j}$$

alors qu'un modèle log-linéaire saturé s'écrit :

$$\log e_{ij} = a'_0 + a'_i + b'_j + u_{ij}$$

Ainsi, l'analyse des correspondances suggère de décomposer le terme d'interaction  $u_{ij}$  sous forme simplement multiplicative dans le cas d'un seul facteur, et plus généralement sous forme de matrice de rang  $q$  dans le cas où l'on retient  $q$  facteurs.

Il est vrai que dans le cas d'une table de contingence à double entrée, le modèle log-linéaire non-saturé est trivial (hypothèse d'indépendance) et le

<sup>1</sup> Distance du  $\chi^2$ , critère de Kullback-Leibler (cf. section 3.4.3-b), ou encore critère de la déviance, très utilisé pour les modèles logistiques (cf. par exemple Celeux et Nakache, 1994).

modèle saturé aussi (ajustement parfait). D'où les tentatives de donner au terme d'interaction des formes plus simples, avec en particulier les modèles dit *RC*, puis *multifactor* de Goodman (cf. Goodman, 1986). L'analyse des correspondances, qui revient à une décomposition aux valeurs singulières de la matrice normée (que l'on peut appeler matrice d'interaction) :

$$\frac{f_{ij} - f_{i.}f_{.j}}{\sqrt{f_{i.}f_{.j}}}$$

répond à une même préoccupation<sup>1</sup>.

Le cas des tables de contingences multiples est beaucoup plus complexe, et dans les configurations où le modèle log-linéaire peut être appliqué (peu de variables, beaucoup d'individus, des idées a priori sur le rôle de telle ou telle variable) l'approche "analyse de résidus" mentionnée plus haut paraît bien appropriée.

### c – Difficultés de l'articulation exploration-inférence

Lorsque l'on est en situation trop exploratoire pour pouvoir formuler des hypothèses, ou lorsque le nombre de variables est trop élevé par rapport au nombre des individus pour pouvoir construire un modèle pertinent, on a recours à l'analyse des correspondances multiples.

Son utilisation permet d'une part de déceler, dans un premier temps, les liaisons intéressantes entre certaines variables, et d'autre part de sélectionner et réduire les variables et leurs modalités. Rappelons que l'on travaille sur les "faces de l'hypercube" c'est-à-dire sur les cumuls de fréquences correspondant à des effectifs importants.

On pourrait penser tester les liaisons par des modèles log-linéaires afin de préciser et de mesurer le niveau et l'intensité de celles-ci (l'intérieur de l'hypercube, lorsque le nombre d'individus le permet). Cette démarche demande cependant une certaine prudence.

Ce serait en effet une erreur de raisonnement (malheureusement répandue chez les praticiens) de penser que l'on peut tester sur des données un modèle suggéré par les mêmes données.

Comme l'a spécifié Cox (1977) dans un remarquable article de synthèse sur les tests de signification, l'articulation *exploratoire - confirmatoire* pose des problèmes d'une grande complexité, analogues à ceux que nous avons rencontrés dans la section précédente à propos de l'analyse discriminante : tester une fonction discriminante sur l'échantillon d'apprentissage donne une idée trop optimiste de son pouvoir de prédiction.

Dans les deux cas en effet, les échantillons, *et donc les fluctuations qui leurs sont propres*, sont sollicités soit pour construire une fonction ou une règle

<sup>1</sup> Elle effectue cette décomposition dans un cadre géométrique euclidien simple, en produisant des visualisations assorties de règles d'interprétation.

de classement (cas de l'analyse discriminante) soit pour choisir un modèle (cas d'une analyse des correspondances multiples préalable à un modèle log-linéaire).

La difficulté est accentuée par l'effet "comparaisons multiples"<sup>1</sup> que l'on peut craindre dans la mesure où l'analyse des correspondances multiples peut traiter simultanément plusieurs dizaines, voire des centaines de variables.

Même lorsque le tableau contenant  $p$  variables nominales est généré selon un modèle stipulant l'indépendance totale entre les  $p$  variables, un certain nombre de paires de variables (parmi les  $p(p-1)/2$  paires possibles) peut donner lieu à des liaisons significatives selon les valeurs usuelles des seuils, et ceci d'autant plus facilement que  $p$  est grand. Un modèle restreint à cette sélection de variables pourrait effectivement confirmer une structure qui ne serait en fait qu'un artefact.

Il existe au moins deux types de solutions pragmatiques pour contourner ces difficultés : travailler sur un échantillon supplémentaire (échantillon-test, validation croisée) comme dans le cas de la discrimination; travailler avec des seuils de signification plus sévères au niveau de la lecture des modèles (comme dans le cas de comparaisons multiples)<sup>2</sup>.

---

<sup>1</sup> Cf. par exemple la section 1.1, § 1.4.4.a; et ci-dessus, la section 3.4.3.c à propos de la sélection des modèles log-linéaires.

<sup>2</sup> Remarquons que la démarche "analyse des correspondances des résidus d'un modèle log-linéaire" mentionnée plus haut, qui correspond à une articulation en sens inverse : *Inférence -Exploration*, ne prête pas le flanc à ces critiques. Elle correspond à une situation méthodologique plus particulière, pour laquelle les modèles log-linéaires pouvaient être utilisés d'emblée. L'approche exploratoire est cependant, en général, et presque par nature, la première phase des investigations.

---

## Section 3.5

---

# Segmentation

Les méthodes de segmentation cherchent à résoudre les problèmes de discrimination et de régression en segmentant de façon progressive l'échantillon pour obtenir un *arbre de décision binaire*. La voie a été ouverte par Sonquist et Morgan (1964) et Morgan et Messenger (1973) avec la méthode dite AID (*Automatic Interaction Detection*)<sup>1</sup>. De nombreuses contributions ont suivi, mais les travaux de Breiman, Friedman, Olshen et Stone (1984) ont renouvelé l'approche et suscité un regain d'intérêt pour la segmentation. Leur méthode, connue sous le nom de CART (*Classification And Regression Tree*), diffère de l'AID par le mode de construction de l'arbre et la technique d'*élagage* conduisant à un sous-arbre exploitable ayant des propriétés satisfaisantes<sup>2</sup>.

La segmentation par la méthode CART vient donc concurrencer les méthodes plus classiques que sont la régression multiple, l'analyse discriminante et la régression logistique. Elle présente des avantages importants dont le premier est sans doute la lisibilité des règles d'affectation, l'interprétation des résultats étant directe et intuitive. Par ailleurs la technique est non-paramétrique et peu contrainte par la nature des données. On peut en effet utiliser en même temps comme variables explicatives, des variables continues, ordinales et nominales sans codage préalable. De plus, la technique fournit d'office la sélection des variables à utiliser en tenant compte d'éventuelles interactions. Elle est robuste vis-à-vis de données erronées ou de valeurs aberrantes et gère les données manquantes aussi bien dans la construction de l'arbre et l'estimation de son risque que dans l'application de la règle à un nouveau sujet. Enfin c'est le même principe, la même méthode, le même algorithme qui sont mis en œuvre pour analyser une variable nominale (discrimination) et une variable continue (régression).

Cependant, les règles d'affectation pourront paraître parfois "abruptes" et trop sensibles à de légères perturbations des données. Il apparaîtra parfois difficile de décider quel est l'arbre "optimal". On peut également regretter l'absence d'une fonction globale mettant en jeu l'ensemble des variables (fonction linéaire discriminante ou équation de régression) qui prive l'utilisateur d'une représentation géométrique sous forme de configurations de points dans l'espace.

---

<sup>1</sup> Cf. Bouroche et Tenenhaus (1970).

<sup>2</sup> On pourra se reporter pour des éléments théoriques à l'ouvrage cité de Breiman *et al.*, et pour une présentation pratique à l'article de Guegen et Nakache (1988) et aux deux ouvrages édités par Celeux (1990) et Celeux et Nakache (1994).

### 3.5.1 Formulation du problème, principe et vocabulaire

Comme en régression (linéaire ou logistique) et en discrimination, on est en présence d'un tableau de données contenant une variable privilégiée  $y$  "à expliquer" par les autres variables du tableau  $x_1, x_2, \dots, x_p$ .

Il s'agit d'une part de sélectionner parmi les variables explicatives celles qui sont les plus discriminantes pour la variable nominale  $y$  (ou celles qui sont le plus liées au phénomène décrit par la variable continue  $y$ ), et d'autre part de construire une règle de décision permettant d'affecter un nouvel individu à l'une des  $k$  classes (cas de la discrimination) ou de lui affecter une valeur  $y$  (cas de la régression).

La méthode de segmentation consiste à rechercher d'abord la variable  $x_j$  qui explique le mieux la variable  $y$ . Cette variable définit une première division de l'échantillon en deux sous-ensembles, appelés *segments*. Puis on réitère cette procédure à l'intérieur de chacun de ces deux segments en recherchant la deuxième meilleure variable, et ainsi de suite <sup>1</sup>.

On construit ainsi un *arbre de décision binaire* par divisions successives de l'échantillon en deux sous-ensembles (figure 3.5 - 1) où l'on distingue :

- les *segments intermédiaires* ou *nœuds* qui engendrent deux segments descendants immédiats,
- les *segments terminaux* qui ne sont plus divisés,
- une *branche* d'un segment  $t$  qui comprend tous les segments descendant de  $t$ ,  $t$  n'étant pas inclus dans la branche,
- l'*arbre binaire complet* noté  $A_{\max}$  pour lequel chaque segment terminal contient un seul individu,
- un *sous-arbre*  $A$  qui est obtenu à partir de  $A_{\max}$  par *élagage* d'une ou de plusieurs branches.

Par ailleurs, la méthode CART, contrairement aux autres méthodes de segmentation, n'impose aucune règle (fondée sur un seuil) d'arrêt de division des segments. Elle fournit, à partir de l'arbre binaire complet, la séquence des sous-arbres obtenue en utilisant une *procédure d'élagage*. Celle-ci est basée sur la suppression successive des branches les moins informatives en terme de discrimination entre les classes ou en terme d'explication de la variable  $y$ .

Au cours de la phase d'élagage, la méthode sélectionne un sous-arbre "optimal" en se fondant sur l'estimation de l'erreur théorique d'affectation ou de prévision à l'aide, soit d'un *échantillon-test* (technique présentée ci-après) quand l'échantillon est suffisamment important, soit de la *validation croisée*.

<sup>1</sup> Notons que cette méthode, contrairement aux autres méthodes multidimensionnelles, ne considère pas simultanément l'ensemble des variables explicatives mais les examine une par une. Cependant, les liaisons entre variables explicatives sont prises en compte aux différentes étapes de la construction de l'arbre.

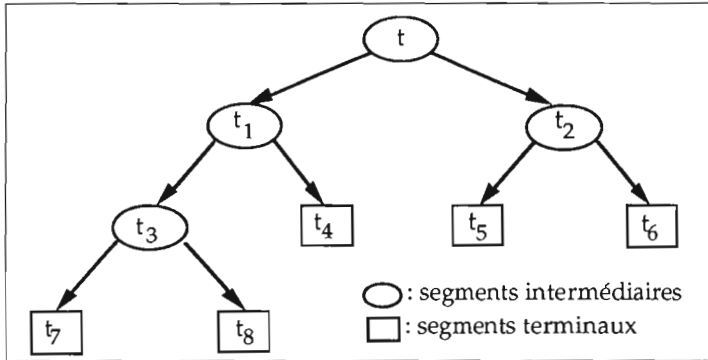


Figure 3.5 - 1  
 Arbre de décision binaire

### 3.5.2 Construction d'un arbre de décision binaire

L'idée de base est d'effectuer la division d'un nœud de telle sorte que les deux segments descendants soient plus homogènes que le nœud parent et qu'ils soient les plus différents possible entre eux vis-à-vis de la variable  $y$ .

Cette procédure nécessite de définir un critère permettant de sélectionner la "meilleure" division d'un nœud. Le critère de la régression différera de celui de la discrimination, mais le principe de construction reste le même dans les deux cas.

Les différentes phases de construction de l'arbre sont les suivantes :

- 1- établir pour chaque nœud l'ensemble des divisions admissibles.
- 2- définir un critère permettant de sélectionner la "meilleure" division d'un nœud.
- 3- définir une règle permettant de déclarer un nœud comme terminal ou intermédiaire.
- 4- affecter chaque nœud terminal à l'un des groupes (cas de la discrimination), ou affecter une valeur à  $y$  pour chaque nœud terminal (cas de la régression).
- 5- estimer le risque d'erreur de classement (cas de la discrimination) ou de prévision (cas de la régression) associé à l'arbre.

#### a – Algorithme général de segmentation

Les variables explicatives peuvent être de nature quelconque. Dans un premier temps, considérons le cas des variables continues. Les étapes de l'algorithme sont les suivantes :

- 1- Au départ, on dispose d'un seul segment contenant l'ensemble des individus.

2 - A la première étape, la procédure de construction de l'arbre examine une par une toutes les variables explicatives.

Pour une variable  $x_j$  donnée, elle passe alors en revue toutes les divisions possibles de la forme  $x_j < \alpha$  où  $\alpha$  est une valeur quelconque contenue dans l'étendue de la variable  $x_j$  considérée.

Chaque division scinde l'échantillon en segments descendants : le segment de gauche  $t_g$  contient les sujets vérifiant  $x_j < \alpha$  et le segment de droite  $t_d$  contient les autres ( $x_j \geq \alpha$ ). De toutes les divisions  $d_j^m$  possibles de  $x_j$ , où  $m$  représente la  $m^{\text{ième}}$  division (soit encore la  $m^{\text{ième}}$  valeur classée de  $x_j$ ), la procédure sélectionne la "meilleure"  $d_j^*$ , au sens d'un critère de division à préciser <sup>1</sup>.

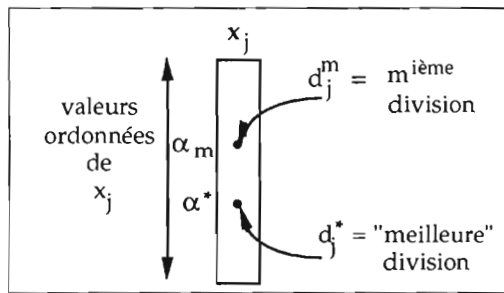


Figure 3.5 - 2  
Divisions possibles pour la variable  $x_j$

On obtient ainsi, pour chacune des  $p$  variables, la meilleure division et l'on retiendra finalement, parmi ces  $p$  divisions, celle, notée  $d^*$ , qui fournit les deux segments les plus "typés" vis-à-vis de  $y$ .

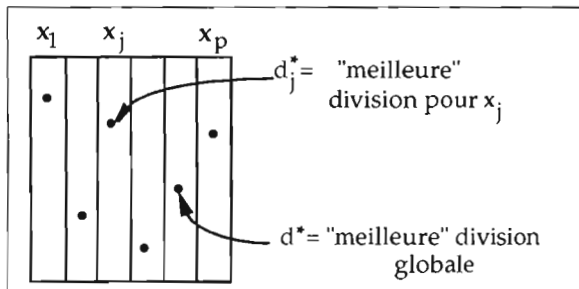


Figure 3.5 - 3  
Meilleures divisions pour l'ensemble des variables

<sup>1</sup> Critère de la variance résiduelle pour la régression (cas d'une variable à expliquer continue), cf. § 3.5.2.b. Critère de la pureté maximale dans le cas de la discrimination, cf. § 3.5.2.c.



- 3 - A l'étape suivante, on applique la même procédure à chacun des deux segments descendants obtenus. Les variables explicatives peuvent être différentes selon les segments.
- 4 - On arrête la procédure lorsque tous les segments sont déclarés terminaux : soit parce qu'ils ne nécessitent plus de divisions soit parce que leur taille est inférieure à un effectif fixé.

Pour un nouvel individu, on définit une règle d'affectation simple en le faisant *descendre* dans l'arbre.

Si, parmi les variables explicatives, certaines sont nominales, elles sont prises en compte de la manière suivante :

- une variable à deux modalités ne peut fournir qu'une seule division,
- une variable à  $k$  modalités ordonnées fournit  $k - 1$  divisions,
- une variable à  $k$  modalités non ordonnées fournit  $2^{k-1} - 1$  divisions; toutes les divisions correspondant aux différents sous-ensembles de modalités sont examinées.

Par exemple, à partir d'une variable **a** à deux modalités, d'une variable **b** à 4 modalités ordonnées et d'une variable **c** à 3 modalités non ordonnées, les divisions possibles d'un nœud en deux segments descendants  $t_g$  (celui de gauche) et  $t_d$  (celui de droite) sont les suivantes<sup>1</sup> :

**Tableau 3.5 - 1**  
Divisions possibles d'un segment par une variable nominale

	$t_g$	$t_d$
<i>var. binaire</i>	(a <sub>1</sub> )	(a <sub>2</sub> )
<i>variable ordonnée (ordinaire)</i>	(b <sub>1</sub> )	(b <sub>2</sub> , b <sub>3</sub> , b <sub>4</sub> )
	(b <sub>1</sub> , b <sub>2</sub> )	(b <sub>3</sub> , b <sub>4</sub> )
	(b <sub>1</sub> , b <sub>2</sub> , b <sub>3</sub> )	(b <sub>4</sub> )
<i>variable non ordonnée</i>	(c <sub>1</sub> )	(c <sub>2</sub> , c <sub>3</sub> )
	(c <sub>2</sub> )	(c <sub>1</sub> , c <sub>3</sub> )
	(c <sub>3</sub> )	(c <sub>1</sub> , c <sub>2</sub> )

## b – Cas de la régression

Lorsque la variable à expliquer  $y$  est continue, le critère de sélection de la "meilleure" division d'un nœud est fondé sur la variance de  $y$  dans les segments descendants. Cette variance doit être plus faible que la variance de  $y$  dans le nœud parent.

<sup>1</sup> Remarquons que la segmentation effectue simultanément un découpage sur la population observée et sur les valeurs des variables explicatives.

- Critère de la variance résiduelle minimale

Pour toute division  $d_j^m$  d'un nœud  $t$  par une variable  $x_j$ , on calcule la moyenne pondérée des variances de  $y$  à l'intérieur de chacun de ses segments descendants  $t_g$  et  $t_d$ , c'est-à-dire la variance résiduelle du nœud  $t$  :

$$\text{var}(d_j^m, t) = \left(\frac{n_g}{n_t} s_g^2\right) + \left(\frac{n_d}{n_t} s_d^2\right)$$

où  $n_g, n_d, n_t$  sont respectivement les effectifs des segments  $t_g, t_d, t$  et  $s_g^2, s_d^2$  sont les variances de la variable continue  $y$  à l'intérieur des segments  $t_g$  et  $t_d$ <sup>1</sup>.

On retient la "meilleure" division  $d_j^*$  réalisée par la variable  $x_j$  qui correspond à la variance résiduelle minimale :

$$\text{var}(d_j^*, t) = \min_{\text{med}_j} \{\text{var}(d_j^m, t)\}$$

où  $d_j$  est l'ensemble des divisions de la variable  $x_j$ .

Parmi toutes les meilleures divisions  $d_j^*$  obtenues à partir des  $p$  variables explicatives, la meilleure division (globale) du nœud  $t$  est effectuée à l'aide de la variable qui assure :

$$\text{var}(d^*, t) = \min_{j=1, \dots, p} \{\text{var}(d_j^*, t)\}$$

- Les étapes de l'algorithme

Considérons un ensemble d'individus sur lesquels on relève les informations concernant une variable continue  $y$  et  $p = 8$  variables explicatives  $x_1, \dots, x_8$ . On suppose que les valeurs de  $y$  ont pour moyenne  $m = 10$  et pour variance  $s^2 = 60$ .

On commence par examiner la variable continue  $x_1$ .

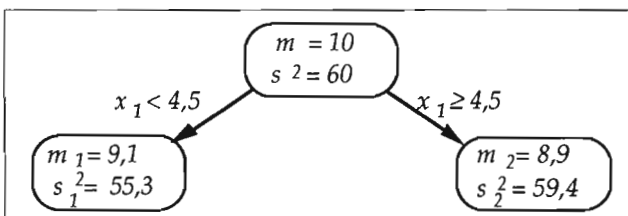


Figure 3.5 - 4  
Régression : meilleure division pour la variable  $x_1$

<sup>1</sup> Il s'agit de la variance interne ou *intra* introduite notamment en analyse discriminante, c'est-à-dire de la variance non expliquée par la coupure.

On retient la valeur de coupure qui minimise la variance à l'intérieur des deux segments descendants, par exemple la division associée à la valeur 4,5 (cf. figure 3.5 - 4)

Mais cette meilleure division obtenue avec  $x_1$  n'est peut-être pas la plus efficace en terme de réduction de la variance. Il faut étudier les autres variables. On recherche, de la même manière, la meilleure division de l'échantillon pour chacune des  $p - 1 = 7$  autres variables. On choisira alors la division qui présente la plus faible moyenne pondérée des variances de  $y$  à l'intérieur des deux segments descendants, par exemple la variable continue  $x_5$  pour la valeur  $\alpha = 7,2$ .

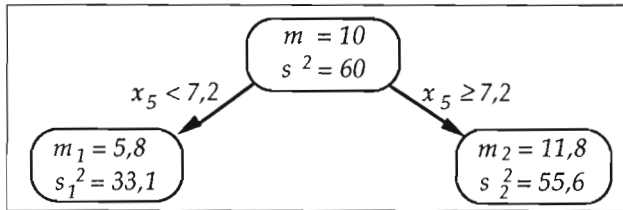


Figure 3.5 - 5  
Régression : meilleure division pour toutes les  $p$  variables

On réitère cette procédure à l'intérieur de chacun des deux segments obtenus  $t_1$  et  $t_2$ . Pour le segment  $t_1$ , ce sera par exemple la variable nominale  $x_7$  à deux modalités ; la meilleure division sera obtenue pour les valeurs  $x_7=1$  (segment  $t_3$ ), et  $x_7=2$  (segment  $t_4$ ). On sélectionnera la variable  $x_2$  à deux modalités, pour le segment  $t_2$ . On aboutit ainsi à l'arbre à deux niveaux représenté sur la figure 3.5 - 6. (Sur cette figure, l'indice bas des variances est celui des segments correspondants :  $s_i^2$  correspond au segment  $t_i$ ).

On pourrait arrêter là la procédure de division et produire l'arbre de prédiction à 4 segments terminaux.

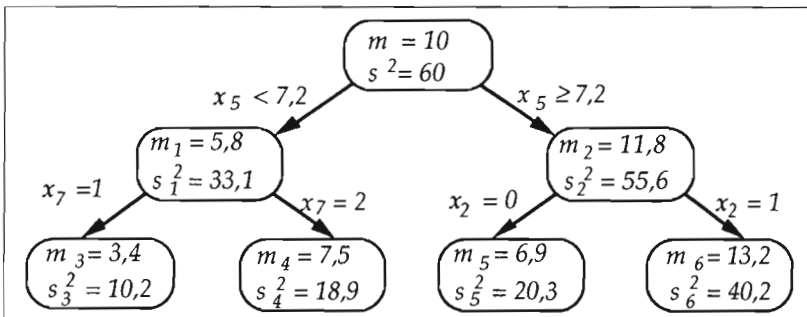


Figure 3.5 - 6  
Régression : Arbre à deux niveaux

### - Règle d'affectation

Considérons alors un nouvel individu  $i$  dont on cherche à prévoir la valeur de  $y_i$ . Il tombera dans un de ces 4 segments terminaux après avoir parcouru un chemin de l'arbre suivant les valeurs qu'il présente pour  $x_5$ ,  $x_7$  et  $x_2$ . La valeur affectée à  $y_i$  sera la moyenne dans le segment et l'écart-type correspondra à celui du segment.

### - Erreur Apparente de Prévion associée à un arbre A

Si certaines variances des segments sont encore importantes, on peut continuer les divisions dans le but de réduire davantage les variances des segments terminaux.

Ainsi on associe à chaque segment terminal  $t$  de l'arbre A l'erreur  $R_t$  suivante :

$$R_t = \frac{n_t}{n} \times s_t^2$$

où  $n$  est le nombre total d'individus,  $n_t$  est le nombre d'individus du segment  $t$ ,  $s_t^2$  est la variance de la variable  $y$  à l'intérieur du segment  $t$  c'est-à-dire :

$$s_t^2 = \frac{1}{n_t} \sum_i (y_i - \bar{y}_t)^2$$

avec  $\bar{y}_t$ , la moyenne des valeurs  $y_i$  des individus du segment  $t$ .

L'Erreur Apparente de Prévion (EAP) associée à l'arbre A vaut :

$$EAP(A) = \sum_{t \in A} R_t \quad [3.5 - 1]$$

et correspond à la moyenne pondérée des variances de  $y$  dans chacun des segments terminaux de l'arbre A. Le rapport  $EAP(A)/s^2$  est l'équivalent de l'expression  $(1 - R^2)$  de la régression linéaire multiple<sup>1</sup> et représente le pourcentage de la variance totale non expliquée par les variables  $x_1, x_2, \dots, x_p$ .

Plus on divise, plus les variances décroissent pour être finalement nulles quand chaque segment terminal contient un seul individu. Au grand arbre complet noté  $A_{\max}$  ainsi obtenu est affectée une Erreur Apparente de Prévion nulle.

### c - Cas de la discrimination

Lorsque la variable  $y$  est nominale et répartit les individus en  $k$  classes, la sélection d'une division doit être telle que les segments descendants soient plus "purs" que le nœud parent. Autrement dit, il faut que le mélange des

<sup>1</sup> Dans la régression linéaire multiple, on suppose que la variance de la réponse  $y$  conditionnellement aux covariables (variables explicatives) est constante, ce qui n'est pas le cas pour la régression par arbre.

classes soit moins important dans les segments descendants que dans le nœud parent.

- *Critère de la pureté maximale*

A chaque segment  $t$  est donc associée une mesure de l'impureté  $i(t)$  définie par :

$$i(t) = \sum_r^k \sum_s^k P(r|t)P(s|t)$$

avec  $r \neq s$  et où  $P(r|t)$  et  $P(s|t)$  sont les proportions d'individus dans les classes  $c_r$  et  $c_s$  dans le segment<sup>1</sup>  $t$ .

Un segment est *pur* s'il ne contient que des individus d'une seule classe, dans un tel cas :  $i(t) = 0$ . Plus le mélange des classes dans le segment  $t$  est important, plus l'impureté  $i(t)$  est élevée.

Chaque division  $d_j^m$  du nœud  $t$  par la variable  $x_j$  entraîne une réduction de l'impureté qui s'exprime par :

$$\Delta_j^m = i(t) - p_g i(t_g) - p_d i(t_d)$$

où  $p_g$  et  $p_d$  sont les proportions d'individus du nœud  $t$  respectivement dans les segments descendants  $t_g$  et  $t_d$  (la fonction  $i(t)$  étant concave, l'impureté moyenne ne peut que décroître par division d'un nœud).

Par conséquent pour chaque variable  $x_j$ , la meilleure division  $d_j^*$  est telle que la réduction de l'impureté  $\Delta_j^*$  est maximale :

$$\Delta_j^* = \max_{m \in d_j} \{\Delta_j^m\}$$

où  $d_j$  est l'ensemble des divisions de la variable  $x_j$ .

Sur l'ensemble des  $p$  variables, la division du nœud  $t$  est effectuée à l'aide de la variable qui assure :

$$\Delta^* = \max_{j=1, \dots, p} \{\Delta_j^*\}$$

- *Les étapes de l'algorithme*

Considérons maintenant 300 individus répartis en 3 classes  $c_1, c_2, c_3$  de même taille et sur lesquels 10 mesures quantitatives ont été relevées.

On procède comme dans le cas de la régression par segmentation en examinant toutes les variables.

<sup>1</sup> La fonction  $i(t)$  est l'indice de diversité de Gini (cf. Goodman et Kruskal, 1954). On aurait pu également utiliser l'entropie de Shannon :  $i(t) = - \sum_r^k P(r|t) \log P(s|t)$ .

Pour la variable  $x_1$ , on aboutit par exemple à la meilleure division (qui n'est pas nécessairement la plus discriminante) observable sur la figure 3.5 - 7.

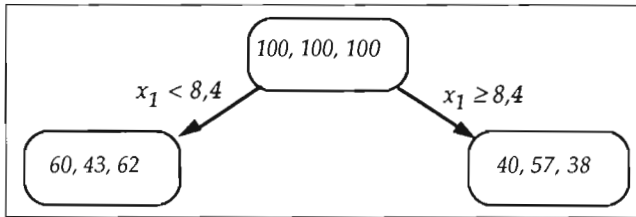


Figure 3.5 - 7  
Discrimination : meilleure division pour la variable  $x_1$

On retient finalement, parmi toutes les variables, celle qui produit la meilleure "meilleure division", par exemple la variable continue  $x_8$  pour  $\alpha = 3,5$ .

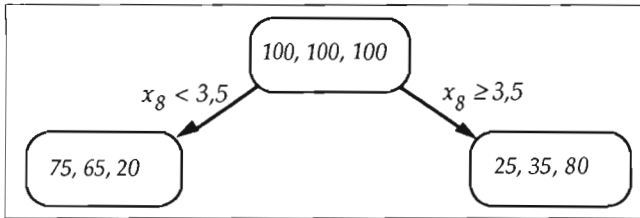


Figure 3.5 - 8  
Discrimination : meilleure division pour toutes les variables

On obtient ainsi la meilleure séparation entre les 3 classes, ce qui se traduit par le schéma de la figure 3.5 - 8. On applique cette même procédure aux deux segments descendants obtenus.

#### - Règle d'affectation

Si on considère le segment terminal  $t$  de taille  $n_t$ , il contient  $n_1(t)$  sujets appartenant à la classe 1, ...,  $n_r(t)$  sujets de la classe  $r$ , ...,  $n_k(t)$  sujets de la classe  $k$ .

Chaque segment terminal est affecté à la classe qui  $y$  est la mieux représentée. Par exemple, les segments 1 et 4 de la figure 3.5 - 9 sont affectés à la classe 2. Un nouvel individu qui *descend* dans l'arbre arrive dans un segment terminal et sera affecté à la classe correspondante.

#### - Taux d'Erreur Apparente de classement

A tout segment terminal  $t$  de l'arbre  $A$  associé à une classe  $c_s$  correspond une erreur de classement de la forme :

$$R(s|t) = \sum_{r=1}^k P(r|t)$$

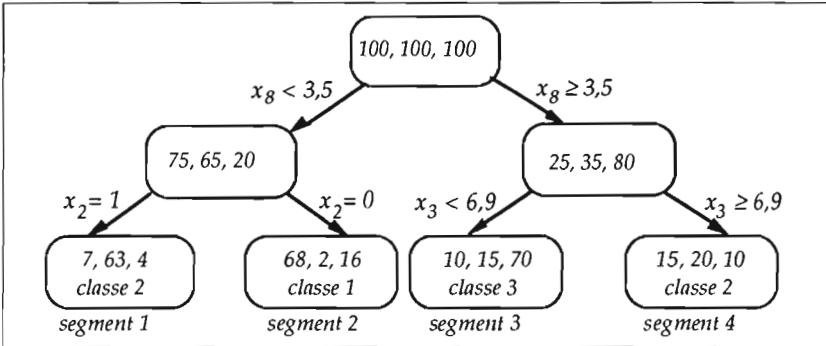


Figure 3.5 - 9  
Discrimination : Arbre à deux niveaux

avec  $r \neq s$  et où  $P(r|t) = \frac{n_r(t)}{n_t}$  est la proportion d'individus du segment  $t$  affectés à la classe  $c_s$  et qui appartiennent à la classe  $c_r$ .

Le Taux d'Erreur Apparent de classement (TEA) associé à l'arbre vaut :

$$TEA(A) = \sum_{t \in A} \frac{n_t R(s|t)}{n} = \sum_{t \in A} \sum_{r=1}^k \frac{n_r(t)}{n} \quad [3.5 - 2]$$

avec  $r \neq s$ . Il représente la proportion d'individus mal classés dans l'ensemble des segments terminaux.

Ainsi, l'arbre de la figure 3.5 - 9 ne fournit pas une bonne règle de décision en terme d'erreur de classement. En effet, un sujet qui parcourt l'arbre et qui tombe dans le segment 1 est affecté à la classe 2 avec une erreur de classement de 14,9 %; celui qui tombe dans le segment 4 est affecté à la classe 2 avec une erreur de classement de 55,5 %.

Le Taux d'Erreur Apparent de classement associé à l'arbre est la moyenne des erreurs de classement dans les différents segments terminaux, soit :

$$TEA = \frac{(74 \times 14,9\% + 86 \times 20,9\% + 95 \times 26,3\% + 45 \times 55,5\%)}{300} = 26,3\%$$

On a sans doute intérêt à continuer à diviser les segments. La question est de savoir à quel moment il faut arrêter la procédure de division.

### 3.5.3 Sélection du "meilleur sous-arbre"

Par "meilleur" sous-arbre, on entend un arbre qui contient le moins de segments terminaux et dont l'erreur apparente de prévision ou de classement est la plus petite possible, tout en fournissant une estimation correcte de l'erreur théorique.

Un sous-arbre ayant peu de segments terminaux entraîne une erreur apparente qui, bien que reflétant l'erreur théorique, est trop importante. En effet, si l'arbre est trop petit, on peut être conduit à perdre de bonnes divisions et à ne pas utiliser toute l'information contenue dans l'échantillon. Inversement, à un arbre trop grand (avec de nombreuses divisions) est associée une erreur apparente faible mais qui donne une estimation trop optimiste de l'erreur théorique. C'est donc entre ces deux extrêmes que doit être choisi le "meilleur" sous-arbre.

La méthode proposée par Breiman *et al.* (*op. cit.*) est fondée sur l'utilisation d'un échantillon-test et présente un double avantage :

- déterminer le "meilleur" sous-arbre sans employer de tests statistiques pour définir une règle d'arrêt de la procédure de division,
- obtenir une estimation précise de l'erreur théorique de prévision ou de classement.

### a – Procédures de sélection

Il est nécessaire de diviser l'échantillon de base en deux parties, l'échantillon d'apprentissage (par exemple constitué par les 2/3 de l'échantillon de base) et l'échantillon-test (le tiers restant). La recherche du "meilleur" sous-arbre  $A^*$  se fait alors de la façon suivante :

- A partir de l'échantillon d'apprentissage, on construit l'arbre complet  $A_{\max}$  ou un arbre tel que chaque segment terminal contienne peu d'individus.

Puis l'opération d'élagage de l'arbre  $A_{\max}$  consiste à construire une séquence optimale de sous-arbres emboîtés  $\{A_H, \dots, A_h, \dots, A_1\}$  où  $A_H$  coïncide avec  $A_{\max}$ ,  $A_h$  est le sous-arbre ayant  $h$  segments terminaux et  $A_1$  est l'échantillon total. Chaque sous-arbre  $A_h$  de cette séquence est optimal au sens suivant : son Erreur Apparente (EA) est minimale parmi les sous-arbres ayant le même nombre de segments terminaux<sup>1</sup>.

Si  $S_h$  est l'ensemble des sous-arbres de  $A_{\max}$  ayant  $h$  segments terminaux alors :

$$EA(A_h) = \min_{A \in S_h} \{EA(A)\}$$

- A partir de l'échantillon-test, on sélectionne, parmi les sous-arbres de la séquence optimale, le meilleur sous-arbre  $A^*$ . C'est celui qui présente la plus petite erreur théorique (ET) :

$$ET(A^*) = \min_{1 \leq h \leq H} \{ET(A_h)\}$$

<sup>1</sup> En fait, des algorithmes appropriés permettent de choisir une séquence sous-optimale, mais accessible par le calcul (cf. Breiman *et al.*, 1984; Celeux et Lechevallier in : Celeux, 1990).



Les individus de l'échantillon-test parcourent chacun des sous-arbres de la séquence optimale et tombent dans un segment terminal, ce qui entraîne une estimation de l'erreur théorique pour chaque sous-arbre.

En pratique, l'estimation de l'erreur théorique décroît rapidement à mesure que le nombre de segments terminaux des sous-arbres augmente, puis elle passe par un palier et croît ensuite lentement. Le sous-arbre  $A^*$  sélectionné comme optimal est le plus petit sous-arbre associé à l'estimation la plus petite de l'erreur théorique.

### b – Estimation de l'Erreur Théorique de Prévision

L'estimation de l'*Erreur Théorique de Prévision* pour un sous-arbre A de la séquence optimale,  $E\hat{T}P(A)$ , est calculée sur l'échantillon-test suivant la formule utilisée pour l'*Erreur Apparente de Prévision* [3.5 - 1] :

$$E\hat{T}P(A) = \sum_{t \in A} \tilde{R}_t$$

avec  $\tilde{R}_t = \frac{\tilde{n}_t}{\tilde{n}} \times \tilde{s}_t^2$  et où  $\tilde{n}$  est la taille de l'échantillon-test,  $\tilde{n}_t$  est le nombre d'individus de l'échantillon-test qui appartiennent au segment t et  $\tilde{s}_t^2$  est la variance de la variable y à l'intérieur du segment t.

### c – Estimation du Taux d'Erreur Théorique de classement

Les appellations de *Taux d'Erreur Apparent* ou *Théorique de Classement* n'ont de sens que dans le cas le plus simple c'est-à-dire si les probabilités *a priori* des classes sont estimées par les fréquences des classes dans l'échantillon et si les coûts de mauvaise classification sont tous égaux. Dans le cas général, on utilise un *Coût d'Erreur Apparent* ou *Théorique* pour lesquels les formules de calcul sont plus complexes.

#### - Cas le plus simple

L'estimation du *Taux d'Erreur Théorique de classement* se calcule comme le *Taux d'Erreur Apparent* [3.5 - 2] à partir de l'échantillon-test. Elle est égale à la proportion  $\tilde{p}_t$  d'individus mal classés par le sous-arbre A dans l'échantillon-test (cf. formule [3.5 - 2]).

$$T\hat{E}A(A) = \sum_{t \in A} \sum_{r=1}^k \frac{\tilde{n}_r(t)}{\tilde{n}} = \tilde{p}_t$$

avec  $r \neq s$ , où  $\tilde{n}$  est l'effectif de l'échantillon-test et  $\tilde{n}_r(t)$  est le nombre d'individus de l'échantillon-test affectés à la classe  $c_s$  et qui appartiennent à la classe  $c_r$  dans le segment terminal t.

Il est possible de fournir un intervalle de confiance associé à cette proportion  $\tilde{p}_t$  à partir de l'estimation de la variance de cette proportion :

$$\widehat{\text{Var}}(\tilde{p}_t) = \frac{\tilde{p}_t(1-\tilde{p}_t)}{\tilde{n}}$$

### - Cas général

La règle de décision la plus générale est celle qui tient compte des probabilités *a priori*  $\pi_r$  ( $r = 1, 2, \dots, k$ ) des  $k$  classes à discriminer et des *coûts de mauvais classement* notés  $C(r|s)$  où  $r \neq s = 1, 2, \dots, k$ .

$C(r|s)$  désigne le coût<sup>1</sup> entraîné par l'affectation d'un individu à la classe  $c_s$  alors qu'il appartient à la classe  $c_r$ . La règle générale d'affectation d'un segment terminal  $t$  à une classe est fondée sur le *coût moyen d'erreur de classement* (appelé aussi *risque d'erreur*).

Si  $n_r(t)$  désigne le nombre d'individus de la classe  $c_r$  du segment  $t$  et  $n_r$  l'effectif total de la classe  $c_r$ , on a :

$$P(r|t) = \frac{p_r \frac{n_r(t)}{n_r}}{P(t)}$$

où  $P(t) = \sum_{r=1}^k p_r \frac{n_r(t)}{n_r}$  est la probabilité d'"aboutir" au segment  $t$ .

Le coût moyen d'erreur de classement  $R(s|t)$  entraîné par l'affectation du segment  $t$  à la classe  $c_s$  est égal à :

$$R(s|t) = \sum_{r=1}^k C(r|s)P(r|t)$$

Ainsi le segment terminal  $t$  est affecté à la classe  $c_j$  si :

$$R(j|t) = \min_{s=1, \dots, k} \{R(s|t)\}$$

### Remarque

Si la probabilité  $\pi_r$  d'appartenance *a priori* à la classe  $c_r$  est égale à la proportion d'individus de cette classe dans l'échantillon :

$$\pi_r = \frac{n_r}{n}$$

alors  $P(t)$  tel que :

$$P(t) = \sum_{r=1}^k p_r \frac{n_r(t)}{n_r}$$

est simplement la proportion d'individus composant le segment terminal  $t$ .

<sup>1</sup> Les différents coûts  $C(s|s)$  sont nuls et en général  $C(r|s) \neq C(s|r)$ .

### 3.5.4 Divisions équi-réductrices et équi-divisantes

La *meilleure* division  $d^*$  d'un nœud est celle qui assure la plus grande réduction de la variance résiduelle ou de l'impureté en passant du nœud à ses segments descendants. Cette notion de maximum absolu est très stricte. Il peut exister en effet des divisions presque aussi bonnes, pouvant jouer un rôle important au niveau des interprétations.

Par extension, on définit, à côté de  $d^*$ , deux autres types de divisions :

- les divisions *équi-réductrices* qui assurent, après  $d^*$ , les plus fortes réductions de l'impureté ou les plus faibles variances résiduelles. Elles permettent d'intervenir sur le choix de la "meilleure" variable explicative.
- les divisions *équi-divisantes* qui fournissent les répartitions les plus proches de la meilleure division  $d^*$ . Elles permettent de gérer l'existence de données manquantes dans l'affectation d'un nouvel individu à une classe ou à une valeur de  $y$ .

#### a – Divisions équi-réductrices

La procédure de division d'un nœud fournit les premières *meilleures divisions* d'un nœud pour lesquelles la réduction de la variance résiduelle ou de l'impureté  $\Delta_j^*$  est élevée (cf. 3.5.2.b et c).

Si la meilleure division  $d^*$  du nœud  $t$  est obtenue à partir de la variable  $x^*$ , on définit la première division équi-réductrice  $d_i^*$  effectuée sur la variable  $x_i$  ( $x_i \neq x^*$ ) avec  $i = 1, \dots, p$ . C'est celle qui correspond à une réduction des segments descendants la plus proche de celle de la *meilleure division*  $d^*$ .

En d'autres termes, c'est la deuxième meilleure division du nœud  $t$ . On définit par extension les 2ème, 3ème, ..., divisions équi-réductrices<sup>1</sup>.

#### b – Divisions équi-divisantes

Les divisions équi-divisantes<sup>2</sup> permettent de classer un nouvel individu présentant une donnée manquante pour la variable définissant la division. L'idée est la suivante : on cherche une variable qui remplace au mieux la variable divisant le nœud, c'est-à-dire qui assure presque la même séparation des individus. De la même manière, on peut définir la seconde, troisième, ..., meilleure division équi-divisante.

<sup>1</sup> Les divisions équi-réductrices sont parfois appelées concurrentes. Il est possible ainsi d'intervenir sur le choix des variables associées aux "meilleures" divisions suivant la perception personnelle qu'a l'utilisateur du problème. En effet, à la variable produisant la "meilleure" division, on peut préférer une autre variable que l'on sait plus pertinente pour l'étude.

<sup>2</sup> Les divisions équi-divisantes sont parfois appelées suppléantes.

Ainsi, si la valeur de  $x_j$  est manquante pour un nouvel individu, on l'affectera à un des segments descendants en utilisant la meilleure division équi-divisante de  $d^*$ . Si la valeur de la variable associée à la meilleure division équi-divisante est manquante, on aura recours à la deuxième meilleure division équi-divisante, etc.

### 3.5.5 Lien avec les méthodes de classement

Segmentation, discrimination, classement, classification ou classification supervisée, régression linéaire multiple, régression logistique, régression pas-à-pas, ..., le vocabulaire ne manque pas pour désigner, suivant le domaine d'application, des opérations qui sont souvent proches, sinon identiques. On veut, dans ce paragraphe, brièvement situer la segmentation parmi les outils répondant à des préoccupations voisines de la part de l'utilisateur.

La segmentation, bien que travaillant par divisions de l'échantillon en classes, est plus proche des techniques de régression pas à pas (qu'il s'agisse de régression linéaire multiple ou de régression logistique) et de discrimination pas-à-pas que des méthodes de classification automatique. En effet, il ne s'agit pas de faire apparaître des classes, mais de chercher les groupes d'individus les plus "explicatifs" des modalités d'une variable qualitative particulière (ou des valeurs d'une variable continue). Le principe est, on l'a vu, de chercher la dichotomie (induite à chaque pas par une des variables) la plus "liée" à la variable privilégiée.

La segmentation n'est pas vraiment multidimensionnelle au sens géométrique du terme (on ne calcule pas de distances dans  $\mathbb{R}^p$  ni dans  $\mathbb{R}^n$  comme pour les méthodes factorielles ou de classification), mais on utilise les variables explicatives conditionnellement les unes par rapport aux autres. On peut donc parfois atteindre des effets d'interaction assez difficiles à saisir par d'autres méthodes, sans prétendre d'ailleurs les atteindre tous.

La parenté avec les méthodes descriptives reste forte, dans la mesure où les aspects "contrôle des opérations par l'utilisateur", "transparence du fonctionnement", voire "ergonomie des résultats" occupent une position de premier plan. L'arbre de décision binaire est lisible par tout utilisateur. Autre avantage déjà évoqué dans l'introduction de cette section, la mixité des variables qu'accepte la procédure : nominales, ordinales, continues peuvent être mélangées au niveau des variables explicatives, et peuvent constituer la variable à expliquer.

La validation par une méthode de rééchantillonnage (limité aux échantillons-test dans l'exposé qui précède) est elle-même une des techniques de validation les plus transparentes pour l'utilisateur.

Pour conclure, on doit cependant reconnaître quelques défauts à la segmentation par arbre binaire, qui rendent son *utilisation exclusive* insuffisante.

L'aspect séquentiel est redoutable, car les covariations qui servent à sélectionner les variables ne mesurent pas un lien causal et une variable peut en cacher une autre, beaucoup plus fondamentale, qui n'a plus aucune chance d'apparaître dans la suite du processus. Les divisions de réserve (équi-réductrices et équi-divisantes) sont là pour pallier partiellement cet inconvénient. Mais l'arbre binaire perd alors une partie de sa séduisante simplicité.

L'absence de visualisation globale, propice à une réflexion critique sur le recueil de données et à une observation simultanée de l'ensemble des covariations, est également une faiblesse par rapport aux méthodes factorielles.

Enfin, il se peut que la nature du phénomène étudié fasse que des combinaisons linéaires (après éventuel recodage) soient optimales pour prévoir la variable étudiée (ou son *logit* ou toute autre fonction). Dans ce cas, la segmentation progressive sera surclassée.

Ces quelques critiques ne portent cependant que sur l'usage exclusif de la segmentation par arbre binaire. Une démarche impliquant plusieurs points de vue (visualisation préalable des variables explicatives avec positionnement *a posteriori* de la variable à expliquer, régression ou discrimination) permet d'éviter la plupart des écueils mentionnés.

## Analyses partielles et projetées

Ces méthodes se proposent d'analyser les associations existant entre des variables et des individus, non seulement après élimination d'effets de niveaux ou d'échelle, mais également après avoir tenu compte de l'influence éventuelle de "variables exogènes".

A l'origine et au centre de ces techniques se trouve l'analyse en composantes principales partielle ou sur variables instrumentales selon la terminologie de Rao (1964).

### 3.6.1 Définition du coefficient de corrélation partielle

Deux variables aléatoires  $X_1$  et  $X_2$  sont supposées dépendre d'une même variable aléatoire  $Z$ . On dispose d'un échantillon de chacune de ces variables. On peut mesurer directement le coefficient de corrélation  $r(x_1, x_2)$  sur deux échantillons de taille  $n$  représentés dans  $\mathbb{R}^n$  par les vecteurs à  $n$  composantes  $x_1$  et  $x_2$ . Mais nous voulons en fait connaître la liaison existant entre  $x_1$  et  $x_2$  en éliminant l'effet de la variable  $Z$  dont les  $n$  observations sont les composantes du vecteur  $z$ .

Pour prendre un exemple élémentaire classique<sup>1</sup>,  $X_1$  est la plus grande dimension d'un œuf,  $X_2$  la plus petite et  $Z$  son poids. Sur un échantillon de  $n = 100$  œufs, on trouvera un coefficient  $r(x_1, x_2)$  fortement positif, car il existe de gros œufs, pour lesquels  $X_1$  et  $X_2$  ont des valeurs élevées, et des petits, pour lesquels ces valeurs sont faibles. Par contre, si le poids  $Z$  est fixé, la liaison observée sera inverse car, à poids égal, les œufs sont plus ou moins sphériques.

Comment mesurer cette liaison entre  $X_1$  et  $X_2$  "à  $Z$  constant"? Une première méthode consiste à regrouper les observations en classes à l'intérieur desquelles les valeurs de  $Z$  sont peu différentes. On calcule alors le coefficient de corrélation entre  $X_1$  et  $X_2$  dans chaque classe et l'on fait, par exemple, une moyenne pondérée de ces coefficients, pour avoir une idée d'ensemble de la liaison. Cette méthode est excellente et il est conseillé de l'employer chaque fois que la taille des échantillons permet une division en classes d'effectifs suffisants.

Une autre méthode va nous permettre de calculer la liaison entre  $X_1$  et  $X_2$  "à  $Z$  constant" de façon simple, même lorsque les échantillons sont petits

---

<sup>1</sup> Cf. Darrois (1957).

(mais au prix d'une hypothèse sur la linéarité des liaisons). Ce coefficient de corrélation entre  $X_1$  et  $X_2$  "à  $Z$  constant" s'appellera le coefficient de *corrélation partielle* entre  $X_1$  et  $X_2$ , et on le notera  $\rho(X_1, X_2|Z)$ . Son calcul repose sur l'hypothèse que l'effet de  $Z$  sur les variables  $X_1$  et  $X_2$  se manifeste par des relations du type<sup>1</sup> :

$$\begin{cases} X_1 = \alpha_1 Z + \varepsilon_1 \\ X_2 = \alpha_2 Z + \varepsilon_2 \end{cases}$$

Une fois ôtée l'influence de la variable  $Z$ , les variables aléatoires  $X_1$  et  $X_2$  deviennent  $X_1 - \alpha_1 Z = \varepsilon_1$  et  $X_2 - \alpha_2 Z = \varepsilon_2$ . Le coefficient de *corrélation partielle* théorique  $\rho(X_1, X_2|Z)$  est par définition le coefficient de corrélation usuel entre  $\varepsilon_1$  et  $\varepsilon_2$  :

$$\rho(X_1, X_2|Z) = \frac{\text{cov}(\varepsilon_1, \varepsilon_2)}{\sqrt{\text{var}(\varepsilon_1)\text{var}(\varepsilon_2)}}$$

On définit de façon analogue une matrice des covariances partielles  $\mathbf{V}(X|Z)$  et une matrice des corrélations partielles  $\mathbf{C}(X|Z)$  entre  $p$  variables  $X_1, X_2, \dots, X_p$ , lorsque  $q$  variables  $Z_1, Z_2, \dots, Z_q$  sont supposées fixées. On a alors le système suivant :

$$\begin{cases} X_1 = \alpha_{11}Z_1 + \alpha_{12}Z_2 + \dots + \alpha_{1q}Z_q + \varepsilon_1 \\ X_2 = \alpha_{21}Z_1 + \alpha_{22}Z_2 + \dots + \alpha_{2q}Z_q + \varepsilon_2 \\ \dots \\ X_p = \alpha_{p1}Z_1 + \alpha_{p2}Z_2 + \dots + \alpha_{pq}Z_q + \varepsilon_p \end{cases}$$

$\mathbf{V}(X|Z)$  et  $\mathbf{C}(X|Z)$  sont respectivement les matrices des covariances et des corrélations théoriques entre les variables résiduelles :  $\varepsilon_1, \varepsilon_2, \dots, \varepsilon_p$ .

### 3.6.2 Calcul des covariances et corrélations partielles

#### a – Cas de deux variables

Pour les  $n$  observations des trois variables  $X_1, X_2, Z$ , qui sont les composantes, supposées ici centrées, des 3 vecteurs  $\mathbf{x}_1, \mathbf{x}_2, \mathbf{z}$ , ces relations d'ajustement s'écrivent, avec les notations de la section 3.2 (mais la lettre  $x$  désigne maintenant des variables endogènes ou expliquées) :

$$\begin{cases} \mathbf{x}_1 = a_1 \mathbf{z} + e_1 \\ \mathbf{x}_2 = a_2 \mathbf{z} + e_2 \end{cases}$$

<sup>1</sup> Comme pour tout modèle linéaire, les variables entre lesquelles existe une relation linéaire peuvent être des variables transformées construites à partir des variables réellement observées. Le caractère linéaire de la relation n'est donc pas une contrainte importante.

où  $a_1$  et  $a_2$  sont respectivement les estimations de  $\alpha_1$  et  $\alpha_2$  par la méthode des moindres-carrés alors que  $e_1$  et  $e_2$  sont les résidus observés. La covariance partielle expérimentale s'écrit :

$$\text{Cov}(x_1, x_2 | z) = \frac{1}{n} \mathbf{e}'_1 \mathbf{e}_2 = \frac{1}{n} (\mathbf{x}_1 - a_1 \mathbf{z})' (\mathbf{x}_2 - a_2 \mathbf{z})$$

soit :

$$\text{Cov}(x_1, x_2 | z) = \frac{1}{n} \{ \mathbf{x}'_1 \mathbf{x}_2 - a_1 \mathbf{z}' \mathbf{x}_2 - a_2 \mathbf{x}'_1 \mathbf{z} + a_1 a_2 \mathbf{z}' \mathbf{z} \}$$

On remplace les coefficients de régression par leur valeur  $a_1 = \mathbf{x}_1 \mathbf{z}' / \mathbf{z}' \mathbf{z}$  et  $a_2 = \mathbf{x}_2 \mathbf{z}' / \mathbf{z}' \mathbf{z}$  et l'on obtient après simplification :

$$\text{Cov}(x_1, x_2 | z) = \frac{1}{n} \left\{ \mathbf{x}'_1 \mathbf{x}_2 - \frac{(\mathbf{x}'_1 \mathbf{z})(\mathbf{x}'_2 \mathbf{z})}{\mathbf{z}' \mathbf{z}} \right\}$$

expression que l'on peut écrire :

$$\text{Cov}(x_1, x_2 | z) = \text{Cov}(x_1, x_2) - \frac{\text{Cov}(x_1, z) \text{Cov}(x_2, z)}{\text{Var}(z)} \quad [3.6-1]$$

Les variances résiduelles se calculent de façon analogue et l'on a pour  $\mathbf{e}_1$  :

$$\frac{1}{n} \mathbf{e}'_1 \mathbf{e}_1 = \text{Var}(x_1) - \frac{\text{Var}^2(x_1, z)}{\text{Var}(z)} = (1 - r^2(x_1, z)) \text{Var}(x_1)$$

Le coefficient de corrélation partielle  $r(x_1, x_2 | z)$  s'écrit alors, en faisant apparaître les coefficients de corrélation usuels :

$$r(x_1, x_2 | z) = \frac{r(x_1, x_2) - r(x_1, z)r(x_2, z)}{\sqrt{(1 - r^2(x_1, z))(1 - r^2(x_2, z))}}$$

### b – Cas de $p$ variables (X) et de $q$ variables (Z)

Nous disposons maintenant de  $p$  vecteurs  $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_p$  auxquels correspondent  $p$  points dans  $\mathbb{R}^n$ . On peut mesurer la covariance (ou la corrélation) entre ces variables après élimination de l'effet de  $q$  autres variables représentées dans  $\mathbb{R}^n$  par les vecteurs  $\mathbf{z}_1, \mathbf{z}_2, \dots, \mathbf{z}_q$ . On désignera par  $\mathbf{X}$  la matrice  $(n, p)$  et par  $\mathbf{Z}$  la matrice  $(n, q)$  qui rassemblent en colonne ces divers vecteurs.

Pour la  $k^{\text{ième}}$  variable, l'ajustement des moindres-carrés entre  $\mathbf{x}_k$  et les variables exogènes  $\mathbf{z}_1, \mathbf{z}_2, \dots, \mathbf{z}_q$  s'écrit :

$$\mathbf{x}_k = a_{k1} \mathbf{z}_1 + a_{k2} \mathbf{z}_2 + \dots + a_{kq} \mathbf{z}_q + \mathbf{e}_k$$

où  $\mathbf{e}_k$  est le vecteur résiduel. Nous appellerons  $\mathbf{a}_k$  le vecteur-colonne de ces  $q$  coefficients. Après avoir effectué les  $p$  ajustements similaires concernant  $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_p$  nous rassemblons dans la matrice  $\mathbf{A}$  de dimension  $(q, p)$  les  $p$  vecteurs-colonnes de coefficients  $(\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_p)$  et dans la matrice  $\mathbf{E}$  de dimension  $(n, p)$  les  $p$  vecteurs résiduels  $(\mathbf{e}_1, \mathbf{e}_2, \dots, \mathbf{e}_p)$ . Le système des ajustements s'écrit alors de façon synthétique :



$$\underset{(n,p)}{X} = \underset{(n,q)}{Z} \underset{(q,p)}{A} + \underset{(n,p)}{E}$$

Dans la matrice  $A$ , la  $k^{\text{ième}}$  colonne est :

$$a_k = (Z'Z)^{-1}Z'x_k.$$

Il est donc possible d'écrire  $A$  sous la forme :

$$A = (Z'Z)^{-1}Z'X \quad [3.6 - 2]$$

Avec ces notations, la matrice  $(p, p)$  qui définit les covariances partielles expérimentales sur les  $X$  "à  $Z$  constant" s'écrira :

$$\begin{aligned} V(X|Z) &= \frac{1}{n} E'E = \frac{1}{n} (X - ZA)'(X - ZA) \\ &= \frac{1}{n} (X'X - A'Z'X - X'ZA + A'Z'ZA) \end{aligned}$$

En remplaçant  $A$  par son expression [3.6 - 2] et après simplification :

$$V(X|Z) = \frac{1}{n} \{X'X - X'Z'(Z'Z)^{-1}Z'X\} \quad [3.6 - 3]$$

Imaginons que soient rassemblés dans un tableau  $T$  à  $n$  lignes et  $p + q$  colonnes les tableaux centrés  $X$  et  $Z$  :

$$T = [X, Z]$$

Alors la matrice des covariances entre les colonnes de  $T$  peut être partitionnée en quatre sous-matrices de covariances :

$$V(T) = \begin{bmatrix} V_{XX} & V_{ZX} \\ V_{XZ} & V_{ZZ} \end{bmatrix}$$

avec :

$$V_{XX} = \frac{1}{n} X'X ; \quad V_{ZZ} = \frac{1}{n} Z'Z ; \quad V_{ZX} = V'_{XZ} = \frac{1}{n} Z'X.$$

Alors la relation [3.6 - 3] prend la forme :

$$V(X|Z) = V_{XX} - V_{ZX} V_{ZZ}^{-1} V_{XZ} \quad [3.6 - 4]$$

où elle apparaît comme une généralisation, pour  $q \geq 1$ , de [3.6 - 1].

La matrice des *corrélations partielles* se calcule aisément à partir de la matrice des covariances partielles  $V(X|Z)$  comme une matrice des corrélations ordinaires à partir d'une matrice des covariances.

### 3.6.3 Analyse du nuage résiduel ou analyse partielle

L'analyse du tableau  $X$  lorsque les variables  $z_1, z_2, \dots, z_q$  sont fixées, se ramène donc à l'analyse générale (cf. section 1.1) du tableau des écarts  $E$ . Ainsi les points représentant les variables dans  $\mathbb{R}^n$  auront pour

coordonnées (à une homothétie près), sur l'axe factoriel  $\alpha$ , les composantes du  $\alpha^{\text{ième}}$  vecteur-propre  $u_\alpha$  de la *matrice des covariances-partielles*<sup>1</sup> :

$$V(X|Z) = \frac{1}{n} E'E$$

Poursuivant l'interprétation géométrique de l'ajustement des moindres-carrés, on peut remarquer que :

$$nV(X|Z) = X'(I - Z(Z'Z)^{-1}Z')X = X'(I - P_Z)X = X'P_Z^*X$$

où  $P_Z^*$  ( $= I - P_Z$ ) est une matrice ( $n, n$ ) symétrique et idempotente, analogue à la matrice  $P_X$  définie par la formule [3.2 - 3] au paragraphe 3.2.2.b. Ici  $P_Z^*$  effectue la projection de tout vecteur de  $\mathbb{R}^n$  sur le sous-espace à  $(n - q - 1)$  dimensions, orthogonal au sous-espace engendré par  $(z_1, z_2, \dots, z_q)$ . C'est cette projection que l'on analyse lorsqu'on opère la transformation des données  $E = P_Z^* X$ .

Ainsi, dans l'hypothèse où les régressions traduisent effectivement l'effet des variables que l'on désire fixer, il est possible d'étudier *a posteriori* les liaisons et les associations existant entre des variables et des observations, "toutes choses égales par ailleurs".

Dans certains cas, on peut au contraire (cf. paragraphe 3.6.4 ci-après) être intéressé par la projection du nuage sur le sous-espace engendré par  $Z$ , le tableau analysé étant alors le tableau  $F = P_Z X$ . On réservera le nom d'*analyse projetée* à l'analyse de  $F$ .

### 3.6.4 Autres analyses partielles ou projetées

Il existe plusieurs variantes de méthodes impliquant des projections sur des sous-espaces. Une vue générale ainsi que des extensions de ce type d'approche sont données par Sabatier (1984, 1987).

On a vu que l'analyse canonique (section 3.1) part d'une situation analogue, c'est-à-dire d'un tableau de la forme  $R = (X, Z)$ , mais cherche le plus petit angle entre les sous-espaces engendrés par les colonnes de  $X$  et de  $Z$  dans  $\mathbb{R}^n$ . Ceci a conduit à diagonaliser une matrice du type :

$$S = (X'X)^{-1}X'Z(Z'Z)^{-1}Z'X = (X'X)^{-1}X'P_ZX$$

où  $X'P_ZX = (P_ZX)'P_ZX$  est proportionnel à la matrice d'inertie du nuage projeté sur le sous-espace engendré par les colonnes de  $Z$ .

<sup>1</sup> Pour une analyse *normée*, on utiliserait la matrice des *corrélations partielles*.

Dans l'équation  $\mathbf{S}\mathbf{u} = \lambda\mathbf{u}$ , posons  $\mathbf{u} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{v}$ . On obtient :

$$\mathbf{X}'\mathbf{P}_Z\mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{v} = \lambda\mathbf{v}$$

ce qui montre que  $\mathbf{v}$  est bien un axe principal de l'analyse du nuage projeté<sup>1</sup> avec la métrique  $(\mathbf{X}'\mathbf{X})^{-1}$  (cf. § 1.1.6).

On a également vu que l'analyse discriminante est un cas particulier de ce type d'analyse lorsque  $\mathbf{Z}$  est le tableau de codage disjonctif d'une variable nominale.

### a – Analyse canonique des correspondances

Une technique voisine, qui aurait pu avoir sa place dans les sections consacrées à l'analyse canonique ou à l'analyse discriminante, est l'analyse canonique des correspondances, proposée par Ter Braak (1986, 1987), étudiée et appliquée par Chessel *et al.* (1987), Lebreton *et al.* (1988), et étendue par Ter Braak (1988) à l'analyse canonique partielle des correspondances.

On a fait allusion au début de ce chapitre aux dangers d'une prolifération indéfinie de méthodes spécifiques, en reconnaissant cependant que si des situations typiques ou des structures typiques de tableaux se présentent avec une certaine fréquence, il est loisible de forger des instruments *ad hoc*.

En écologie précisément, les observations se présentent souvent sous la forme d'un tableau  $\mathbf{R} = (\mathbf{X}, \mathbf{Z})$  où, pour  $n$  sites (lignes de  $\mathbf{X}$  et de  $\mathbf{Z}$ ), on dispose d'un tableau numérique  $\mathbf{X}$  (qui peut aussi être une autre table de contingence ou un tableau disjonctif complet) décrivant les sites (variables géologiques, climatiques, pétrochimiques, botaniques, etc.) et d'une table de contingence (ou parfois de présence-absence)  $\mathbf{Z}$  donnant le nombre ou la présence de  $q$  espèces animales ou végétales sur les  $n$  sites.

Si l'on appelle  $\mathbf{D}_n$  et  $\mathbf{D}_q$  les matrices diagonales d'ordres  $(n, n)$ , et  $(q, q)$  contenant les marges de la table  $\mathbf{Z}$ , on munira les  $n$  lignes de  $\mathbf{X}$  de masses proportionnelles à la diagonale de  $\mathbf{D}_n$  (en particulier pour centrer les  $p$  colonnes de  $\mathbf{X}$ ). On notera encore  $\mathbf{X}$  dans la suite la matrice centrée de cette façon. L'analyse canonique des correspondances revient à diagonaliser :

$$\mathbf{S} = (\mathbf{X}'\mathbf{D}_n\mathbf{X})^{-1}(\mathbf{X}'\mathbf{Z}\mathbf{D}_q^{-1})\mathbf{D}_q(\mathbf{D}_q^{-1}\mathbf{Z}'\mathbf{X})$$

Si le tableau  $\mathbf{Z}$  est un tableau disjonctif complet (une seule espèce et un seul spécimen par site),  $\mathbf{Z}'\mathbf{Z} = \mathbf{D}_q$  et la matrice  $\mathbf{D}_n$  est une matrice scalaire ; l'analyse canonique des correspondances est alors simplement l'analyse discriminante visant à prédire les espèces à partir des caractéristiques des sites<sup>2</sup>.

<sup>1</sup> On vérifie que  $\mathbf{v}$  est bien de norme 1 pour la métrique  $(\mathbf{X}'\mathbf{X})^{-1}$  puisque  $\mathbf{u}'\mathbf{X}'\mathbf{X}\mathbf{u} = 1$ .

<sup>2</sup> Comme le remarquent Lebreton *et al.* (1988), on peut se ramener aux calculs d'une analyse discriminante dans le cas général en multipliant les lignes de  $\mathbf{Z}$  de façon à ne laisser qu'un spécimen d'une seule espèce par ligne et en répétant de façon similaire les

La matrice  $A = X'ZD_q^{-1}$  d'ordre  $(p, q)$  contient les moyennes des variables par espèces.

Comme on vient de le voir à propos de l'analyse canonique, il s'agit ici d'une analyse en axes principaux de  $A$  dans la métrique définie par  $(X'D_n X)^{-1}$ , inverse de la matrice des covariances totales pondérées des variables-colonnes de  $X$ .

Réécrivons une matrice du type de  $S$  dans le cas où  $D_n$  est une matrice scalaire (nombre constant d'espèces par site) et en posant  $Y = ZD_q^{-1/2}$  :

$$S = (X'X)^{-1}(X'Y Y'X)$$

Remarquons que si le vecteur  $u$  est vecteur propre de  $S$  relatif à la valeur propre  $\lambda$ , alors :

$$v = Y'Xu$$

est vecteur propre de :

$$S_1 = Y'X(X'X)^{-1}X'Y = Y' P_X Y = (P_X Y)'(P_X Y)$$

relatif à la même valeur propre  $\lambda$ .

Or  $S_1$  correspond à l'analyse en axes principaux de la projection de la table de contingence (normalisée)  $Y$  sur le sous-espace engendré par les colonnes du tableau  $X$  dans l'espace  $\mathbb{R}^n$ .

L'analyse canonique des correspondances peut donc être considérée comme une analyse partielle particulière. Elle diffère de l'analyse canonique en ce sens qu'elle traite de façon dissymétrique les deux tableaux  $X$  et  $Z$  (elle ne fait jamais intervenir la matrice  $(Z'Z)^{-1}$ , c'est-à-dire finalement la structure interne du tableau  $Z$ , indépendamment de  $X$ ).

### **b — Analyse non-symétrique des correspondances**

On a vu plus haut que, en présence d'un tableau de données  $R = (X, Z)$ , comprenant deux groupes de variables, l'analyse canonique conduisait à diagonaliser la matrice :

$$S = (X'X)^{-1}X'Z(Z'Z)^{-1}Z'X$$

alors que l'analyse du nuage des lignes de  $X$  projeté sur le sous-espace engendré par les colonnes de  $Z$  conduit à diagonaliser :

$$S_1 = X'Z(Z'Z)^{-1}Z'X = X' P_Z X$$

---

lignes de  $X$ . Cette dilatation de  $Z$  supprime les cooccurrences d'espèces à l'intérieur d'un même site.

Si les matrices  $X$  et  $Z$  sont des tableaux disjonctifs complets, la diagonalisation de  $S$  est celle impliquée dans l'analyse des correspondances de la table de contingence  $C = X'Z$ .

La diagonalisation de  $S_1$  correspond (à un centrage près) à l'*analyse non-symétrique des correspondances* de cette même table  $C$ , introduite et développée par Lauro et D'Ambra (1984) pour traiter les situations où les variables lignes et colonnes jouent des rôles dissymétriques<sup>1</sup>.

---

<sup>1</sup> Cette méthode a connu des développements parallèles à ceux de l'analyse des correspondances : généralisations au cas multiple, liens avec les modèles log-linéaires, études de validation et de stabilité (pour une vue générale de ces travaux, cf. Balbi, 1994).

## Section 3.7

## Structures de graphe, analyses locales

La nature ou l'origine du recueil de données suggèrent souvent une structure *a priori* de l'ensemble des individus ou observations, avant toute analyse statistique.

On peut voir sur la figure 3.7 - 1 des représentations qui correspondent à trois structures distinctes de l'ensemble des observations. La structure de partition, qui correspond à un graphe formé de cliques disjointes, peut être décrite par une simple variable nominale, et entre donc dans le cas des analyses partielles présentées plus haut. Elle fera cependant l'objet d'un traitement particulier qui fait intervenir les matrices de covariances intra-classes et inter-classes, comme en analyse factorielle discriminante.

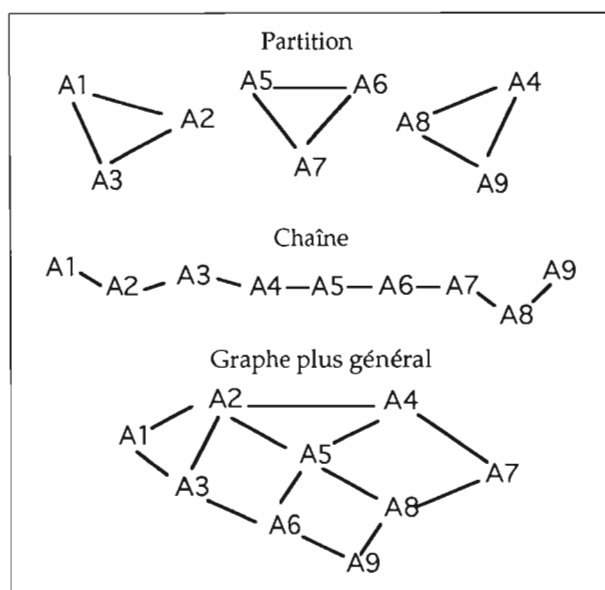


Figure 3.7 - 1

**Graphes correspondant à trois types de structures courantes :**  
**Partition, chaîne (chronologie), graphe non orienté.**

La structure de chaîne correspond le plus souvent à des observations consécutives dans le temps, alors que la structure plus générale de graphe non orienté peut schématiser un système d'observations géographiques, pour lequel il existe une certaine dépendance entre observations contiguës. Ces structures ne peuvent pas être prises en compte par des variables nominales car elles concernent des couples d'observations.

### 3.7.1 Variance locale et covariance locale d'une variable

La décomposition de la variance en *variance entre classes* et *variance dans les classes* n'est plus possible dans le cas d'une structure de graphe.

On peut faire intervenir une autre décomposition, fondée sur la propriété de la covariance empirique<sup>1</sup> entre deux variables  $x$  et  $y$  d'être également une covariance entre tous les couples d'observations :

$$\text{cov}(x, y) = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y}) = \frac{1}{2n(n-1)} \sum_{i=1}^n \sum_{i'=1}^n (x_i - x_{i'})(y_i - y_{i'})$$

On a évidemment l'expression de la variance si  $x = y$  :

$$\text{var}(x) = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2 = \frac{1}{2n(n-1)} \sum_{i=1}^n \sum_{i'=1}^n (x_i - x_{i'})^2 \quad [3.7-1]$$

S'il existe une relation binaire symétrique sur l'ensemble  $I$  des individus, définie par une partie symétrique  $R$  de l'ensemble produit  $I \times I$  ( $R$  sera l'ensemble des couples contigus  $(i, i')$ ), on peut écrire, dans le cas de la variance par exemple :

$$\text{var}(x) = \frac{1}{2n(n-1)} \left\{ \sum_{(i,i') \in R} (x_i - x_{i'})^2 + \sum_{(i,i') \notin R} (x_i - x_{i'})^2 \right\}$$

c'est-à-dire séparer, dans le dénominateur de la variance, les contributions des couples contigus (ou adjacents sur le graphe) et des autres couples.

#### Matrice de contiguïté

Un couple de sommets adjacents du graphe est relié par une *arête*. Le nombre des arêtes attachées à un même sommet  $i$  est appelé le *degré* de ce sommet. Ce nombre est noté  $m_i$ . Le nombre d'arêtes du graphe s'écrit alors :

$$m_a = \frac{1}{2} \sum_{i=1}^n m_i$$

Si tous les sommets sont reliés par une arête, le graphe est dit *complet*. Un tel graphe possède  $n(n-1)/2$  arêtes (on ne distingue pas l'arête  $(i, i')$  de l'arête  $(i', i)$ ). On construit une matrice carrée  $\mathbf{M}$ , d'ordre  $(n, n)$ , dite *matrice de contiguïté*, ou matrice associée au graphe telle que  $m_{ii'} = 1$  si  $i$  est voisin de  $i'$  et  $m_{ii'} = 0$  sinon<sup>2</sup>.

<sup>1</sup> La covariance empirique sera calculée ici en divisant la somme des produits par  $(n-1)$  (au lieu de  $n$ ). On obtient ainsi une estimation sans biais de la covariance théorique.

<sup>2</sup> On peut également travailler sur des structures de contiguïté qui incluent des proximités à distance  $1, 2, \dots, n$  les matrices de contiguïté correspondantes étant construites à partir des puissances booléennes de la matrice  $\mathbf{M}$  (cf. Lebart, 1969-a). Nous nous limiterons ici aux structures de contiguïté pour lesquelles deux parties sont immédiatement contiguës (distance 1) ou disjointes bien que les résultats soient également généralisables à des structures de contiguïté plus complexes.

Notons qu'avec les notations précédentes :

$$m_i = \sum_{i'=1}^n m_{ii'}$$

On voit que cette matrice est symétrique du fait de la symétrie de la relation de contiguïté. On adoptera la convention selon laquelle une observation n'est pas contiguë à elle-même, ce qui implique que les termes  $m_{ii}$  situés sur la diagonale principale de la matrice  $\mathbf{M}$  sont tous nuls. On peut réécrire, dans la dernière formule donnant la variance, le terme faisant intervenir les couples contigus sous la forme :

$$\sum_{(i,i') \in \mathbf{R}} (x_i - x_{i'})^2 = \sum_{i=1}^n \sum_{i'=1}^n m_{ii'} (x_i - x_{i'})^2$$

On appelle variance locale  $v_l(x)$  d'une variable  $x$  la demi-moyenne des carrés des accroissements correspondant à des observations contiguës. Posant :

$$m = \sum_{i=1}^n \sum_{i'=1}^n m_{ii'}$$

on a :

$$v_l(x) = \frac{1}{2m} \sum_{i=1}^n \sum_{i'=1}^n m_{ii'} (x_i - x_{i'})^2 \quad [3.7 - 2]$$

La variance totale  $var(x)$  donnée par la formule [3.7 - 1] est donc la variance locale correspondant au graphe complet.

### 3.7.2 Coefficient de contiguïté de Geary (1954)

Il est clair que si la variable  $x$  est indépendante de la structure de graphe, la variance locale est une estimation de la variance totale. Si les valeurs voisines de  $x$  sont corrélées positivement, alors la variance locale sous-estime la variance totale. Le coefficient de contiguïté  $c(x)$  est défini comme le rapport de la variance locale à la variance totale.

$$c(x) = v_l(x) / var(x) \quad [3.7 - 3]$$

Sous l'hypothèse selon laquelle les valeurs  $x_i$  sont des réalisations de variables aléatoires normales indépendantes, on peut calculer les quatre premiers moments du coefficient  $c(x)$  en fonction de la trace des puissances de la matrice  $\mathbf{M}$  associée au graphe<sup>1</sup>. On voit ainsi que pour le graphe des départements français (pour lequel deux sommets-départements sont joints par une arête s'ils ont une frontière commune) la distribution de  $c(x)$  est très proche d'une distribution normale.

<sup>1</sup> Pour un exposé plus complet cf. Lebart et Tabard (1973). Pour d'autres applications de la notion de contiguïté, cf. Aluja Banet et Lebart. (1984).



*Analyse des correspondances des matrices associées M*

On verra au § 4.1.5, à propos de l'interprétation des taux d'inertie, un exemple d'analyse des correspondances de la matrice  $\mathbf{M}$  associée à un cycle. Montrons que ce type d'analyse a des propriétés optimales en terme de contiguïté : le premier facteur direct  $\varphi$  (cf. § 4.1.5) associé à la plus grande valeur propre  $\lambda$  est la fonction de contiguïté minimale et le coefficient de Geary correspondant vaut :

$$c(\varphi) = 1 - \sqrt{\lambda}$$

Cette propriété est vraie pour les graphes réguliers (de degrés  $m_i$  constants) et s'étend à tous les graphes sous réserve d'une modification de la définition du numérateur de  $c(x)$  : pour le calcul de la moyenne de  $x$  et de sa variance totale, on pondère chaque sommet par son degré.

Dans ces conditions, si  $\mathbf{N}$  désigne la matrice diagonale telle que  $n_{ii} = m_i$  et si  $\mathbf{x}$  désigne le vecteur des observations de  $x$ , supposé centré, alors  $c(x)$  s'écrit :

$$c(x) = \frac{\mathbf{x}'(\mathbf{N} - \mathbf{M})\mathbf{x}}{\mathbf{x}'\mathbf{N}\mathbf{x}}$$

Le minimum  $\mu$  de  $c(x)$  est alors aussi le minimum de  $\mathbf{x}'(\mathbf{N} - \mathbf{M})\mathbf{x}$  avec la contrainte  $\mathbf{x}'\mathbf{N}\mathbf{x} = 1$  c'est-à-dire la plus petite valeur propre  $\mu$  donnée par le système (cf. § 1.1.7) :

$$(\mathbf{N} - \mathbf{M})\mathbf{x} = \mu \mathbf{N}\mathbf{x}$$

que l'on écrit :

$$\mathbf{N}^{-1}\mathbf{M}\mathbf{x} = (1 - \mu) \mathbf{x}$$

On reconnaît dans cette dernière formule la formule de transition de l'analyse des correspondances de la matrice symétrique  $\mathbf{M}$ .

On posera donc  $\sqrt{\lambda} = \varepsilon(1 - \mu)$ , avec  $\varepsilon = 1$  si le facteur est direct,  $\varepsilon = -1$  s'il est inverse. Le minimum de  $\mu$  correspond donc au maximum de  $\lambda$  pour  $\varepsilon > 0$ . Ainsi :

$$\text{Min}\{c(x)\} = 1 - \sqrt{\lambda_{\max}}$$

Les deux premiers facteurs sont donc les deux fonctions ( $\mathbf{N}$ -orthogonales) de contiguïté minimale, propriété qui assure une représentation graphique plane du graphe respectant le mieux possible les voisinages entre sommets<sup>1</sup> (le critère étant le coefficient de contiguïté, c'est-à-dire un critère de moindres carrés appliqué aux couples d'observations).

<sup>1</sup> On trouvera dans Benzécri (1973, Tome II B, n°10 : "Sur l'analyse de la correspondance définie par un graphe") des exemples donnant lieu à des résolutions numériques ou analytiques de description de graphes particuliers (cartes géographiques, réseaux à mailles carrées, produits tensoriels de réseaux, etc.). On observe en particulier dans ces cas des "effets Guttman à plusieurs dimensions", ce qui se traduit par des vecteurs propres de rangs élevés dont les composantes sont des fonctions polynomiales de celles des premiers vecteurs propres.

### 3.7.3 Analyse locale

Généralisons les résultats précédents au cas de plusieurs variables <sup>1</sup>. La covariance locale entre deux variables  $x_j$  et  $x_{j'}$  s'écrit :

$$\text{cov}_l(x_j, x_{j'}) = \frac{1}{2m} \sum_{i=1}^n \sum_{i'=1}^n m_{ii'} (x_{ji} - x_{ji'}) (x_{j'i} - x_{j'i'}) .$$

Si  $\mathbf{X}$  désigne la matrice d'ordre  $(n, p)$  de terme général  $x_{ji}$  ( $n$  observations de  $p$  variables), la matrice des covariances locales  $\mathbf{V}_l$  s'écrit :

$$\mathbf{V}_l = \frac{1}{2m} \mathbf{X}' (\mathbf{N} - \mathbf{M}) \mathbf{X} \quad [3.7-4]$$

Si le graphe est formé de cliques disjointes de mêmes tailles (structure de partition en classes d'égales importances), la matrice  $\mathbf{V}_l$  est proportionnelle à la matrice  $\mathbf{D}$  de variance intra-classes, qu'elle généralise dans ce cas. Si les classes sont d'effectifs inégaux, le système de pondération par le nombre de couples fait qu'il n'y a plus proportionnalité entre ces deux matrices <sup>2</sup>.

On peut définir une matrice des corrélations locales, de terme général :

$$c_l(x_j, x_{j'}) = \frac{\text{cov}_l(x_j, x_{j'})}{\sqrt{v_l(x_j) v_l(x_{j'})}}$$

La diagonalisation de cette matrice nous fournit, comme en analyse en composantes principales, une image des liaisons existant au niveau local, qu'il peut être intéressant de confronter aux liaisons globales (ainsi, dans le cas de données géographiques, l'opposition entre grandes régions très différentes peut masquer des covariations que l'analyse de la matrice des corrélations locales restitue).

### 3.7.4 Analyse de contiguïté et projections révélatrices

#### a – Analyse de contiguïté

La variance locale  $v_l(\mathbf{u})$  d'une combinaison linéaire  $u(i)$  des  $p$  variables s'écrit en fonction de la matrice de contiguïté, avec les notations habituelles :

<sup>1</sup> Alors que le coefficient de contiguïté de Geary est l'analogue, dans le cas d'un ensemble fini, d'un point du *variogramme* (correspondant à la distance "1" dans le cas isotropique) utilisé en géostatistique (Matheron, 1963), la matrice des covariances locales est l'analogue, dans le cas fini ou discret, de la matrice de *codispersion intrinsèque* qui intervient dans la théorie des variables régionalisées (Matheron, 1965).

<sup>2</sup> Cette proportionnalité est rétablie si, comme l'a remarqué Mom (1988), on modifie le coefficient de contiguïté en prenant la moyenne des carrés des différences entre chaque sommet et la *moyenne des sommets* qui lui sont contigus sur le graphe. La variance locale devient alors la variance des différences locales, selon les termes d'Escofier (1989).

$$v_l\{u(i)\} = v_l\left\{\sum_{j=1}^p u_j x_{ji}\right\} = \frac{1}{2m} \sum_{i=1}^n \sum_{i'=1}^n m_{ii'} (u(i) - u(i'))^2 = \mathbf{u}' \mathbf{V}_l \mathbf{u}$$

Si  $\mathbf{V}$  désigne la matrice des covariances totales, le coefficient de contiguïté de la combinaison linéaire  $u(i)$  s'écrit comme le quotient des deux formes quadratiques :

$$c(\mathbf{u}) = \frac{\mathbf{u}' \mathbf{V}_l \mathbf{u}}{\mathbf{u}' \mathbf{V} \mathbf{u}} \quad [3.7 - 5]$$

La recherche des combinaisons linéaires de contiguïté minimale (analyse de contiguïté) constitue, dans le cas de graphes réguliers, une généralisation de l'analyse factorielle discriminante, qui se réduit à celle-ci lorsque le graphe est formé de cliques disjointes. Dans le cas de graphes quelconques, il faut modifier la variance locale selon les préconisations de Mom (1988) pour avoir une généralisation stricte.

L'analyse de contiguïté est beaucoup moins utilisée que l'analyse discriminante qui a le mérite de rapprocher des données complexes et une structure très simple (la structure de partition)<sup>1</sup>. Elle peut être utilisée dans le cadre d'une démarche s'apparentant aux techniques dites de *projections révélatrices* (cf. Caussinus, 1992) qui, très schématiquement, cherchent des directions "intéressantes" plutôt que des dimensions principales au sens des moindres carrés<sup>2</sup>. Il existe autant de variantes de la méthode qu'il existe de façons de définir l'intérêt d'une projection.

## b – Représentation de groupes par projection

Si l'on veut déterminer une projection qui sépare le mieux possible des groupes existant dans l'ensemble des observations (sans connaître *a priori* ces groupes, sinon l'analyse factorielle discriminante classique répond à la question), on peut procéder de la façon suivante. On part d'un tableau de données  $\mathbf{X}$  d'ordre  $(n, p)$  pour lequel on n'a aucune information externe. On définit une relation de contiguïté sur l'ensemble des lignes de  $\mathbf{X}$  (il s'agit ici d'une contiguïté *a posteriori*) à partir d'un seuil de distance  $d_0$ . Parmi les  $n(n-1)$  couples d'observations (lignes de  $\mathbf{X}$ ) dans l'espace  $\mathbb{R}^p$ , les couples d'observations  $(i, i')$  tels que  $d(i, i') \leq d_0$  sont déclarés "contigus". On définit donc la matrice de contiguïté  $\mathbf{M}$  par les relations :

$$m_{ii'} = 1 \quad \text{si } d(i, i') \leq d_0 \quad \text{et } m_{ii'} = 0 \quad \text{sinon}$$

Une seconde façon de définir une relation de contiguïté *a posteriori* est de considérer comme contigus le pourcentage  $s_0$  ( $s_0 = 10$  par exemple) des couples les plus proches au sens de  $d(i, i')$ , ce qui permet de définir un seuil  $d_0$  après le calcul des  $n(n-1)/s_0$  plus petites distances.

<sup>1</sup> Pour des programmes de calcul et des applications de l'analyse de contiguïté, cf. Lebart et Tabard (1973).

<sup>2</sup> L'expression "Projection révélatrice" est la traduction, par des auteurs français (Escoufier, Caussinus) de l'expression "projection pursuit" (cf. Friedman et Tukey, 1974; Friedman, 1987; Jones et Sibson, 1987).

Une troisième façon utilise les  $k$  plus proches voisins : sont considérés comme contigus à la ligne  $i$  de  $X$  les  $k$  lignes les plus proches au sens de la distance  $d(i, i')$ . Cette méthode permet d'obtenir un graphe régulier, (avec les notations précédentes :  $m_i = k$ ) mais peut rattacher artificiellement des points isolés ou des petits groupes de points, au graphe d'ensemble qui est nécessairement connexe.

Une fois déterminée la matrice  $M$ , l'analyse de contiguïté, qui calcule les combinaisons linéaires réalisant les minima de  $c(u)$  donné par la formule [3.7 - 5], va produire une représentation qui respectera au mieux la structure de graphe et donc les plus fortes proximités entre points. En revanche, les distances moyennes ou grandes joueront un rôle moins important, ce qui a pour effet de "déplier" une éventuelle structure continue (cf. figure 3.7 - 2).

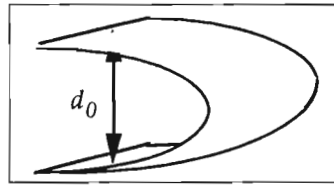


Figure 3.7 - 2

**Exemple de dépliage d'une structure par analyse de contiguïté**  
*Si le seuil est inférieur à la distance  $d_0$ , aucune arête ne joindra les deux plis; le graphe épousera donc la forme de la surface, qui sera dépliée dans les premiers plans de l'analyse.*

On peut imaginer qu'au lieu de sélectionner les arêtes les plus courtes, on garde toutes les arêtes, que l'on pondère par une fonction décroissante de la distance (le graphe de contiguïté devenant un graphe complet valué). On rejoint alors une série de travaux sur ce sujet plus proches des approches classiques de directions révélatrices.

Les premiers travaux sur ces thèmes sont ceux de Art *et al.* (1982), de Gnanadesikan *et al.* (1982). Ils ont été suivi des travaux de Yenyukov (1988), Caussinus et Ruiz (1990)<sup>1</sup>.

### c – Liens avec les analyses partielles

Comme indiqué au paragraphe *d* ci-dessus, on peut définir (au moins de trois façons différentes) une matrice de contiguïté  $M$  d'ordre  $(n, n)$  à partir d'un tableau de données que l'on appellera maintenant  $Z$  d'ordre  $(n, q)$ . Si l'on désire étudier des corrélations partielles entre les  $p$  colonnes d'une matrice  $X$  d'ordre  $(n, p)$  "à  $Z$  constant", on peut calculer la matrice des

<sup>1</sup> L'approche par analyse de contiguïté permet de mettre en évidence les deux structures qui sont confrontées: la structure locale, traduite sous forme de graphe (qui peut lui même être décrit par analyse des correspondances; cf. § 3.7-2 ci-dessus), la structure globale (analyse en composantes principales de  $X$ ), et le compromis entre les deux structures, décrit par l'analyse de contiguïté (cf. Burtschy et Lebart, 1991).

covariances partielles par la formule [3.6 - 3], mais on peut également calculer la matrice des covariances locales données par la formule [3.7 - 4] où  $\mathbf{M}$  est une matrice de contiguïté issue de  $\mathbf{Z}$  (et  $\mathbf{N}$  la matrice diagonale des degrés de  $\mathbf{M}$ ). Cette mesure de covariance partielle à l'avantage d'être non-linéaire (vis-à-vis des colonnes à fixer). Elle a l'inconvénient d'exiger des calculs de distances entre les lignes de  $\mathbf{Z}$  et dépend donc des échelles des mesures ou des poids des colonnes de  $\mathbf{Z}$ , ce qui n'est évidemment pas le cas pour des covariances calculées sur des résidus obtenus par régression multiple.

### 3.7.5 Extensions, généralisations, applications

Plusieurs variantes ou généralisations sont possibles autour de la notion de contiguïté. Déjà, à l'origine de ces travaux, les coefficients de Geary (1954) et de Moran (1948, 1954) constituaient deux mesures possibles (et très voisines) du degré de contiguïté<sup>1</sup>.

Citons brièvement, sans être exhaustif, quelques extensions ou applications : Le Foll (1982) introduit une pondération des sommets du graphe (les arêtes sont alors valuées par les produits des masses des sommets adjacents); Le Foll et Burtschy (1983) confrontent l'analyse locale et l'analyse des correspondances classique pour décrire des tableaux d'échanges; Carlier (1985) étudie les évolutions de tables de contingence par plusieurs méthodes dont l'analyse locale; Sabatier (1987) situe l'analyse locale dans un formalisme qui intègre les analyses partielles. Les travaux de Mom (1988) ont été cités précédemment.

Escofier (1989) introduit, dans la lignée de Mom, mais indépendamment, les intéressantes notions d'analyse lissée et d'analyse des différences locales, qui généralisent les notions d'analyse inter et intra (voir § 3.7.6 ci-dessous).

Dans l'analyse lissée, chaque point-individu  $i$  (ligne  $i$  de  $\mathbf{X}$ ) est remplacé par le barycentre de ses voisins sur le graphe. Ceci revient, avec nos notations (qui ne concernent que le cas où les sommets ont des poids *a priori* identiques, mais peuvent avoir des degrés différents) à remplacer  $\mathbf{X}$  par  $\mathbf{N}^{-1}\mathbf{MX}$ . Ce lissage a pour effet d'éliminer les variations locales.

Dans l'analyse des différences locales, on retranche à chaque vecteur-individu le barycentre des points voisins (on remplace  $\mathbf{X}$  par  $\mathbf{X} - \mathbf{N}^{-1}\mathbf{MX}$ ). On est ici dans une situation très voisine de l'analyse locale. On reviendra sur ces analyses au paragraphe suivant.

Cazes et Moreau (1991), Moreau (1992) considèrent le cas d'une double structure de graphe, présente à la fois sur les lignes et les colonnes d'une table de contingence. Faraj (1993) utilise l'analyse locale comme une analyse

<sup>1</sup> On pourra consulter les ouvrages de Cliff et Ord (1981), Ripley (1981), pour une vue plus large de la panoplie des outils disponibles.

partielle pour fixer l'effet de plusieurs variables nominales. Mentionnons enfin une synthèse de travaux sur ces thèmes par Méot *et al.* (1993)<sup>1</sup>.

### 3.7.6 Cas particuliers : Structure de partition

Il est fréquent que l'ensemble des individus ou observations soit partitionné en  $q$  classes connues *a priori* et jouant un rôle privilégié par rapport aux variables-colonnes du tableau de données  $X$  d'ordre  $(n, p)$ . Cette situation a été rencontrée en analyse factorielle discriminante (section 3.3) : il s'agissait alors de prédire l'appartenance d'un individu à une classe à partir des valeurs des variables pour cet individu.

Selon la formule de Huygens (formule [2.2 - 1] du § 2.2.3.a), l'inertie totale du nuage se décompose en inertie inter-classes (variabilité entre les classes) et inertie intra-classes (variabilité à l'intérieur des classes)<sup>2</sup> :

$$I = I_{inter} + I_{intra}$$

A cette décomposition est associée deux analyses : l'analyse *inter-classes* qui décrit les positions relatives des classes et ignore les individus, et l'analyse *intra-classes* qui s'attache à décrire les différences de comportement à l'intérieur des classes ce qui revient à éliminer l'effet dû à la structure de partition.

#### a – Analyse inter-classes

L'analyse inter-classes est simplement l'analyse du tableau agrégé d'ordre  $(q, p)$ . On a vu que l'analyse factorielle discriminante est une analyse inter-classes particulière (§3.3.4.c)<sup>3</sup>. Dans le cas où les variables sont nominales, on réalise l'analyse des correspondances du tableau des centres de gravité (ou tableau des barycentres) des  $q$  groupes d'individus, obtenu en croisant les classes de la partition avec les modalités des autres variables<sup>4</sup>. L'analyse inter-classes correspond dans ce cas à l'analyse discriminante barycentrique (cf. §3.3.7.b). L'analyse inter-classes est clairement un cas particulier de l'analyse lissée précitée lorsque le graphe est associé à une partition.

<sup>1</sup> Il faudrait citer, dans ce survol des utilisations de la notion de contiguïté, les méthodes de classification faisant appel aux contraintes de contiguïté. Une revue en est faite par Gordon et Finden (1985).

<sup>2</sup> Dans le cas de variables continues, il s'agit plus spécifiquement de la décomposition de la matrice de covariance (ou de corrélation si les variables sont réduites) en variance inter-classes (variance des moyennes des classes) et en variance intra-classes (variance de chaque classe autour de sa moyenne) donnée par la formule [3.3 - 1].

<sup>3</sup> Elle peut en effet être décrite comme une analyse en axes principaux des points-moyens de chacune des classes dans la métrique définie par l'inverse de la matrice des covariances "intra-classes".

<sup>4</sup> Il s'agit en fait d'une bande d'un tableau de Burt (cf. §1.4.7.b).

## b – Analyse intra-classes

L'analyse intra-classes permet d'étudier les différences de comportement à l'intérieur des classes en analysant la dispersion des individus à l'intérieur de leurs classes d'appartenance (cf. Benzécri, 1983; Cazes, 1986-a; Benali et Escofier, 1990).

Chaque individu est représenté par un point dont les coordonnées expriment l'écart entre ses propres coordonnées et celles du centre de gravité de sa classe. L'inertie inter-classes est ainsi éliminée. On ne cherche donc plus à savoir de quelle manière un individu se différencie de l'ensemble du nuage mais comment il se différencie de l'ensemble des individus appartenant à la même classe. On s'affranchit ainsi de l'influence de la variable de partition en étudiant les liaisons entre les variables à analyser, conditionnellement à la variable de partition.

L'analyse intra-classes est un cas particulier de l'analyse des différences locales (graphe associé à une partition) et également un cas particulier de l'analyse partielle (cf. section 3.6) lorsque la variable exogène  $z$  est nominale.

Escofier (1987) introduit une méthode d'analyse intra-classes dans le cas de variables nominales, appelée *analyse des correspondances multiples conditionnelles*, qui est en fait un cas particulier de la généralisation de l'analyse des correspondances proposée également par Escofier (1984). L'influence de la variable de partition est éliminée ; le nuage des individus est recentré par classe, et le nuage des modalités est projeté sur l'orthogonal du sous-espace engendré par les modalités de la variables de partition<sup>1</sup>.

Une extension de l'analyse des correspondances multiples conditionnelles, est étudiée par Piron (1990, 1992) lorsque les variables sont des fréquences. Dans ce cas, la structure induite sur les individus relève d'une série de partitions emboîtées (structure fréquente dans les relevés géographiques).

Pour le cas de doubles partitions (partition  $Q$  sur les lignes et partition  $S$  sur les colonnes d'une table de contingence) Cazes (1986-a et 1986-b), et Cazes, Chessel et Doledec (1988) proposent l'*analyse des correspondances internes* qui consiste à réaliser l'analyse intra-classes en considérant un double centrage dans l'espace des lignes et dans celui des colonnes. On projette d'une part le nuage des points-lignes sur l'orthogonal du sous-espace engendré par les modalités de la variable de partition  $Q$ , et d'autre part le nuage des points-colonnes sur l'orthogonal du sous-espace engendré par les modalités de la variable de partition  $S$ .

---

<sup>1</sup> L'analyse des correspondances multiples conditionnelle conserve toutes les propriétés de l'analyse des correspondances. Elle est implémentée dans le logiciel SPAD.N sous forme de procédure.

## Tableaux multiples, groupes de variables

L'analyse des tableaux multiples est un très vaste domaine de recherche que l'on ne fera qu'effleurer dans cette section, en se limitant à quelques situations spécifiques, proches de la démarche exploratoire.

Le théorème d'Eckart et Young (décomposition aux valeurs singulières étudiée en section 1.1) qui est à la base des méthodes factorielles, n'admet pas de généralisation au sens suivant : il n'existe pas de décomposition optimale unique d'un tableau à trois entrées (empilement de  $q$  tableaux  $X_k$ , chacun d'ordre  $(n, p)$ ) en tableaux de rangs 1.

En revanche, il existe des modèles particuliers, qui varient selon les disciplines et la nature des tableaux, pour aborder ce type de données.

### 3.8.1 Quelques travaux de référence

Commençons par évoquer quelques travaux de référence sur le thème des tableaux à plusieurs dimensions<sup>1</sup>.

Les premiers travaux sur ce thème sont ceux de Tucker (1964, 1966) puis ceux de Harshman (1970), tous les deux dans le cadre de l'analyse factorielle classique. Montrons brièvement quelles sont les relations qui sont à la base de ces modèles.

L'un des modèles de Tucker, dit TUCKALS-3 (Kroonenberg et de Leeuw, 1980), s'applique à une séquence de matrices symétriques d'ordre  $(p, p)$   $S_1, \dots, S_q$  (qui sont par exemple des matrices de distances entre individus). Il conduit à la relation ( $s_{ijk}$  désignant une estimation, par le modèle, de l'élément  $(i, j)$  de la matrice  $S_k$ ) :

$$s_{ijk} = \sum_{u=1}^p \sum_{v=1}^p \sum_{t=1}^r a_{iu} a_{jv} b_{kt} c_{uvt}$$

Le modèle dit PARAFAC, de Harshman, donne lieu à une relation analogue, mais plus simple.

---

<sup>1</sup> On trouvera une synthèse et une classification des principales démarches dans l'ouvrage de Kroonenberg (1983) qui a prolongé les travaux de Tucker. On pourra aussi consulter la revue comparative de Carlier *et al.* (1988), qui fait d'ailleurs partie d'un recueil entièrement consacré à ce thème (Coppi et Bolasco, 1989). Une revue se trouve également dans Kiers (1989). Sur le thème plus circonscrit des évolutions de tables de contingence, cf. Carlier (1985), van der Heijden (1987).



Pour une série de matrices  $X_k$  d'ordre  $(n, p)$ , le terme général  $x_{ijk}$  peut s'écrire :

$$x_{ijk} = \sum_{t=1}^r a_{it} b_{jt} c_{kt}$$

Ces formules peuvent être vues comme des généralisations possibles de la formule de reconstitution de données<sup>1</sup>.

Une autre méthode très utilisée dans le contexte des méthodes de *multidimensional scaling* est la méthode INDSCAL de Carroll et Chang (1970) qui est un cas particulier de la méthode PARAFAC de Harshman.

Ces exemples laissent imaginer le nombre de modèles et de variantes possibles.

Les quatre paragraphes de cette section seront tous consacrés à une structure de tableaux multiples très particulière, mais fréquente en pratique : il s'agit d'un tableau  $X$  d'ordre  $(n, p)$  tel que :

$$X = (X_1, X_2, \dots, X_k, \dots, X_q)$$

Les différents blocs n'ont pas forcément le même nombre de colonnes et cette structure est par conséquent plus générale qu'un tableau à trois entrées.

Selon les cas, les lignes seront des individus ou observations, les colonnes de chaque bloc des variables. Les blocs peuvent correspondre à des instants ou des contextes différents pour les mêmes variables, ou à des groupes de variables différents.

La section 3.6 a abordé le cas de l'analyse d'un tableau de données de type  $R = (X, Z)$  dans laquelle les deux ensembles de colonnes (colonnes de  $X$  et de  $Z$ ) jouaient des rôles dissymétriques. Il existe des circonstances dans lesquelles les rôles sont parfaitement symétriques. C'est le cas notamment des méthodes d'analyses procrustéennes orthogonales qui visent à comparer deux structures de distances sur les mêmes objets, ceux-ci étant décrits successivement par deux ensembles différents de variables (§ 3.8.2).

La méthode STATIS (§ 3.8.3) et l'analyse factorielle multiple (§3.8.4) sont proches à bien des égards dans leurs procédures mais se différencient dans les options de traitements. Elles procèdent en trois étapes : la comparaison globale des tableaux, la représentation du nuage moyen et la représentation simultanée des tableaux.

Brièvement évoquée à propos de l'analyse canonique, l'analyse canonique généralisée (on désigne sous ce nom l'une des généralisations possibles de l'analyse canonique, en fait la plus mentionnée et utilisée) sera présentée dans un cadre plus général (§ 3.8.5). Cette méthode, assez délicate à utiliser directement en pratique, fournit un cadre théorique simple commun aux

<sup>1</sup> Cf. par exemple: Hayashi et Hayashi (1982) pour un algorithme d'estimation des coefficients du modèle.

principales méthodes factorielles exploratoires et aux méthodes explicatives de base des sections 3.1 à 3.3, qu'elle contient toutes comme cas particulier.

### 3.8.2 Analyses procrustéennes <sup>1</sup>

Les méthodes d'analyse procrustéennes tentent de répondre à une préoccupation fréquente en statistique multidimensionnelle :  $n$  individus ou observations sont décrits d'une part par  $p$  variables (colonnes de  $\mathbf{X}$ ), d'autre part par  $q$  autres variables (colonnes de  $\mathbf{Z}$ ). Comment comparer les deux nuages d'individus, les deux systèmes de distances entre individus ?

C'est Tucker (1958) qui proposa à l'origine une telle méthode pour comparer deux batteries de tests passés sur les mêmes individus<sup>2</sup>. La technique a ensuite été étudiée par Cliff (1966), Schönemann (1968), Schönemann et Carrol (1970), puis généralisée par Gower (1975, 1984)<sup>3</sup>.

#### a – Analyse procrustéenne orthogonale

Fixons  $\mathbf{X}$  par exemple (les rôles de  $\mathbf{X}$  et  $\mathbf{Z}$  sont symétriques) et supposons  $p = q$ . Ceci n'est pas une restriction car, si par exemple  $p > q$ , on peut toujours compléter le tableau  $\mathbf{Z}$  par  $p - q$  colonnes nulles.

Si les lignes de  $\mathbf{Z}$ , d'ordre  $(n, p)$ , subissent toutes un même déplacement (translation et rotation dans  $\mathbb{R}^p$ ),  $\mathbf{Z}$  est transformé en  $\mathbf{ZB} + \mathbf{T}$ , où  $\mathbf{T}$  est une matrice d'ordre  $(n, p)$  dont les colonnes peuvent être différentes, mais constantes (translation) et où  $\mathbf{B}$   $(p, p)$  est une matrice orthogonale (rotation ou symétrie par rapport à l'origine).

On cherchera à rendre minimale la somme des carrés  $s$  des écarts entre  $\mathbf{X}$  et  $(\mathbf{ZB} + \mathbf{T})$ , qui peut s'écrire<sup>4</sup> :

$$s = \text{trace} (\mathbf{X} - \mathbf{ZB} - \mathbf{T})' (\mathbf{X} - \mathbf{ZB} - \mathbf{T})$$

Le critère  $s$  s'écrit encore, si les tableaux  $\mathbf{X}$  et  $\mathbf{Z}$  sont centrés en colonnes :

$$s = \text{trace} (\mathbf{X} - \mathbf{ZB})' (\mathbf{X} - \mathbf{ZB}) + \mathbf{T}'\mathbf{T} \quad [3.8 - 1]$$

La recherche d'un minimum pour  $s$  implique  $\mathbf{T} = \mathbf{0}$  (aucune translation n'est requise quand les tableaux sont centrés).

<sup>1</sup> Procruste (ou Procuste, par altération) est un aubergiste de la mythologie grecque qui raccourcissait ou allongeait ses clients ( $\mathbf{X}$ , par exemple) pour les ajuster à la longueur de son lit ( $\mathbf{Z}$ ). Thésée mettra fin à ses jours en lui infligeant le même supplice.

<sup>2</sup> On peut de la même façon comparer un même ensemble de variables sur deux ensembles d'individus différents. C'est le cas si l'on veut comparer deux matrices des corrélations (une matrice des corrélations globales, par exemple, à confronter à une matrice des corrélations locales).

<sup>3</sup> Cf. également Lafosse (1985), Fichet et al. (1990).

<sup>4</sup> Rappelons que  $\text{trace}(\mathbf{A}'\mathbf{A}) = \sum_{i,j} a_{ij}^2$ ; que  $\text{trace} \mathbf{A} = \text{trace} \mathbf{A}'$ ; et que, lorsque les opérations sont possibles,  $\text{trace} (\mathbf{A} + \mathbf{C}) = \text{trace} \mathbf{A} + \text{trace} \mathbf{C}$ ;  $\text{trace} \mathbf{AC} = \text{trace} \mathbf{CA}$ .

Développant l'expression du critère  $s$  et en tenant compte du fait que :

$$\text{trace } \mathbf{B}'\mathbf{Z}'\mathbf{Z}\mathbf{B} = \text{trace } \mathbf{Z}'\mathbf{Z}\mathbf{B}\mathbf{B}' = \text{trace } \mathbf{Z}'\mathbf{Z}$$

il vient :

$$s = \text{trace } (\mathbf{X}'\mathbf{X} + \mathbf{Z}'\mathbf{Z} - 2 \mathbf{B}'\mathbf{Z}'\mathbf{X})$$

Rendre minimal le critère  $s$  revient à rendre maximal  $\text{trace } (\mathbf{B}'\mathbf{Z}'\mathbf{X})$ .

Ecrivons la formule de reconstitution des données (cf. section 1.1, formule [1.1-6]) issue de l'analyse générale (décomposition aux valeurs singulières) du tableau  $\mathbf{Z}'\mathbf{X}$  :

$$\mathbf{Z}'\mathbf{X} = \sum_{\alpha=1}^p \sqrt{\lambda_{\alpha}} \mathbf{v}_{\alpha} \mathbf{u}'_{\alpha}$$

d'où :

$$\text{trace}(\mathbf{B}'\mathbf{Z}'\mathbf{X}) = \text{trace} \left( \sum_{\alpha=1}^p \sqrt{\lambda_{\alpha}} \mathbf{B}'\mathbf{v}_{\alpha} \mathbf{u}'_{\alpha} \right) = \sum_{\alpha=1}^p \sqrt{\lambda_{\alpha}} (\mathbf{u}'_{\alpha} \mathbf{B}'\mathbf{v}_{\alpha})$$

$\mathbf{B}$  étant orthogonal et  $\mathbf{v}_{\alpha}$  unitaire,  $\mathbf{B}'\mathbf{v}_{\alpha}$  est unitaire et on aura toujours  $\mathbf{u}'_{\alpha} \mathbf{B}'\mathbf{v}_{\alpha} \leq 1$ . On aura  $\mathbf{u}'_{\alpha} \mathbf{B}'\mathbf{v}_{\alpha} = 1$  si et seulement si  $\mathbf{B}'\mathbf{v}_{\alpha} = \mathbf{u}_{\alpha}$ .

D'où la relation  $\mathbf{B}'\mathbf{V} = \mathbf{U}$  et la solution cherchée  $\mathbf{B} = \mathbf{V}\mathbf{U}'$ .

L'analyse procrustéenne orthogonale implique donc la décomposition aux valeurs singulières de  $\mathbf{Z}'\mathbf{X}$  et donc la diagonalisation de la matrice  $\mathbf{X}'\mathbf{Z}\mathbf{Z}'\mathbf{X}$ .

### *Autre présentation de l'analyse procrustéenne orthogonale*

On peut donner une autre présentation de cette méthode, en procédant de façon hiérarchique, par extraction progressive d'axes procrustéens. La méthode est analogue à l'analyse canonique, aux contraintes de normalisation près.

Les tableaux  $\mathbf{X}$  et  $\mathbf{Z}$  étant centrés, elle consiste à chercher deux combinaisons linéaires  $\mathbf{X}\mathbf{u}$  et  $\mathbf{Z}\mathbf{v}$ , à coefficients normés ( $\mathbf{u}'\mathbf{u} = 1$ ,  $\mathbf{v}'\mathbf{v} = 1$ ), de covariances maximales, c'est-à-dire telles  $\mathbf{v}'\mathbf{Z}'\mathbf{X}\mathbf{u}$  soit maximale.

Une démonstration en tout point analogue à celle du paragraphe 3.1.2.a (comme dans le cas de l'analyse canonique, les deux multiplicateurs de Lagrange sont égaux à une même valeur  $\lambda$ ) nous montre alors que  $\mathbf{u}$  et  $\mathbf{v}$  sont solutions de :

$$\mathbf{X}'\mathbf{Z}\mathbf{Z}'\mathbf{X}\mathbf{u} = \lambda^2 \mathbf{u} \quad \text{et} \quad \mathbf{Z}'\mathbf{X}\mathbf{X}'\mathbf{Z}\mathbf{v} = \lambda^2 \mathbf{v}$$

qui sont bien les équations de l'analyse générale du tableau  $\mathbf{Z}'\mathbf{X}$ .

En extrayant les différents axes (avec des contraintes d'orthogonalité usuelles), et en notant  $\mathbf{U}$  et  $\mathbf{V}$  les matrices orthogonales contenant en colonnes les vecteurs  $\mathbf{u}_{\alpha}$  et  $\mathbf{v}_{\alpha}$  correspondant aux différents axes indexés par  $\alpha$ , on aura rendu maximal le critère :  $\text{trace}(\mathbf{V}'\mathbf{Z}'\mathbf{X}\mathbf{U})$  (les éléments diagonaux de cette matrice sont en effet les covariances maximales trouvées).

Remarquons que, jusqu'ici, on n'a pas supposé  $p = q$  dans cette présentation.

Or rendre maximale cette trace revient à rendre minimal le critère lorsque  $p = q$  ( $\mathbf{U}$  et  $\mathbf{V}$  étant deux matrices orthogonales) :

$$s_1 = \text{trace}(\mathbf{XU} - \mathbf{ZV})(\mathbf{XU} - \mathbf{ZV})'$$

On peut écrire cette quantité :

$$s_1 = \text{trace} \mathbf{U}(\mathbf{XU} - \mathbf{ZV})(\mathbf{XU} - \mathbf{ZV})\mathbf{U}'$$

Finalement :

$$s_1 = \text{trace}(\mathbf{X} - \mathbf{ZVU}')(\mathbf{X} - \mathbf{ZVU}')$$

qui coïncide avec le critère  $s$  de la première approche pour  $\mathbf{B} = \mathbf{VU}'$  (formule [3.8-1], avec  $\mathbf{T} = \mathbf{0}$ )

### b – Analyse procrustéenne sans contrainte<sup>1</sup>

On cherchera à rendre minimale la somme des carrés  $s$  des écarts entre  $\mathbf{X}$  et  $(\mathbf{ZA} + \mathbf{T})$ , ce qui revient à rendre minimal (si les tableaux  $\mathbf{X}$  et  $\mathbf{Z}$  sont centrés en colonnes), sans contrainte sur la matrice  $\mathbf{A}$ , le critère :

$$s = \text{trace}(\mathbf{X} - \mathbf{ZA})(\mathbf{X} - \mathbf{ZA})'$$

On trouve après un calcul analogue à celui du calcul des coefficients de régression multiple<sup>2</sup> :

$$\mathbf{A} = (\mathbf{Z}'\mathbf{Z})^{-1}\mathbf{Z}'\mathbf{X}$$

C'est la matrice des coefficients d'une régression simultanée, qui revient à effectuer séparément  $p$  régressions indépendantes pour chacune des  $p$  colonnes de  $\mathbf{X}$ . Dans ce cas, une analyse des résidus  $\mathbf{X} - \mathbf{ZA}$  (analyse partielle, cf. section 3.6) nous renseignera sur les éventuels traits structuraux de  $\mathbf{X}$  non-expliqués par  $\mathbf{Z}$ . Notons que l'analyse procrustéenne sans contrainte fait intervenir de façon dissymétrique les tableaux  $\mathbf{X}$  et  $\mathbf{Z}$ .

Il existe de nombreuses autres variantes des analyses procrustéennes (impliquant des dilatations, des axes obliques) pour lesquelles on pourra consulter les références citées.

### c – Formulaire de quelques méthodes d'analyse impliquant deux groupes de variables

Récapitulons quelques unes des méthodes d'analyse de tableaux du type  $\mathbf{R} = (\mathbf{X}, \mathbf{Z})$ , en donnant le formulaire des matrices à diagonaliser ou des matrices de coefficients :

<sup>1</sup> C'est l'approche initiale de Hurley et Cattell (1962) qui sont d'ailleurs à l'origine du nom de cette méthode.

<sup>2</sup> Le problème a été résolu au § 3.6.2 b ci-dessus.

$(X'X)^{-1} X'Z (Z'Z)^{-1} Z'X$	ou	$(Z'Z)^{-1} Z'X (X'X)^{-1} X'Z$	(Analyse canonique)
$X'Z (Z'Z)^{-1} Z'X$	et	$Z'X (X'X)^{-1} X'Z$	(Analyses projetées)
$X'(I - Z (Z'Z)^{-1} Z') X$	et	$Z'(I - X (X'X)^{-1} X') Z$	(Analyses partielles)
$X'Z Z'X$	ou	$Z'X X'Z$	(Analyse procrustéenne orthogonale)
$(Z'Z)^{-1} Z'X$	et	$(X'X)^{-1} X'Z$	(Analyse procrustéenne sans contrainte)

Dans les cas où  $X$  et  $Z$  sont des tableaux de variables numériques, celles-ci sont centrées. Si l'on excepte les cas de l'analyse canonique et de l'analyse procrustéenne sans contrainte (ou régression multiple simultanée), il est en général souhaitable de réduire les variables. Notons également que les analyses projetées et l'analyse procrustéenne orthogonale sont équivalentes à des analyses en composantes principales lorsque  $X = Z$ .

### 3.8.3 Méthode STATIS

La méthode STATIS<sup>1</sup> a été proposée par l'équipe d'Escoufier (1980, 1985 a)<sup>2</sup> pour permettre l'analyse conjointe de plusieurs tableaux de données. Elle s'applique à des tableaux de mesures dans lesquels un ensemble d'individus est décrit par plusieurs groupes de variables ou pour lequel les mêmes variables sont observées sur plusieurs groupes d'individus. L'objet est de comparer les tableaux, puis de décrire l'éventuelle structure commune aux différents tableaux, enfin d'appréhender les différences entre tableaux. Nous présentons seulement les principes de la méthode STATIS sans entrer dans les détails, renvoyant le lecteur à l'ouvrage de Lavit (1988). Nous nous plaçons dans le cadre de  $q$  tableaux de mesures de variables centrées-réduites observées sur les mêmes individus de poids égaux à 1.

#### a – Notations

On note  $n$  le nombre d'individus ;  $p$  le nombre total de variables (supposées ici centrées-réduites) tous groupes confondus ;  $p_k$  le nombre des variables du  $k^{\text{ième}}$  groupe ;  $X$  le tableau complet de terme général  $x_{ij}$  valeur de l'individu  $i$  pour la variable  $j$  ;  $X_k$  le sous-tableau de  $X$  associé au groupe  $k$  ;  $q$  représente le nombre de groupes.

L'individu  $i$  correspond à une ligne du tableau  $X = (X_1, X_2, \dots, X_k, \dots, X_q)$  ; à cet individu, dit "moyen", sont associés  $q$  individus dits "partiels", notés  $i^k$ , correspondant aux lignes des divers tableaux  $X_k$ .

#### b – Comparaison globale entre les tableaux : l'interstructure

On s'intéresse ici aux relations entre les  $q$  tableaux  $X_k$  d'ordre  $(n, p_k)$ . On considère les matrices de produits scalaires entre les individus  $X_k X'_k$  (ou

<sup>1</sup> Le sigle STATIS signifie "Structuration des Tableaux À Trois Indices de la Statistique.

<sup>2</sup> cf. L'Hermier des Plantes (1976), Pagès et al. (1976).

$X_k Q_k X'_k$  si l'on introduit une métrique particulière à chaque tableau  $Q_k$ , mais dans cet exposé schématique,  $Q_k = I$ ) de dimension  $(n, n)$  et l'on cherche à décrire les distances entre ces matrices. On considère pour cela chaque matrice  $X_k X'_k$ , notée  $W_k$ , comme un point dans l'espace  $\mathbb{R}^{n^2}$  obtenu en empilant les colonnes de cette matrice. On définit ainsi un nuage de  $q$  points-tableaux dans  $\mathbb{R}^{n^2}$  et le tableau associé  $W_{n^2}$  de dimension  $(n^2, q)$ .

L'analyse générale du tableau  $W_{n^2}$ , qui revient à diagonaliser la matrice  $S$  d'ordre  $(q, q)$  de terme général  $s_{kk'} = \text{trace}(W_k W_{k'})$ , permet de représenter les  $q$  points-tableaux dans un espace de faible dimension et de comparer globalement les tableaux entre eux. Si tous les tableaux sont voisins, ils seront concentrés près 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.

Si le nombre  $p_k$  de variables du tableau  $k$  n'est pas constant, on a intérêt à normer les termes de  $S$  en analysant la matrice  $\hat{S}$  de terme général  $\hat{s}_{kk'}$ , qui n'est autre que le coefficient  $Rv$  de Robert et Escoufier (1976) :

$$\hat{s}_{kk'} = \frac{\text{trace } W_k W_{k'}}{\sqrt{\text{trace } W_k^2 \text{ trace } W_{k'}^2}}$$

**Remarque :**

Dans le cas où l'on dispose d'un ensemble de variables observées sur  $q$  groupes d'individus, on considère les matrices de covariances (ou de corrélations si les variables sont réduites) de dimension  $(p, p)$ . On calculera alors, à partir d'un nuage de  $q$  points-tableaux dans l'espace  $\mathbb{R}^{p^2}$ , le tableau  $W_{p^2}$  de dimension  $(p^2, q)$ .

**c – Le nuage moyen ou compromis : l'intrastructure**

On cherche à construire un nuage moyen qui soit un compromis des  $q$  nuages correspondants aux tableaux  $X_k$ . Le compromis peut être calculé de différentes façons, en fonction de la nature des données et des connaissances a priori. Ce peut être une simple moyenne pondérée  $C_1$  des tableaux  $X_k$ , lorsqu'il s'agit par exemple de l'évolution d'un tableau impliquant les mêmes individus et les mêmes variables :

$$C_1 = \sum_{k=1}^q \alpha_k X_k$$

Si le nombre des variables  $p_k$  varie avec  $k$ , le compromis pourra toujours être calculé au niveau des produits scalaires (éventuellement normés) :

$$W_1 = \sum_{k=1}^q \alpha_k X_k X'_k = \sum_{k=1}^q \alpha_k W_k$$

Les promoteurs de cette stratégie d'analyse recommandent de prendre comme poids  $\alpha_k$  la coordonnée du tableau  $k$  sur le premier axe de l'analyse de l'interstructure : un tableau aura ainsi un poids d'autant plus élevé qu'il est représentatif de la tendance moyenne.

L'analyse du compromis revient ensuite à effectuer l'analyse en composantes principales ou l'analyse générale du tableau  $C_1$  ou  $W_1$  selon le cas. Elle permet donc de dégager la structure du nuage des individus commune aux tableaux.

#### d – Représentation simultanée des nuages partiels : les trajectoires

L'analyse de l'interstructure met en évidence les écarts entre les tableaux. L'intrastructure est décrite par le ou les compromis. Il reste à décrire les écarts par rapport au compromis, au niveau des variables et des individus. Si le tableau compromis est du type  $C_1$ , il est aisé de représenter en éléments supplémentaires, à partir des tableaux  $X_k$ , les trajectoires d'individus (un individu  $i$  est représenté par les  $q$  points  $i_k$ ) et, de façon similaire, les trajectoires de variables.

Dans le cas d'un compromis de type  $W_1$ , on peut toujours représenter les trajectoires d'individus (lignes des tableaux  $W_k$ ).

### 3.8.4 Analyse factorielle multiple

L'analyse factorielle multiple (Escofier et Pagès, 1983), traite des tableaux dans lesquels un ensemble d'individus est décrit par plusieurs groupes de variables. Les variables peuvent être continues, nominales et même, sous certaines conditions, de type fréquence. Toutefois, à l'intérieur d'un groupe, elles doivent être de même type.

Nous nous contentons ici d'esquisser les principales caractéristiques de la méthode, en nous plaçant dans le cas particulier de variables continues centrées-réduites de poids 1. Nous renvoyons le lecteur désireux d'approfondir l'analyse factorielle multiple à l'ouvrage de Escofier et Pagès (1988). Les notations de base sont les mêmes que pour la méthode STATIS.

#### a – Une analyse en composantes principales pondérée

Le fait de vouloir introduire plusieurs groupes de variables en tant qu'éléments actifs dans une même analyse factorielle impose d'équilibrer leur influence *a priori* dans cette analyse. Une analyse simultanée de plusieurs groupes dont les premiers facteurs seraient engendrés par un seul d'entre eux ne présenterait en effet que peu d'intérêt.

En analyse factorielle multiple, chaque variable du groupe  $k$  est pondérée par  $1/\sqrt{\lambda_1^k}$  où  $\lambda_1^k$  est la première valeur propre de l'analyse en composantes principales effectuées sur les variables de ce groupe  $k$ . A l'intérieur d'un

groupe, toutes les variables ont le même poids : la structure de chaque groupe est respectée. Géométriquement, cela revient à rendre égale à 1 l'inertie axiale maximum de chacun des  $k$  sous-nuages. Du fait de cette pondération, aucun groupe ne peut engendrer à lui seul le premier axe ; en revanche, un groupe multidimensionnel contribue à un plus grand nombre d'axes qu'un groupe unidimensionnel.

Le principe de l'analyse factorielle multiple repose sur une analyse en composantes principales du tableau complet  $\mathbf{X} = (\mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_k, \dots, \mathbf{X}_q)$ , les variables étant ainsi pondérées. Cette analyse permet d'équilibrer le rôle des groupes de variables et fournit une représentation des individus et des variables qui s'interprète selon les règles usuelles de l'analyse en composantes principales. Au-delà de cette analyse en composantes principales pondérée, la prise en considération de groupes de variables augmente les possibilités d'interprétation des résultats.

Le  $\alpha$ ïème facteur de l'analyse factorielle multiple de  $\mathbf{X}$  est noté  $\psi_\alpha$  dans  $\mathbb{R}^p$  et  $\varphi_\alpha$  dans  $\mathbb{R}^n$  ; il est associé à la valeur propre  $\lambda_\alpha$  ; la  $\alpha$ ïème valeur propre de l'analyse en composantes principales séparée de  $\mathbf{X}_k$  est notée  $\lambda_\alpha^k$ .

### b – Recherche de facteurs communs (intrastructures)

Au groupe de variables  $k$  correspond dans  $\mathbb{R}^n$  un sous-espace  $V_k$  à  $k$  dimensions ; un facteur commun est une dimension commune à ces sous-espaces. Cette idée est présente dans les analyses canoniques et *multicanoniques* (cas de plus de deux groupes). Mais ces analyses considèrent chaque nuage  $k$  uniquement au travers du sous-espace qu'il engendre, sans prendre en compte la répartition de l'inertie dans ce sous-espace. Comparée à ces méthodes, l'analyse factorielle multiple recherche des facteurs à la fois communs aux groupes de variables et représentant des directions de forte inertie de ces groupes.

Du fait de la pondération des variables, l'analyse factorielle multiple peut être interprétée comme une analyse multicanonique. En effet, dans ce cas l'inertie projetée des variables du groupe  $k$  sur la direction  $\mathbf{z}$  constitue une mesure de liaison entre la variable  $\mathbf{z}$  et le groupe de variables  $k$ . Cette mesure, notée  $L(\mathbf{z}, k)$ , possède les propriétés suivantes :

- $0 \leq L(\mathbf{z}, k) \leq 1$ .
- $L(\mathbf{z}, k) = 0 \Leftrightarrow \mathbf{z}$  est non corrélée avec chaque variable du groupe  $k$ .
- $L(\mathbf{z}, k) = 1 \Leftrightarrow \mathbf{z}$  est la première composante principale de  $k$ .

Le critère satisfait par la  $\alpha$ ïème composante principale (notée  $\mathbf{z}_\alpha$ ) de l'analyse factorielle multiple peut s'écrire, compte tenu des contraintes d'orthogonalité avec les  $\alpha - 1$  premières composantes principales :

$$\text{Max} \left( \sum_k L(\mathbf{z}_\alpha, k) \right)$$



Du point de vue de ce critère, les composantes principales de l'analyse factorielle multiple composent la suite de variables orthogonales les plus liées aux groupes de variables. En ce sens ce sont les facteurs communs à ces groupes.

### c – Représentation des groupes de variables (interstructure)

La mise en évidence de facteurs communs est une voie commode pour analyser les liaisons entre groupes de variables. On peut chercher à visualiser globalement ces liaisons par un graphique dans lequel chaque groupe est représenté par un point.

Au groupe de variables  $k$  on peut associer, comme dans la méthode STATIS, la matrice  $\mathbf{W}_k = \mathbf{X}_k \mathbf{X}'_k$  des produits scalaires entre individus. Toutes ces matrices sont d'ordre  $(n, n)$ . Ce sont des éléments de l'espace  $\mathbb{R}^{n^2}$ ; ces éléments constituent le nuage des  $k$  groupes de variables. L'analyse factorielle multiple fait intervenir d'autres éléments de  $\mathbb{R}^{n^2}$ : les matrices de produits scalaires associées à chaque composante principale normée  $\mathbf{z}_\alpha$ ; ces éléments, que l'on peut écrire  $\mathbf{z}_\alpha \mathbf{z}'_\alpha$  forment une base orthonormée d'un sous-espace de  $\mathbb{R}^{n^2}$ . C'est sur cette base que l'on projettera les  $k$  points-groupes, pour visualiser leurs proximités.

Cette représentation possède quelques propriétés remarquables. En particulier la projection de  $\mathbf{W}_k = \mathbf{X}_k \mathbf{X}'_k$  sur  $\mathbf{z}_\alpha \mathbf{z}'_\alpha$  est égale à  $L(\mathbf{z}, k)$ . Il est ainsi possible d'interpréter axe par axe les proximités entre les points-groupes.

### d – Représentations superposées des nuages partiels des groupes actifs (trajectoires)

A chaque groupe de variables est associé un nuage partiel d'individus. La comparaison directe des représentations issues des analyses en composantes principales séparées des  $\mathbf{X}_k$  ne répond pas directement à cet objectif car ces analyses, étant effectuées séparément, ne tiennent pas compte d'éventuelles structures communes. Il faudrait en fait une analyse procrustéenne généralisée pour résoudre ce problème.

En analyse factorielle multiple on projette les nuages partiels sur les axes principaux du nuage total. Bien qu'ils n'interviennent pas directement dans la construction des axes, les nuages partiels ne sont pas véritablement considérés comme supplémentaires puisque leurs données sont incluses dans le nuage total analysé. Il en résulte deux propriétés utiles lors de l'interprétation :

$$\psi_{\alpha i} = \frac{1}{p} \sum_k \psi_{\alpha i k}$$

le point "moyen"  $i$  est au centre de gravité  $\psi_{\alpha i}$  des points "partiels"  $\psi_{\alpha i k}$  qui lui sont homologues.

$$\Psi_{\alpha j^k} = \frac{1}{\lambda_1^k} \frac{1}{\sqrt{\lambda_\alpha}} \sum_{j \in J_k} x_{ij} \Phi_{\alpha j}$$

Cette relation n'est autre que la restriction au groupe  $k$  de l'une des relations usuelles de transition. L'individu partiel  $i^k$  apparaît du côté des variables pour lesquelles il a de fortes valeurs (les  $x_{ij}$  sont des valeurs centrées-réduites) et à l'opposé de celles pour lesquelles il a de faibles valeurs.

### Cas particuliers

Lorsque chaque groupe ne comporte qu'une seule variable quantitative, l'analyse factorielle multiple se confond avec une analyse en composantes principales. Lorsque chaque groupe ne comporte qu'une seule variable qualitative, l'analyse factorielle multiple se confond avec une analyse des correspondances multiples.

## 3.8.5 Analyse canonique généralisée

L'analyse canonique généralisée<sup>1</sup> est une méthode d'analyse de tableaux  $X$  d'ordre  $(n, p)$  qui peuvent s'écrire, comme aux paragraphes précédents, sous la forme :

$$X = (X_1, X_2, \dots, X_k, \dots, X_q)$$

On note encore  $n$  le nombre d'individus ;  $p$  le nombre total de variables, tous groupes confondus ;  $p_k$  le nombre des variables du  $k^{\text{ième}}$  groupe ;  $q$  le nombre de groupes.

L'analyse canonique généralisée a la vertu de contenir comme cas particulier une grande partie des méthodes descriptives et explicatives qui ont été présentées.

Si  $q = 2$ , l'analyse de  $X = (X_1, X_2)$  coïncide avec l'analyse canonique des deux groupes. On a vu qu'à ce titre, elle contient comme cas particulier l'analyse discriminante (cas où l'un des deux blocs est un tableau disjonctif complet) et donc l'analyse des correspondances des tables de contingence ( $X_1$  et  $X_2$  sont tous deux disjonctifs complets).

<sup>1</sup> L'analyse canonique généralisée a été présentée dans Horst (1961), où elle figure au troisième rang parmi quatre généralisations possibles de l'analyse canonique. Elle a été reprise ou développée par Carrol (1968) dont le nom est souvent attaché à la méthode, Kettenring (1971), Saporta (1975 a), Masson (1980). Casin et Turlot (1986) ont montré qu'elle pouvait être considérée comme une analyse discriminante particulière, et en déduisent des règles d'interprétation nouvelles. Ballif (1986) a développé sous le nom d'AMDG (Analyse multivariée descriptive généralisée) toute une méthodologie de traitement de données, pour laquelle la notion de variable est remplacée par celle plus large de "SEHO" (sous-espace homogène d'observables) et où l'analyse canonique généralisée joue un rôle central.

Toujours si  $q = 2$ , en tant qu'analyse canonique classique, elle contient également la régression multiple (si par exemple  $X_1$  n'a qu'une seule colonne), et donc l'analyse de la variance et de la covariance ( $X_2$  disjonctif complet ou mixte, après régularisation).

Si  $q \geq 2$ , et si chaque bloc  $X_k$  est un tableau disjonctif complet, l'analyse canonique généralisée n'est autre que l'analyse des correspondances multiples de  $X$ . Enfin, toujours si  $q \geq 2$ , si chaque bloc  $X_k$  n'est formé que d'une seule colonne ( $p_k = 1$  pour tout  $k$ ), elle n'est autre que l'analyse en composantes principales normée de  $X$ .

### a – Formulation générale

L'analyse canonique généralisée a déjà été présentée au § 3.1.3.b dans le cas particulier où les blocs sont des tableaux disjonctifs complets. Il convient de donner ici une formulation plus générale, qui puisse englober toutes les méthodes précitées.

Dans l'espace  $\mathbb{R}^n$ , où les  $p$  variables (colonnes de  $X$ ) sont des points, on désigne par  $V_k$  le sous-espace engendré par les colonnes de  $X_k$ .

La projection  $y_k$  d'une variable  $y$  quelconque (point de  $\mathbb{R}^n$ ) sur le sous-espace  $V_k$  s'écrit, si les colonnes de  $X_k$  sont linéairement indépendantes<sup>1</sup> :

$$y_k = X_k(X_k'X_k)^{-1}X_k'y = P_k y \quad [3.8 - 2]$$

Remarquons que si  $\mathbb{R}^n$  était muni d'un produit scalaire associé à une matrice diagonale  $M$ , la formule précédente s'écrirait sous la forme plus générale :

$$y_k = X_k(X_k'MX_k)^{-1}X_k'M y = Q_k y$$

où l'opérateur idempotent de projection  $Q_k$  sur  $V_k$  vaut<sup>2</sup>:

$$Q_k = X_k(X_k'MX_k)^{-1}X_k'M$$

Ce cadre plus général alourdirait les notations sans changer la substance de l'exposé, qui se poursuivra donc avec  $M = I$ , comme dans la formule [3.8 - 2].

Le carré du cosinus de  $y$  avec  $V_k$  (et donc de  $y$  avec  $P_k y$ ) que l'on notera  $R^2(y, k)$  s'écrit :

$$R^2(y, k) = \frac{y'P_k y}{y'y} = \frac{y'X_k(X_k'X_k)^{-1}X_k'y}{y'y} \quad [3.8 - 3]$$

<sup>1</sup> Si les  $p_k$  colonnes de  $X_k$  ne sont pas linéairement indépendantes, il suffit de les remplacer par les  $r_k$  colonnes de  $V$  correspondant à des valeurs propres non nulles dans la décomposition aux valeurs singulières de  $X$  qui s'écrit :  $X = V \Lambda^{1/2} U$ .

<sup>2</sup> Alors que  $P_k$  est symétrique, l'opérateur-projection  $Q_k$  est  $M$ -symétrique, c'est-à-dire que l'on a la relation :  $MQ_k = Q_k'M$ .

On définit le premier axe de l'analyse canonique généralisée comme un vecteur  $\mathbf{y}$  tel que la quantité  $s$  :

$$s = \sum_{k=1}^q R^2(\mathbf{y}, k)$$

soit maximale.

Notons que si les  $\mathbf{X}_k$  sont centrés, le coefficient  $R^2(\mathbf{y}, k)$  est le carré du coefficient de corrélation multiple  $R^2(\mathbf{y}, k)$  entre  $\mathbf{y}$  et  $\mathbf{X}_k$ .

Chaque cosinus carré  $R^2(\mathbf{y}, k)$  est une mesure de proximité entre le vecteur  $\mathbf{y}$  et le sous-espace  $V_k$  engendré par les colonnes de  $\mathbf{X}_k$ . La maximisation du critère  $s$  fait en sorte que le vecteur  $\mathbf{y}$  soit le plus près possible de l'ensemble des groupes de variables.

Il s'agit donc de rendre maximale la somme :

$$s = \sum_{k=1}^q \mathbf{y}' \mathbf{X}_k (\mathbf{X}_k' \mathbf{X}_k)^{-1} \mathbf{X}_k' \mathbf{y}$$

avec la contrainte :  $\mathbf{y}' \mathbf{y} = 1$

Le vecteur  $\mathbf{y}$  de  $\mathbb{R}^n$  sera donc le vecteur propre correspondant à la plus grande valeur propre  $\lambda$  de la matrice  $\mathbf{S}$  d'ordre  $(n, n)$  :

$$\mathbf{S} = \sum_{k=1}^q \mathbf{X}_k (\mathbf{X}_k' \mathbf{X}_k)^{-1} \mathbf{X}_k' \quad [3.8 - 4]$$

Les axes suivants s'obtiennent en rendant maximal le même critère  $s$ , avec la même contrainte de norme, et des contraintes d'orthogonalité par rapport à l'ensemble des axes précédents.

### b – Propriétés de l'Analyse Canonique Généralisée

On va montrer successivement que l'analyse canonique (et donc tout l'éventail des méthodes qui en sont des cas particuliers), l'analyse en composantes principales normée et l'analyse des correspondances multiples sont des cas particuliers de l'Analyse canonique généralisée.

- Pour  $q = 2$ , l'analyse canonique généralisée est une analyse canonique classique.

L'équation donnant  $\mathbf{y}$  s'écrit, pour  $q = 2$  :

$$\mathbf{X}_1 (\mathbf{X}_1' \mathbf{X}_1)^{-1} \mathbf{X}_1' \mathbf{y} + \mathbf{X}_2 (\mathbf{X}_2' \mathbf{X}_2)^{-1} \mathbf{X}_2' \mathbf{y} = \lambda \mathbf{y} \quad [3.8 - 5]$$

Posons<sup>1</sup> :  $(\mathbf{X}_1' \mathbf{X}_1)^{-1} \mathbf{X}_1' \mathbf{y} = \mathbf{a}$  et  $(\mathbf{X}_2' \mathbf{X}_2)^{-1} \mathbf{X}_2' \mathbf{y} = \mathbf{b}$ .

La relation [3.8 - 5] devient simplement :

<sup>1</sup> On note que  $\mathbf{a}$  et  $\mathbf{b}$  sont les vecteurs de coefficients de régression de  $\mathbf{y}$  expliqué respectivement par  $\mathbf{X}_1$  et  $\mathbf{X}_2$ .

$$\mathbf{X}_1 \mathbf{a} + \mathbf{X}_2 \mathbf{b} = \lambda \mathbf{y} \quad [3.8-6]$$

Prémultiplions ensuite les deux membres de la relation [3.8 - 6] par  $(\mathbf{X}'_1 \mathbf{X}_1)^{-1} \mathbf{X}'_1$ , il reste :

$$(\mathbf{X}'_1 \mathbf{X}_1)^{-1} \mathbf{X}'_1 \mathbf{X}_2 \mathbf{b} = (\lambda - 1) \mathbf{a} \quad [3.8 - 7]$$

On obtient de la même façon, en prémultipliant les deux membres de la relation [3.8-6] par  $(\mathbf{X}'_2 \mathbf{X}_2)^{-1} \mathbf{X}'_2$  :

$$(\mathbf{X}'_2 \mathbf{X}_2)^{-1} \mathbf{X}'_2 \mathbf{X}_1 \mathbf{a} = (\lambda - 1) \mathbf{b} \quad [3.8 - 8]$$

On obtient finalement, par substitution :

$$(\mathbf{X}'_2 \mathbf{X}_2)^{-1} \mathbf{X}'_2 \mathbf{X}_1 (\mathbf{X}'_1 \mathbf{X}_1)^{-1} \mathbf{X}'_1 \mathbf{X}_2 \mathbf{b} = (\lambda - 1)^2 \mathbf{b}$$

La matrice à diagonaliser n'est autre que celle donnée par la formule [3.1 - 4] du paragraphe 3.1.2.a. On note également la relation entre valeurs propres :  $\beta = \lambda - 1$ .

*- Si  $q \geq 2$  et si les blocs ne comportent chacun qu'une colonne (centrée), l'analyse canonique généralisée est une analyse en composantes principales normée.*

Dans ce cas, on a  $p_k = 1$  pour tout  $k$ , et donc  $p = q$ . On peut maintenant réécrire la formule [3.8 - 4], les  $\mathbf{X}_k$  étant des vecteurs notés  $\mathbf{x}_k$  :

$$\mathbf{S} = \sum_{k=1}^q \mathbf{x}_k (\mathbf{x}'_k \mathbf{x}_k)^{-1} \mathbf{x}'_k = \sum_{k=1}^q \frac{1}{n s_k^2} \mathbf{x}_k \mathbf{x}'_k \quad [3.8 - 9]$$

où :

$$s_k^2 = \frac{1}{n} \mathbf{x}'_k \mathbf{x}_k$$

est la variance empirique de la variable  $k$ .

Si l'on considère la matrice  $\mathbf{T}$  des variables centrées réduites dont la  $k$ ième colonne vaut  $t_k = \frac{1}{s_k} \mathbf{x}_k$ , la matrice  $\mathbf{S}$  prend la forme  $\mathbf{S} = \frac{1}{n} \mathbf{T} \mathbf{T}'$ .

La relation  $\mathbf{S} \mathbf{y} = \lambda \mathbf{y}$  s'écrit alors, en posant  $\mathbf{T} \mathbf{y} = \mathbf{u}$  et en prémultipliant ses deux membres par  $\mathbf{T}'$  :

$$\frac{1}{n} \mathbf{T}' \mathbf{T} \mathbf{u} = \lambda \mathbf{u}$$

soit finalement :

$$\mathbf{C} \mathbf{u} = \lambda \mathbf{u}$$

où  $\mathbf{C}$  est la matrice des corrélations d'ordre  $(p,p)$  correspondant au tableau  $\mathbf{X}$  initial.

Cette présentation a le mérite d'enrichir l'interprétation de l'analyse en composantes principales normée, qui peut être définie comme la recherche

d'une variable artificielle ( $y$ ) qui rend maximale la somme de ses corrélations avec toutes les variables actives.

- Pour  $q \geq 2$ , quand les blocs sont des tableaux disjonctifs complets, l'analyse canonique généralisée est une analyse des correspondances multiples.

Pour retrouver (partiellement) les notations de la section 1.4, changeons les  $X$  en  $Z$ . Posons donc  $Z = X$  et  $Z_k = X_k$  et posons également  $D_k = Z'_k Z_k$ .

$D_k$  est la matrice diagonale d'ordre  $(p_k, p_k)$  correspondant aux marges (sommés des colonnes) du tableau  $Z_k$ . Enfin appelons  $D$  la matrice diagonale d'ordre  $(p, p)$  dont les  $q$  blocs diagonaux sont les  $D_k$ .

L'équation  $Sy = \lambda y$  s'écrit :

$$\sum_{k=1}^q Z_k D_k^{-1} Z'_k y = \lambda y \quad [3.8 - 10]$$

Posons, pour tout entier positif  $h \leq q$ ,  $Z'_h y = u_h$ , ce qui revient également à écrire  $Zy = u$ ,  $u$  étant un vecteur à  $p$  composantes tel que :

$$u' = (u'_1, u'_2, \dots, u'_h, \dots, u'_q)$$

Prémultipliant les deux membres de [3.8 - 10] par  $Z'_h y$ , on peut alors écrire, pour  $h = 1, \dots, q$  :

$$\sum_{k=1}^q Z'_h Z_k D_k^{-1} u_k = \lambda u_h \quad [3.8 - 11]$$

Ces  $q$  équations ne sont autres qu'une écriture par bloc de la relation matricielle :

$$Z' Z D^{-1} u = \lambda u$$

Cette formule est à rapprocher de la formule [1.4 - 1] du § 1.4.3 b, où le paramètre  $s$  est ici noté  $q$  (nombre de tableaux  $X_k$ ). Avec les notations du présent paragraphe, l'équation de l'analyse des correspondances multiples s'écrit :

$$\frac{1}{q} Z' Z D^{-1} u = \lambda' u, \quad \text{d'où : } \lambda = q \lambda'$$

La valeur propre issue de l'analyse canonique généralisée est  $q$  fois plus grande que celle issue de l'analyse des correspondances multiples du même tableau global  $Z$ .

- Pour  $q \geq 2$  dans le cas général, l'analyse canonique généralisée est une analyse générale du tableau  $X$  dans une métrique que l'on peut qualifier de "Mahalanobis par bloc"

Le raisonnement tenu à propos de l'analyse des correspondances multiples (sous-paragraphe précédent ci-dessus) s'applique dans le cas où  $X_k$  est centré, mais quelconque.

La formule [3.8 - 11] prend alors la forme :

$$\sum_{k=1}^q \mathbf{X}'_h \mathbf{X}_k (\mathbf{X}'_k \mathbf{X}_k)^{-1} \mathbf{u}_k = \lambda \mathbf{u}_h \quad [3.8 - 12]$$

Si l'on appelle  $\mathbf{D}$  la matrice diagonale par bloc d'ordre  $(p,p)$  ( $\mathbf{D}$  a  $q^2$  blocs dont  $q$  blocs diagonaux) dont le  $k^{\text{ième}}$  bloc diagonal est :

$$\mathbf{D}_{kk} = (\mathbf{X}'_k \mathbf{X}_k)^{-1}$$

$\mathbf{D}_{kk}$  est la matrice associée à la distance de Mahalanobis interne au groupe  $k$  (cf. § 3.3.4.c).

Les  $q$  formules [3.8 - 12] (pour  $h = 1, \dots, q$ ), s'écrivent :

$$\mathbf{X}' \mathbf{X} \mathbf{D}^{-1} \mathbf{u} = \lambda \mathbf{u}$$

Ce qui établit le résultat annoncé (cf. § 1.1.6).

### c – Utilisation en pratique de l'analyse canonique généralisée

L'analyse canonique généralisée peut s'utiliser comme analyse de compromis dans des approches de type STATIS ou analyse factorielle multiple. Elle n'utilise cependant que les sous-espaces correspondant à chaque groupe, et non la structure interne des nuages dans ces sous-espaces. Ceci peut entraîner les mêmes difficultés d'interprétation que l'analyse canonique.

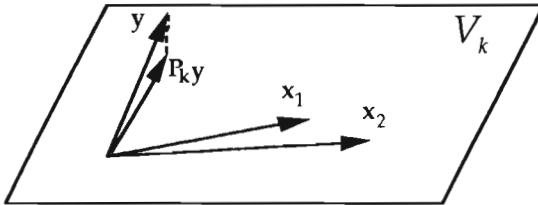


Figure 3.8 - 1

Exemple montrant les insuffisances du coefficient  $R^2(\mathbf{y}, \mathbf{k})$

La figure 3.8 - 1 (cf. Escofier et Pagès, 1988) met ainsi en évidence une faiblesse classique du coefficient de corrélation multiple. Elle montre deux vecteurs  $\mathbf{x}_1$  et  $\mathbf{x}_2$  contenus dans le sous-espace  $V_k$ , et un vecteur  $\mathbf{y}$ , proche du sous espace  $V_k$ , donc proche de sa propre projection  $\mathbf{P}_k \mathbf{y}$  sur  $V_k$ .  $R^2(\mathbf{y}, \mathbf{k})$  est donc voisin de 1, alors que  $\mathbf{y}$  est presque orthogonal à  $\mathbf{x}_1$  et  $\mathbf{x}_2$ .

C'est ce type de difficulté qui a conduit ces auteurs à proposer, pour l'analyse globale de  $\mathbf{X}$ , une métrique diagonale par blocs (le  $k^{\text{ième}}$  bloc  $\mathbf{D}_{kk}$  étant lui-même diagonal et tel que  $\mathbf{D}_{kk} = (1/\lambda_k) \mathbf{I}$ , au lieu de  $\mathbf{D}_{kk} = (\mathbf{X}'_k \mathbf{X}_k)^{-1}$  dans le cas de l'analyse canonique généralisée).

Les cas particuliers pour lesquels l'analyse canonique généralisée, dans le cas  $q \geq 2$ , rejoint des méthodes dont l'interprétation est aisée, sont précisément ceux qui excluent une situation telle que celle de la figure précédente (mauvaise base du sous-espace  $V_k$ ).

En analyse en composantes principales, les  $V_k$  n'ont qu'une dimension, donc  $R^2(y, k)$  est un carré de coefficient de corrélation classique avec la variable  $x_k$  correspondante.

En analyse des correspondances multiples, le codage disjonctif complet fait que chaque  $X_k$  est une base orthogonale du sous-espace  $V_k$  correspondant.

On pourrait penser qu'une généralisation qui n'est utile que dans des cas particuliers n'a pas d'intérêt pour le praticien. On peut en fait aménager l'analyse canonique généralisée en la "régularisant" (cf. § 3.2.5 et § 3.3.6), c'est-à-dire en remplaçant chaque tableau  $X_k$  par le tableau des axes issus d'une analyse en axes principaux de  $X_k$  (qui aura moins de  $p_k$  colonnes s'il y a des colinéarités, ou des quasi-colinéarités, c'est-à-dire des valeurs propres faibles).

Ceci rejoint, en d'autres termes, la démarche de Ballif (*op. cit.*) qui conçoit l'analyse canonique généralisée (désignée, on l'a vu, par AMDG) comme une synthèse d'analyses (c'est-à-dire de sous-espaces stables) plutôt que de tableaux. Le principal intérêt de la méthode est alors de pouvoir traiter simultanément des tableaux très hétérogènes<sup>1</sup>.

Notons que Escofier (1979 *b*) a abordé directement ce problème dans un cas particulier en considérant (sans nommer l'analyse canonique généralisée) un tableau mixte  $X$  (qualitatif-quantitatif) contenant deux sortes de blocs : soit des variables continues isolées  $x_k$ , soit des tableaux disjonctifs complets. Cet auteur a établi un résultat que l'on peut exprimer de cette façon : en remplaçant chaque colonne  $x_k$  de terme général  $x_{ik}$  par un bloc de deux colonnes de termes généraux  $(1 - x_{ik})/2$  et  $(1 + x_{ik})/2$ , il est équivalent de procéder à l'analyse des correspondances de  $X$  ou à l'analyse canonique généralisée de  $X$ , formé des nouveaux blocs .

D'autres propriétés de l'analyse canonique généralisée sont présentées dans les articles cités en début de paragraphe.

---

<sup>1</sup> Les blocs  $V_k$  formés de plusieurs variables nominales sont prétraités par analyse des correspondances multiples, les blocs formés de plusieurs variables continues par analyse en composantes principales, les blocs formés de tables de contingence par analyse des correspondances simple.





## Chapitre 4

---

# **VALIDITÉ ET PORTÉE DES RÉSULTATS**



## Introduction

Au cours des chapitres précédents, on a vu comment fonctionnaient les méthodes de statistique descriptive multidimensionnelle de base (méthodes choisies parmi les plus utilisées) et on a vu à quel point celles-ci pouvaient parfois être proches, dans leur principe mais aussi dans les circonstances de leur utilisation, des méthodes plus explicatives.

Alors que la statistique descriptive élémentaire (unidimensionnelle) n'est qu'une technique de présentation des données (jugée sur ses qualités de fidélité, de précision, d'ergonomie, parfois d'esthétique), les méthodes d'analyse des données produisent en fait plus que des représentations. Elles dévoilent des traits structuraux, permettent d'observer une éventuelle hétérogénéité des données, d'éprouver leur cohérence. Elles supposent une connaissance du domaine étudié, une maîtrise de règles d'interprétation assez complexes, mais ne se réduisent pourtant pas à des tests d'hypothèses ou des validations de modèles.

Devant les résultats d'une analyse factorielle ou d'une classification, on est naturellement conduit à poser un certain nombre de questions sur la qualité des représentations :

- Observe-t-on vraiment quelque chose? Les données ont-elle une structure? Ou, au contraire, de simples fluctuations d'échantillonnage suffiraient-elles à expliquer les valeurs obtenues pour les valeurs propres et les taux d'inertie?
- Les premiers axes principaux indiquent-ils des directions caractéristiques? Les premières valeurs propres sont-elles hautement significatives? Comment apprécier la position d'un point dans l'espace factoriel? Que représente le taux d'inertie en terme d'information?
- A-t-on découvert des classes préexistantes ou au contraire a-t-on découvert une réalité continue en classes ?
- Est-ce que les configurations obtenues sont stables, compte tenu de ce que l'on sait sur la précision des données, la nature du codage, et l'importance relative des différentes variables? Quelle est l'influence sur les résultats d'une modification du tableau de données (ajout ou retrait de certains éléments, modification du codage) ?

Les trois sections qui constituent ce chapitre vont tenter d'apporter des éléments de réponses à chacune de ces questions.

Nous envisagerons tout d'abord le cas des techniques d'axes principaux (méthodes factorielles) : l'analyse en composantes principales, l'analyse des correspondances et ses extensions, l'analyse factorielle discriminante, pour

lesquelles les valeurs propres et les taux d'inertie sont des paramètres permettant de caractériser de façon globale (mais pas simple) les espaces de représentation (section 4.1).

Puis nous présenterons les méthodes de validation plus partielles ou empiriques (calculs de stabilité, zones de confiance, validation par rééchantillonnage) qui concernent aussi bien les méthodes factorielles que les méthodes de classification les plus usuelles (section 4.2).

## Signification des valeurs propres et des taux d'inertie

Pour toute analyse en axes principaux, qu'il s'agisse d'analyse en composantes principales ou d'analyse des correspondances, l'hypothèse d'indépendance des lignes et des colonnes d'un tableau est en général une hypothèse trop sévère pour être réaliste. Il est en effet extrêmement improbable qu'un tableau soumis à l'analyse puisse être aussi dépourvu de structure qu'une table de nombres au hasard.

Bien qu'étant un cas extrême d'une portée pratique limitée, l'hypothèse d'indépendance va cependant nous permettre de définir des *seuils de signification* pour les valeurs propres et les pourcentages d'inertie, qui joueront le rôle de *garde-fou* pour les utilisateurs.

La grande variété des tableaux analysables (tableaux de mesure, de classements, de comptage, etc.) rend extrêmement délicate l'interprétation de ces valeurs propres et des taux d'inertie correspondant, dont on sait qu'ils sont étroitement liés au codage des données.

Sous l'hypothèse d'indépendance des lignes et des colonnes du tableau analysé, les valeurs propres suivent des *lois paramétriques* dans le cas de l'analyse en composantes principales<sup>1</sup>, des *lois non-paramétriques* dans le cas de l'analyse des rangs<sup>2</sup> et de l'analyse des correspondances des tableaux de contingence<sup>3</sup>.

Dans ces situations favorables, il a été possible de procéder à des *tabulations approchées*, et de tracer des *abaques* qui les résument. Nous traiterons principalement le cas de l'analyse des correspondances<sup>4</sup>.

### 4.1.1 Travaux sur la loi des valeurs propres en analyse en composantes principales

La loi de Wishart, établie par Fisher (1915) dans le cas  $p = 2$ , puis par Wishart (1928), généralise la loi du  $\chi^2$ .

---

<sup>1</sup> Il est nécessaire de spécifier la forme analytique de la distribution des variables - loi normale - et d'estimer les paramètres correspondants.

<sup>2</sup> Cf. §1.2.8; la loi de la matrice de corrélation des rangs sous l'hypothèse d'indépendance ne suppose que la continuité des distributions des variables.

<sup>3</sup> Comme dans le test du  $\chi^2$  appliqué aux tables de contingence, la normalité résulte de la convergence de la loi multinomiale vers la loi normale.

<sup>4</sup> On trouvera des abaques approchées relatives à l'analyse des rangs dans Lebart et Fénelon (1971).

La loi du  $\chi^2$ , dans un cadre paramétrique classique, définit la distribution d'une variance empirique sous l'hypothèse d'observations indépendantes identiquement distribuées suivant une loi normale de moyenne nulle et d'écart-type  $\sigma$  connu. La loi de Wishart concerne la distribution d'une matrice des covariances empiriques.

Si les  $n$  vecteurs-lignes d'une matrice  $\mathbf{X}$  d'ordre  $(n, p)$  sont des réalisations indépendantes d'un vecteur multinormal de moyenne théorique nulle, et de matrice des covariances théoriques  $\Sigma$  (non singulière) alors la matrice  $\mathbf{S} = \mathbf{X}'\mathbf{X}$  (qui contient  $p(p+1)/2$  éléments distincts) suit une loi de Wishart, notée  $W(p, n, \Sigma)$  dont la densité  $f(\mathbf{S})$  est donnée par la formule :

$$f(\mathbf{S}) = C(n, p, \Sigma) |\mathbf{S}|^{-\frac{n-p-1}{2}} \exp\left\{-\frac{1}{2} \text{trace}(\Sigma^{-1}\mathbf{S})\right\},$$

la constante  $C(n, p, \Sigma)$  ayant pour valeur :

$$C(n, p, \Sigma) = 2^{-\frac{np}{2}} |\Sigma|^{-\frac{n}{2}} \pi^{-\frac{p(p-1)}{4}} \prod_{k=1}^p \Gamma\left(\frac{1}{2}(n+1-k)\right)$$

On vérifie que pour  $\Sigma = \mathbf{I}$  (matrice unité) et  $p = 1$ , notant  $s = \mathbf{x}'\mathbf{x}$ , on retrouve la densité de probabilité du  $\chi^2$ . En effet :

$$f(s) = C(n, 1, \mathbf{I}) s^{\frac{n}{2}-1} \exp\left\{-\frac{s}{2}\right\}$$

avec :

$$C(n, 1, \mathbf{I}) = 2^{-\frac{n}{2}} \Gamma\left(\frac{n}{2}\right)$$

La loi de la matrice  $\mathbf{S}$  (loi de Wishart<sup>1</sup>) intervient dans l'établissement des tests intervenant en analyse de la variance multidimensionnelle et en analyse discriminante. C'est le cas pour le test d'égalité de plusieurs matrices de covariances [*test de Box*], test d'égalité de vecteurs moyens [*test du A de Wilks*], etc. (cf. Saporta, 1990).

La densité de probabilité des *valeurs propres* issues d'une matrice de Wishart a été explicitée simultanément par Fisher (1939), Girshick (1939), Hsu (1939) et Roy (1939), puis par Mood (1951). On trouve une démonstration donnant la forme de cette densité dans Anderson (1958), Muirhead (1982).

<sup>1</sup> Pour l'établissement de la densité de la loi de Wishart et de certaines lois dérivées, cf. Dugué (1958), Anderson (1958), Muirhead (1982). On note que  $\mathbf{S}$  n'est pas la matrice des covariances empiriques puisque les variables ne sont pas centrées sur la moyenne empirique de l'échantillon. On montre (cf. références ci-dessus) que la loi de  $\mathbf{S}$  après centrage empirique est une loi  $W(p, n-1, \Sigma)$ .

Dans le cas où  $\Sigma = \mathbf{I}$ , la densité de la loi de Wishart s'écrit facilement en fonction de la trace et du déterminant de  $\mathbf{S}$ , c'est-à-dire de la somme et du produit des valeurs propres  $\lambda_k$  :

$$f(\mathbf{S}) = C(n, p, \mathbf{I}) \left( \prod_{k=1}^p \lambda_k \right)^{-\frac{n-p-1}{2}} \exp \left\{ -\frac{1}{2} \sum_{k=1}^p \lambda_k \right\},$$

On retrouvera donc ces éléments (multipliés par le jacobien de la transformation qui est ici le produit de toutes les différences possibles entre valeurs propres) dans l'expression de la densité  $g(\Lambda)$  des valeurs propres :

$$g(\Lambda) = D(n, p) \left( \prod_{k=1}^p \lambda_k \right)^{-\frac{n-p-1}{2}} \exp \left\{ -\frac{1}{2} \sum_{k=1}^p \lambda_k \right\} \prod_{k < j} (\lambda_k - \lambda_j)$$

La constante  $D(n, p)$  ayant pour valeur :

$$D(n, p) = 2^{-\frac{np}{2}} \pi^{\frac{p}{2}} \prod_{k=1}^p \left\{ \Gamma \left( \frac{1}{2}(n+1-k) \right) \Gamma \left( \frac{1}{2}(p+1-k) \right) \right\}$$

L'intégration de cette densité assez complexe a donné lieu à plusieurs publications ; parmi les principales, celles de Pillai (1965), Krishnaiah et Chang (1971), qui s'inspirent des travaux du physicien Mehta (1960, 1967) <sup>1</sup>.

Les distributions ci-dessus s'appliquent à des variables indépendantes de variance théorique égale à 1 (l'hypothèse de moyenne nulle n'est pas nécessaire puisqu'il suffit de travailler avec la matrice des covariances centrées, et de changer  $n$  en  $n-1$  dans la loi de  $\mathbf{S}$ ). Il n'est donc pas facile d'utiliser ces résultats dans les applications usuelles de l'analyse en composantes principales <sup>2</sup>.

#### 4.1.2 Approximation de la distribution des valeurs propres en analyse des correspondances

Nous allons voir que la distribution des valeurs propres en analyse des correspondances sous l'hypothèse d'indépendance des lignes et des

<sup>1</sup> Une table des seuils correspondant aux deux valeurs propres extrêmes a été publiée par Choudary, Hanumara et Thompson (1968) pour des matrices ayant leur plus petit côté  $p$  inférieur à 10 ; par Pillai (1967), Pillai et Chang (1970) et par Clemm, Krishnaiah et Waikar (1973) pour  $p \leq 20$ .

<sup>2</sup> Le fait de réduire les variables ( $\mathbf{X}'\mathbf{X}$  est alors  $n$  fois la matrice des corrélations) ne résout pas le problème car  $\mathbf{X}'\mathbf{X}$ , dont la diagonale est constante et égale à  $n$ , ne suit évidemment pas une loi de Wishart. Les éléments diagonaux d'une matrice de Wishart  $W(p, n-1, \mathbf{I})$  sont en effet des réalisations de  $\chi^2$  à  $n-1$  degrés de liberté.



colonnes peut être approchée par celle des valeurs propres d'une matrice dont la loi est connue (matrice de Wishart)<sup>1</sup>.

Nous reprenons ici les notations du chapitre 1, section 1.3, sur l'analyse des correspondances. L'entier  $k_{ij}$  est le terme général de la table de contingence  $K$  à  $n$  lignes et  $p$  colonnes. On note ici encore :

$$k = \sum_i \sum_j k_{ij} \quad \text{et} \quad f_{ij} = \frac{k_{ij}}{k}$$

Si  $p_{ij}$  désigne la probabilité correspondant à la case  $(i,j)$  et estimée par  $f_{ij}$ , et si l'on note les marges théoriques  $p_i$  et  $p_j$ , l'hypothèse d'indépendance des lignes et des colonnes se traduit par la relation :

$$p_{ij} = p_i \cdot p_j$$

Ainsi  $k_{ij}$  est l'une des  $np$  composantes d'un vecteur multinomial, dont l'espérance mathématique  $E(k_{ij})$  s'écrit :

$$E(k_{ij}) = k p_i \cdot p_j$$

On fera une approximation analogue à celle qui est faite lors de l'établissement de la loi du  $\chi^2$  pour tester l'indépendance des lignes et des colonnes d'un tableau de contingence :  $k$  sera supposé suffisamment grand pour permettre l'utilisation de l'approximation normale de la loi multinomiale.

On considérera d'autre part que les marges observées  $f_i$  et  $f_j$  peuvent être substituées sans dommage aux marges théoriques  $p_i$  et  $p_j$  sans toutefois négliger les contraintes impliquées par cette substitution.

Ces hypothèses permettront d'ailleurs de retrouver le test classique du  $\chi^2$  sur les tables de contingence.

Désignons par  $h$  le vecteur à  $np$  composantes tel que :

$$h_{ij} = \frac{\sqrt{k}(f_{ij} - f_i \cdot f_j)}{\sqrt{f_i \cdot f_j}}$$

<sup>1</sup> Cf. Lebart (1975 b, 1976), Corsten (1976), et dans le cas d'hypothèses plus générales O'Neill (1978, 1981). La loi des valeurs propres issues de l'analyse des correspondances a donné lieu à maintes publications erronées. Ainsi dans le traité classique de statistique de Kendall et Stuart (1961), les valeurs propres sont supposées suivre, comme l'inertie totale, des lois du  $\chi^2$ . Lancaster (1963, 1969) a réfuté ce résultat en montrant que l'espérance mathématique de la première valeur propre est toujours supérieure aux valeurs découlant des assertions de Kendall et Stuart. Les références concernant d'autres approximations peuvent être trouvées dans l'ouvrage de Kshirsagar (1972), où il est suggéré que les valeurs propres, étant des coefficients de corrélations canoniques calculés sur des variables disjonctives (cf. au chapitre 3 du présent ouvrage les sections 3.1.3.a et 3.3.4.b) pourraient suivre une loi très proche de celle de ces mêmes coefficients calculés cette fois sur des variables gaussiennes. Des simulations montrent que cette approximation n'est pas satisfaisante.

Ce vecteur de  $\mathbb{R}^{np}$  a, sous les conditions précédentes, une distribution normale avec  $E(h_{ij}) = 0$  pour tout  $i$  et  $j$ .

Sa matrice des covariances a pour terme général:

$$V_h(ij, i'j') = \delta_{(ij, i'j')} - \sqrt{f_i \cdot f_j \cdot f_{i'} \cdot f_{j'}}$$

où

$$\begin{aligned} \delta_{(ij, i'j')} &= 1 \text{ si } i = i' \text{ et } j = j' \\ \delta_{(ij, i'j')} &= 0 \text{ sinon} \end{aligned}$$

Construisons une matrice orthogonale  $\mathbf{A}$ , d'ordre  $(p, p)$ , telle que sa première colonne ait pour  $j^{\text{ème}}$  élément  $\sqrt{f_j}$  (pour tout  $j \leq p$ ), les  $p-1$  autres colonnes formant avec la première une base orthonormée de  $\mathbb{R}^p$ .

De la même façon, construisons une matrice orthogonale  $\mathbf{B}$  d'ordre  $(n, n)$  telle que sa première ligne ait pour  $i^{\text{ème}}$  élément  $\sqrt{f_i}$  (pour tout  $i \leq n$ ), les  $n-1$  autres lignes formant avec la première une base orthonormée de  $\mathbb{R}^n$ .

La matrice  $\mathbf{B} \otimes \mathbf{A}'$  d'ordre  $(np, np)$ , produit *direct* ou de *Kronecker* des matrices  $\mathbf{B}$  et  $\mathbf{A}'$ , est aussi orthogonale.

Pour tous  $1 < i < n$ ,  $1 < j < p$ ,  $1 < r < n$  et  $1 < s < p$ , on a les relations :

$$\sum_j \sqrt{f_j} h_{ij} = 0 ; \sum_i \sqrt{f_i} h_{ij} = 0 ; \sum_m b_{rm} \sqrt{f_m} = 0 ; \sum_k a_{ks} \sqrt{f_k} = 0 ;$$

De ces relations, on déduit que le vecteur  $\mathbf{y}$  de  $\mathbb{R}^{np}$  tel que:

$$\mathbf{y} = \mathbf{B} \otimes \mathbf{A}' \mathbf{h}$$

a seulement  $(n-1)(p-1)$  composantes non nulles. On a:

$$y_{rs} = 0 \text{ si } r = 1 \text{ ou si } s = 1$$

La matrice des covariances de  $\mathbf{y}$  est :

$$\mathbf{V}_y = (\mathbf{B} \otimes \mathbf{A}') \mathbf{V}_h (\mathbf{B}' \otimes \mathbf{A})$$

Pour tout couple de composantes non nulles, on a:

$$V_y(r s, r' s') = \delta_{rr'} \delta_{ss'}$$

Soit  $\mathbf{Y}$  la matrice d'ordre  $(n, p)$  définie par :

$$\mathbf{Y} = \mathbf{B} \mathbf{H} \mathbf{A}$$

où  $\mathbf{H}$  est la matrice d'ordre  $(n, p)$  de terme général  $h_{ij}$ . La première ligne et la première colonne de  $\mathbf{Y}$  sont nulles.

Les éléments de la sous-matrice  $\hat{\mathbf{Y}}$  d'ordre  $(n-1, p-1)$ , formée des éléments non nuls de  $\mathbf{Y}$ , sont donc distribués indépendamment suivant la loi normale centrée réduite.

La matrice :

$$S = \hat{Y}'\hat{Y}$$

est donc distribuée suivant une loi de Wishart  $W(p-1, n-1, I)$  de paramètres  $(n-1)$  et  $(p-1)$ .

Or  $S$  a les mêmes valeurs propres non nulles que  $Y'Y$  c'est-à-dire que  $A'H'HA$  ; ce sont finalement les mêmes valeurs propres que  $H'H$ , puisque  $A$  est orthogonale.

Remarquons que ceci implique que  $tr(H'H)$  est un  $\chi^2$  à  $(n-1)(p-1)$  degrés de liberté. Or :

$$tr(H'H) = k \sum_i \sum_j \frac{\sqrt{f_{ij} - f_i \cdot f_j}}{f_i \cdot f_j}$$

Il s'agit du test usuel du  $\chi^2$  sur les tableaux de contingence.

La matrice symétrisée  $S^*$  que l'on diagonalise lors de l'analyse des correspondances du tableau  $K$ , est la matrice :

$$S^* = \frac{1}{k} H'H$$

Ainsi, si  $\lambda_\alpha$  est la  $\alpha^{\text{ème}}$  valeur propre issue de l'analyse des correspondances d'un tableau  $K$  d'ordre  $(n, p)$ , de somme totale  $k$ , alors la distribution de  $k\lambda_\alpha$  est approximativement celle de la  $\alpha^{\text{ème}}$  valeur propre d'une matrice de Wishart définie par les paramètres  $W(p-1, n-1, I)$ <sup>1</sup>.

### 4.1.3 Indépendance des taux d'inertie et de la trace

On a vu que la densité  $g(\Lambda)$  de la loi jointe des valeurs propres  $\lambda_1, \lambda_2, \dots, \lambda_p$  d'une matrice de Wishart a la forme :

$$g(\Lambda) = D(n, p) \prod_k \lambda_k^{\frac{n-p-1}{2}} \exp\left\{-\frac{1}{2} \sum_k \lambda_k\right\} \prod_{k < j} (\lambda_k - \lambda_j)$$

Si l'on pose :

$$\begin{cases} \lambda_k = z \tau_k, \text{ pour } k < p \\ \lambda_p = z (1 - \tau_1 - \tau_2 - \dots - \tau_{p-1}) \end{cases}$$

alors  $z$  est la trace de la matrice de Wishart :  $z = \sum_k \lambda_k$

<sup>1</sup> On trouvera une vérification expérimentale de la qualité de l'approximation montrant la concordance entre les lois théoriques des valeurs propres et celles qui résultent de l'approximation ci-dessus dans Lebart (1975 b, 1976).

On trouve aisément une factorisation de la densité (le jacobien de cette transformation vaut  $z^{p-1}$ ) :

$$g(\Lambda) = g_1(z) g_2(\tau_1, \dots, \tau_{p-1})$$

la fonction  $g_1(z)$  s'écrivant :

$$g_1(z) = \frac{1}{2\Gamma(\frac{np}{2})} \left(\frac{z}{2}\right)^{\frac{np}{2}-1} \exp\left(-\frac{z}{2}\right)$$

où l'on reconnaît la densité de la loi du  $\chi^2$  à  $np$  degrés de liberté. La factorisation des densités (et l'indépendance des domaines d'intégration) montrent que les pourcentages de variance  $\tau_1, \tau_2, \dots, \tau_{p-1}$  sont indépendants de la trace  $z$ .

Cette propriété (qui suppose vraie l'hypothèse d'indépendance) semble encore valable dans le cas de l'analyse des correspondances, pour laquelle la loi de Wishart est seulement une loi approchée (les simulations extensives entreprises pour construire les abaques ont permis de vérifier cette indépendance, que nous avons d'ailleurs conjecturée à partir de résultats empiriques, puis démontrée). Elle avait en fait été établie (dans le cadre de l'analyse en composantes principales) par Bartlett (1951).

En analyse des correspondances, la trace mesure la dilatation générale du nuage de points-profiles, alors que les taux d'inertie mesurent la forme du nuage en termes d'aplatissement et d'allongement. Ainsi, même si la trace ne permet pas de rejeter l'hypothèse d'indépendance (test habituel du  $\chi^2$ ), les premiers taux d'inertie pourront néanmoins être significativement élevés : l'analyse des correspondances pourra être utile même sur les tableaux que le  $\chi^2$  ne désigne pas comme étant très riches d'informations (nuage peu dilaté mais non-sphérique de points-profiles).

Inversement, à une trace significativement élevée pourront correspondre des taux d'inertie non significatifs. Bien que l'hypothèse d'indépendance soit rejetée par le test du  $\chi^2$ , l'analyse des correspondances n'est peut-être pas alors le meilleur outil pour décrire la dépendance entre les lignes et les colonnes de la table (nuage dilaté sphérique de points profils).

Ces situations ont été schématisées à la section 1.3 (analyse des correspondances), par la figure 1.3 - 14 du paragraphe 1.3.4.a : les taux d'inertie significatifs ne concernent que la seconde colonne de cette figure (formes non-sphériques), alors que les  $\chi^2$  significatifs ne concernent que la seconde ligne de la figure (forte inertie correspondant à des nuages dilatés).

Ainsi, le modèle de l'analyse des correspondances, que l'on peut schématiser ci-dessous (avec les relations et contraintes entre  $\alpha, \beta, \varphi, \psi, \lambda$  qui sont les relations et contraintes usuelles de la formule de reconstitution des données) :

$$f_{ij} \approx \alpha_i \beta_j \left( 1 + \sum_{h=1}^m \sqrt{\lambda_h} \varphi_h(i) \psi_h(j) \right)$$

n'est pas à abandonner chaque fois que le  $\chi^2$  ne permet pas de rejeter l'hypothèse d'indépendance, contrairement à la plupart des modélisations concernant les tables de contingence.

#### 4.1.4 Exemples d'abaques et tables statistiques

Les tables statistiques établies par simulation et les abaques qui en résultent permettent d'apprécier le degré de signification de la plus grande valeur propre issue de l'analyse des correspondances de tableaux de contingence depuis la dimension (6×6) jusqu'à la dimension (50×100).

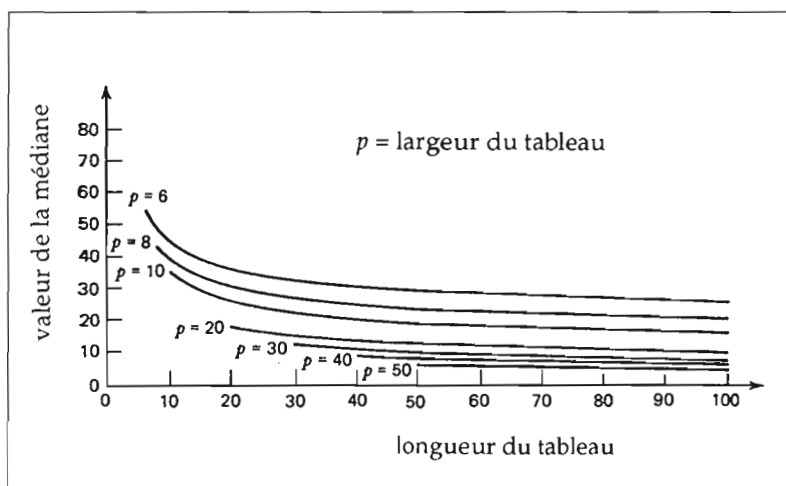


Figure 4.1 - 1  
Valeurs médianes du pourcentage d'inertie  
de la plus grande valeur propre

La figure 4.1 - 1 donne les valeurs médianes du pourcentage d'inertie relatif à la plus grande valeur propre pour les largeurs  $p = 6, 8, 10, 20, 30, 40, 50$ <sup>1</sup>.

Les estimations des valeurs des taux d'inertie correspondant à la première valeur propre apparaissent sur la figure 4.1 - 2 (pour un seuil de 0.05). Les

<sup>1</sup> Des informations plus détaillées concernant la construction de ces abaques (notamment sur les modes de génération de tableaux pseudo-aléatoires) et des *tables approchées*, pour les tableaux dont les dimensions n'excèdent pas 50×100 relatives aux cinq premières valeurs propres sont données dans Lebart (1975a). Les simulations ont mis en jeu des tables de contingences pseudo-aléatoires ayant des marges théoriques et un effectif total donnés en utilisant l'approximation normale de la loi multinomiale. Des expériences ont montré en effet qu'on obtient de cette façon des résultats comparables, en ce qui concerne les valeurs propres, à une procédure ayant recours à une simulation pseudo-aléatoire du schéma multinomial exact.

extrémités des courbes (points :  $6 \times 6$ ,  $8 \times 8$ ,  $10 \times 10$ ) ont été établies à l'aide de 1000 simulations (100 pour les autres points), afin de préciser leur tracé. Ces figures schématiques ne donnent cependant que des ordres de grandeurs.

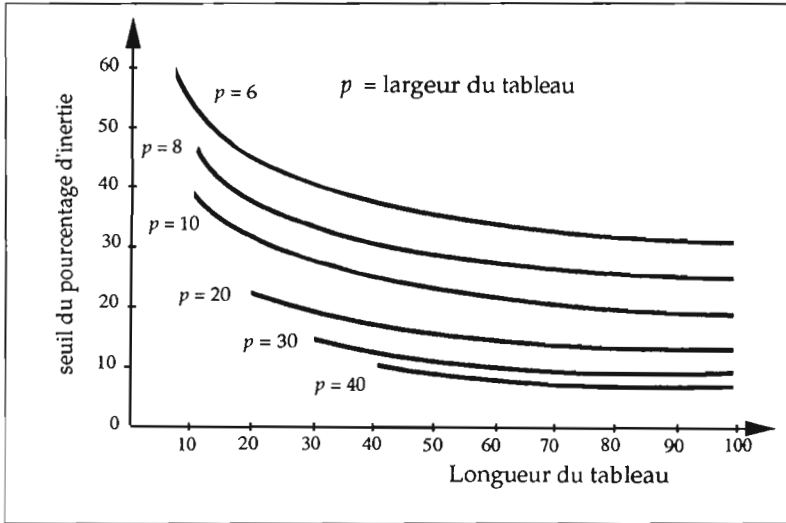


Figure 4.1 - 2  
Seuil (0,05 unilatéral) du pourcentage d'inertie  
de la plus grande valeur propre

Par exemple, on lit sur la figure 4.1 - 2 que, pour un tableau  $10 \times 10$ , la première valeur propre peut atteindre ou dépasser 40% de l'inertie (la loi des taux ne dépendant pas de l'effectif total du tableau) dans 5% des cas, sous l'hypothèse d'indépendance des lignes et des colonnes de la table.

Il s'agit donc ici d'un test de sphéricité du nuage de points-profiles, qui ne remplace pas un test sur les valeurs propres elles-mêmes (il faut alors tabuler  $k\lambda_1$ , car la loi de  $\lambda_1$  seule dépend de  $k$ , effectif total de la table). Ce test donne néanmoins des ordres de grandeur ayant une certaine valeur pédagogique sur l'effet de fluctuations d'échantillonnage sur la forme de nuage de points-profiles.

En revanche, ce type de résultats étendu à l'ensemble des valeurs propres sous l'hypothèse d'indépendance ne peut aider à déterminer le nombre d'axes à retenir, car les valeurs propres ne sont pas indépendantes (même sous l'hypothèse d'indépendance des lignes et des colonnes du tableau, et *a fortiori* si cette hypothèse est rejetée)<sup>1</sup>. Il faudrait donc connaître la loi

<sup>1</sup> Cette forte dépendance entre valeurs propres dans le cas d'une table de contingence générée sous l'hypothèse d'indépendance des lignes et des colonnes se manifeste sous la forme d'autocorrélations entre valeurs propre consécutives et de corrélations négatives entre valeurs propres extrêmes. Ce type de liaison est illustré par la figure 2.4 - 7 du chapitre 2, qui représente graphiquement les corrélations entre les 7 valeurs propres d'une table de contingence aléatoire (8,8) après 1000 simulations.

conditionnelle de la seconde valeur propre, ce qui ne peut donner lieu à des résultats généraux ni à des tables.

### 4.1.5 Taux d'inertie et information

L'utilisation des taux d'inertie (ou pourcentages de variance) comme outil d'évaluation globale de la qualité d'une représentation est très délicate. Les taux d'inertie sont des mesures pessimistes de la qualité d'une représentation (contrairement, par exemple, aux coefficients de corrélation multiple qui sont des mesures optimistes de la qualité d'une régression). La variance brute initiale n'étant pas en général une mesure de référence adéquate, il est souvent injustifié de parler de *part d'information* à propos des *taux d'inertie*.

#### a – Caractère partiel des taux d'inertie

Quelques contre-exemples vont nous montrer que ces coefficients ne sont pas aptes à caractériser de façon satisfaisante la qualité d'une représentation.

##### *Cas du codage disjonctif*

Nous avons vu<sup>1</sup> que, pour une même représentation, l'analyse de deux questions (ou variables) sous codage disjonctif pouvait donner des taux d'inertie considérablement plus faibles que l'analyse, pourtant équivalente, du tableau de contingence croisant les deux variables.

Les taux d'inertie donnent dans ce cas une idée très *pessimiste* de la part d'information représentée. En effet, le codage disjonctif, en introduisant une orthogonalité entre les colonnes (modalités) relatives à une même question, introduit une sorte de sphéricité artificielle du nuage de points-profiles, que l'on retrouve dans la forme du spectre.

Comme cela a été évoqué dans une note au § 1.4.8, Benzécri (1979) a proposé une formule de calcul de taux d'inertie  $\tau(\lambda)$  corrigés sous la forme suivante :

$$\tau(\lambda) = \left(\frac{s}{s-1}\right)^2 \left(\lambda - \frac{1}{s}\right)^2 \quad \text{pour } \lambda > \frac{1}{s}$$

où  $s$  représente le nombre de questions actives,  $\lambda$  représente la valeur propre issu de l'analyse des correspondances du tableau disjonctif complet, ( $\lambda^2$  étant la valeur propre issue de l'analyse des correspondances du tableau de Burt).

---

<sup>1</sup> Cf. § 1.4.6 consacré au cas de deux questions en analyse des correspondances multiples, notamment le tableau 1.4 - 1 et les remarques qui suivent.

Les valeurs propres issues du tableau de Burt dont la diagonale a été annulée sont précisément  $(\lambda - \frac{1}{s})^2$  et seulement celles qui vérifient  $\lambda - \frac{1}{s}$  correspondent à des facteurs directs (cf. paragraphe ci dessous : cas de l'analyse de la matrice associée à un graphe symétrique).

De plus, dans le cas  $s = 2$ , on retrouve les taux d'inertie de l'analyse des correspondances de la vraie table de contingence croisant les deux questions<sup>1</sup>.

### Cas de l'analyse de la matrice associée à un graphe symétrique

Dans plusieurs cas lors de l'analyse des correspondances de la matrice associée à un graphe symétrique (cf. Benzécri, 1973, tome IIB, chapitre 10), un calcul analytique exact peut être fait sans recours à l'ordinateur. Il est alors intéressant d'étudier analytiquement les variations des représentations en fonction des différents codages de la matrice associée.

La relation de transition s'écrit ici :

$$\frac{1}{2} \mathbf{M} \varphi = \varepsilon(\varphi) \sqrt{\lambda} \varphi$$

où  $\mathbf{M}$  est la matrice associée au graphe et  $\varepsilon(\varphi) = 1$  ou  $-1$  selon la parité du facteur  $j$ , c'est-à-dire selon que le facteur est direct<sup>2</sup> ou inverse.

Examinons par exemple le cas de l'analyse d'un cycle simple. La matrice  $\mathbf{M}$  n'a que deux éléments non nuls (égaux à 1) par ligne et par colonne.

Désignons par  $n$  le nombre de sommets du graphe. Pour  $n = 5$ , on a par exemple :

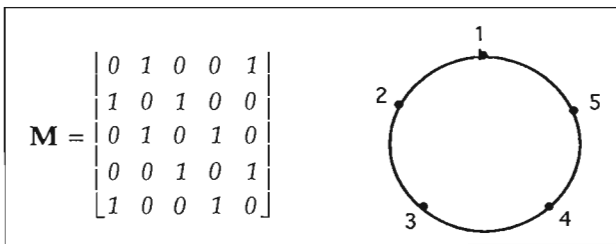


Figure 4.1 - 3  
Exemple de cycle et de sa matrice associée

<sup>1</sup> Dans le cas de l'exemple numérique du paragraphe 1.4.8, le taux correspondant à la première valeur propre (22.77%) devient alors 64%. Greenacre (cf. Greenacre et Blasius, 1994) propose une modification itérative du tableau de Burt qui conduit à des représentations très similaires, mais à des taux intermédiaires entre les taux bruts et les taux rectifiés (sous le nom de *Joint Correspondence Analysis*).

<sup>2</sup> Brièvement ici, un facteur direct est un vecteur propre de  $\mathbf{M}'\mathbf{M} = \mathbf{M}^2$  correspondant à une valeur propre positive de  $\mathbf{M}$ . Seuls les facteurs directs traduisent des similarités.



La relation précédente s'écrit encore pour  $1 < j < n$  :

$$\frac{1}{2}(\varphi(j-1) + \varphi(j+1)) = \varepsilon(\varphi)\sqrt{\lambda}\varphi(j)$$

Les solutions de ce type classique d'équation aux différences finies sont, compte tenu des conditions aux limites :

$$\varphi_\alpha(j) = \cos\left(\frac{2j\alpha\pi}{n}\right) \quad \text{et} \quad \psi_\alpha(j) = \sin\left(\frac{2j\alpha\pi}{n}\right)$$

Ce sont les  $j^{\text{èmes}}$  composantes des deux facteurs associés à la valeur propre double :

$$\lambda_\alpha = \cos^2\left(\frac{2\alpha\pi}{n}\right)$$

On obtient dans le plan des deux premiers facteurs l'équation paramétrique d'un cercle et donc une reconstitution satisfaisante de la structure dont le tableau **M** représente un codage particulier.

La trace de la matrice à diagonaliser s'écrit :

$$\text{tr} \frac{1}{4} \mathbf{M}^2 = \frac{n}{2}$$

Le taux d'inertie correspondant à l'axe  $\alpha$  est donc :

$$\tau_\alpha = \frac{2}{n} \cos^2\left(\frac{2\alpha\pi}{n}\right)$$

Le résultat, en apparence paradoxal, est le suivant :

le taux d'inertie du sous-espace qui "restitue" la structure initiale peut être rendu aussi petit que l'on veut, pourvu de choisir un cycle assez long : si  $n = 10^3$ , alors  $\tau_1 \approx 2 \times 10^{-3}$ .

#### **- Influence du choix des variables en analyse en composantes principales**

Si l'on complète un tableau à  $n$  lignes et  $p$  colonnes, par  $q$  nouvelles colonnes formées de nombres pseudo-aléatoires, l'analyse en composantes principales normées du nouveau tableau à  $p+q$  colonnes donnera les mêmes premiers axes (s'ils prédominent) que l'analyse du tableau initial.

Les pourcentages de variance expliquée seront cependant plus faibles (car la trace qui valait  $p$ , vaut maintenant  $p+q$ ). Pourtant la part d'information dont les axes rendent compte reste naturellement la même.

En pratique, on est dans une situation analogue lorsque le nombre potentiel des variables est très grand (cas par exemple de la présence d'espèces animales ou végétales dans les relevés écologiques). Une certaine discipline dans le choix du recueil des données, dictée par les critères d'homogénéité, devrait en principe permettre d'éviter ces inconvénients.

Mais le statisticien n'a pas toujours la maîtrise de la collecte des données, ni une connaissance suffisante du domaine d'application ; de plus, les critères de choix sont eux-mêmes trop qualitatifs et trop généraux pour définir de façon rigoureuse un tableau optimal parmi tous les tableaux potentiels.

Comme les procédures de codage, le choix proprement-dit des variables a souvent davantage d'influence sur les taux d'inertie que sur les facteurs issus d'une analyse.

### b – Quelle information?

Nous allons voir que la théorie de l'information de Shannon-Wiener (Kullback, 1959) ne nous permet pas de faire apparaître facilement les taux d'inertie comme mesure du degré de "non-sphéricité" d'un nuage.

On utilisera pour le calcul la notion de *divergence* de Jeffreys (1946), qui permet de mesurer la distance entre deux hypothèses  $H_1$  et  $H_2$  dans le cas d'une réalisation d'un vecteur  $x$  issu de l'un des deux schémas relatifs à des lois normales dans  $\mathbb{R}^p$  :

( $H_1$ ) Hypothèse d'indépendance :

$$\begin{cases} \text{Moyenne théorique} = \mu_1 \\ \text{Matrice des covariances théoriques} = \sigma^2 \mathbf{I} \end{cases}$$

( $H_2$ ) Cas général :

$$\begin{cases} \text{Moyenne théorique} = \mu_2 \\ \text{Matrice des covariances théoriques} = \Sigma \text{ (supposée ici régulière)} \end{cases}$$

La divergence va permettre d'exprimer la distance entre les hypothèses  $H_1$  et  $H_2$  en fonction des valeurs propres de  $\Sigma$  et l'on s'apercevra qu'elle met en cause les petites valeurs propres alors que l'analyse factorielle ne retient que les grandes.

Rappelons que l'on définit, pour deux hypothèses  $H_1$  et  $H_2$  pouvant donner lieu à la réalisation d'un vecteur  $x$ , la divergence  $J(H_1, H_2)$  comme la différence :

$$J(H_1, H_2) = \int \log \frac{P(H_1|x)}{P(H_2|x)} dv_1(x) - \int \log \frac{P(H_2|x)}{P(H_1|x)} dv_2(x)$$

$v_1$  et  $v_2$  étant les mesures associées aux hypothèses  $H_1$  et  $H_2$  ; et  $P(H_i|x)$  ( $i = 1, 2$ ) étant la probabilité conditionnelle que  $H_i$  soit vraie connaissant  $x$ .

Dans le cas de densités continues  $f_1(x)$  et  $f_2(x)$ , on a :

$$J(H_1, H_2) = \int (f_1(x) - f_2(x)) \log \frac{f_1(x)}{f_2(x)} dx$$

La densité de probabilité du vecteur  $x$  s'écrit, pour une matrice des covariances théoriques  $\Sigma_i$  et un vecteur  $\mu_i$  :

$$f_i(x) = \frac{1}{\sqrt{2\pi|\Sigma_i|}} \exp\left\{-\frac{1}{2}(x - \mu_i)' \Sigma_i^{-1}(x - \mu_i)\right\}$$

d'où :

$$\log \frac{f_1(\mathbf{x})}{f_2(\mathbf{x})} = \frac{1}{2} \log \frac{|\Sigma_2|}{|\Sigma_1|} - \frac{1}{2} \text{tr}(\Sigma_1^{-1}(\mathbf{x} - \mu_1)(\mathbf{x} - \mu_1)') + \frac{1}{2} \text{tr}(\Sigma_2^{-1}(\mathbf{x} - \mu_2)(\mathbf{x} - \mu_2)')$$

Remplaçant cette valeur dans la formule donnant  $J(H_1, H_2)$ , on voit que le premier terme de  $J(H_1, H_2)$  n'est autre que  $I_{(1;2)}$ , l'information moyenne apportée par l'échantillon  $\mathbf{x}$  sous l'hypothèse  $H_1$ , en vue de discriminer en faveur de  $H_1$  contre  $H_2$  (cf. Kullback, 1959).

On écrira ce premier terme en posant :  $\mathbf{x} - \mu_2 = \mathbf{x} - \mu_1 + \mu_1 - \mu_2$ . Il vient :

$$\begin{aligned} I_{(1;2)} &= \int f_1(\mathbf{x}) \log \frac{f_1(\mathbf{x})}{f_2(\mathbf{x})} d\mathbf{x} \\ &= \frac{1}{2} \log \frac{|\Sigma_1|}{|\Sigma_2|} + \frac{1}{2} \text{tr}(\Sigma_1(\Sigma_2^{-1} - \Sigma_1^{-1})) + \frac{1}{2} \text{tr}(\Sigma_2^{-1}(\mu_1 - \mu_2)(\mu_1 - \mu_2)') \end{aligned}$$

et  $J(H_1, H_2)$  s'écrit donc :

$$\begin{aligned} J(H_1, H_2) &= I_{(1;2)} + I_{(2;1)} \\ &= \frac{1}{2} \text{tr}((\Sigma_1 - \Sigma_2)(\Sigma_2^{-1} - \Sigma_1^{-1})) + \frac{1}{2} \text{tr}((\Sigma_1^{-1} - \Sigma_2^{-1})(\mu_1 - \mu_2)(\mu_1 - \mu_2)') \end{aligned}$$

On s'intéresse au cas<sup>1</sup> pour lequel :

$$\mu_1 = \mu_2, \quad \Sigma_1 = \mathbf{I} \quad \text{et} \quad \Sigma_2 = \Sigma.$$

On notera en abrégé  $J(H_1, H_2) = J(\mathbf{I}, \Sigma)$  avec :

$$J(\mathbf{I}, \Sigma) = \frac{1}{2} \text{tr}((\mathbf{I} - \Sigma)(\Sigma^{-1} - \mathbf{I})) = \frac{1}{2} \text{tr}(\Sigma + \Sigma^{-1}) - p$$

Soit, en faisant apparaître les valeurs propres  $\lambda_\alpha$  de  $\Sigma$  :

$$J(\mathbf{I}, \Sigma) = \frac{1}{2} \left( \sum_{\alpha=1}^p \lambda_\alpha + \sum_{\alpha=1}^p \frac{1}{\lambda_\alpha} \right) - p$$

Si les inerties totales théoriques sont égales sous les hypothèses  $H_1$  et  $H_2$ , on a la relation :

$$\sum_{\alpha=1}^p \lambda_\alpha = p$$

Le seul terme variable dans  $J(\mathbf{I}, \Sigma)$  est donc le terme :

$$\sum_{\alpha=1}^p \frac{1}{\lambda_\alpha}$$

<sup>1</sup> Dans le cas où  $\Sigma_1 = \Sigma_2$ ,  $J(H_1, H_2)$  est proportionnel à la distance de Mahalanobis, ou distance généralisée entre les populations théoriques 1 et 2 (cf. §3.3.4).

On voit que la divergence entre les deux hypothèses sera particulièrement grande dans le cas où certaines valeurs propres de  $\Sigma$  sont voisines de 0.

Dans le cadre de cette formalisation de la théorie de l'information, une valeur propre de  $\Sigma$  infiniment petite jouera un rôle beaucoup plus déterminant que deux valeurs propres expliquant, par exemple, 80% de l'inertie totale, alors que c'est dans le sous-espace des deux facteurs correspondants que l'on observera les principaux traits structuraux.

En fait, comme un *filtre* dans un processus de communication, la représentation des données dans l'espace des premiers axes factoriels a pour effet d'augmenter la *valeur pratique de l'information* au prix d'une perte d'information brute qui peut être considérable. Or cette notion de valeur pratique (Brillouin, 1959) est étrangère à la théorie classique de l'information.

Comme le suggère Thom(1974), on gagnerait souvent à remplacer le mot *information* par le mot *forme* (ici à peu près équivalent au mot anglais *pattern*) lors d'un processus d'observation.

Le meilleur critère de validation sera de vérifier la *stabilité des formes* obtenues à l'issue de cette phase de nos analyses.

#### 4.1.6 Choix du nombre d'axes : quelques résultats utiles

On a vu la difficulté de l'étude distributionnelle et de la signification statistique des valeurs propres et des taux d'inertie.

Cependant qu'il s'agisse d'une simple visualisation des données ou de l'utilisation des axes factoriels en vue d'une analyse ultérieure (classification sur facteurs, régression ou analyse discriminante sur facteurs), il reste important de savoir combien d'axes retenir, autrement dit de connaître la dimension de l'espace de représentation.

Il existe quatre types de procédures pour guider le choix de ce nombre d'axes :

- a - des règles empiriques.
- b - des procédures externes.
- c - des critères fondés sur certaines propriétés statistiques des valeurs propres.
- d - des méthodes de calcul de stabilité, de rééchantillonnage ou de simulation.

Nous évoquerons ici brièvement les points (a) et (b), puis traiterons essentiellement le point (c), dans la ligne des développements de cette section. Le point (d) fera l'objet de la section 4.2 suivante.

### a – Règles empiriques

Les règles empiriques sont fondées sur l'allure de la séquence des valeurs propres, notamment dans le cas de l'analyse en composantes principales<sup>1</sup>. Deux règles, attribuées respectivement à Cattell et Kaiser, seront citées, à titre historique.

Lorsqu'un tableau est généré suivant un modèle stipulant l'indépendance de ses lignes et de ses colonnes, on observe une décroissance régulière des valeurs propres. Cette remarque est à l'origine de procédures empiriques pour juger du nombre d'axes à retenir dans une analyse factorielle. On étudie l'histogramme de décroissance des valeurs propres pour y déceler un changement de pente. Chaque fois que l'histogramme des valeurs propres présente un "décrochage" ou une discontinuité, on peut supposer que quelque chose de non aléatoire intervient. Ce repérage des "coudes" était préconisé par Cattell (1966).

Le second critère empirique est le critère de Kaiser (1961), qui stipule de ne retenir que les valeurs propres supérieures à la moyenne des valeurs propres (c'est-à-dire à 1 dans le cas d'une analyse en composantes principales sur matrices de corrélation), en s'appuyant notamment sur des travaux de Guttman (1954)<sup>2</sup>. D'un emploi très répandu à cause de son extrême simplicité, il peut être facilement mis en défaut. Ainsi, une analyse en composantes principales sur matrice des corrélations en biométrie peut produire un facteur de taille très dominant. Comme la trace est constante, les autres valeurs propres sont condamnées à être très petites, ce qui pourrait interdire l'interprétation d'autres dimensions.

### b – Procédures externes

Les procédures externes sont fondées sur des connaissances extérieures au tableau de données (interprétabilité de certains résultats, informations apportées par le positionnement de certaines variables supplémentaires<sup>3</sup>). Elles sont justifiées par la faible pertinence des valeurs propres et des taux d'inertie, soulignée dans la section 4.1.5 précédente, notamment dans le cas des codages clairsemés (tableaux disjonctifs complets, matrices associées à des graphes, certains tableaux de présence-absence).

De telles procédures externes peuvent être couplées avec les procédures de rééchantillonnage dont on parlera dans la section 4.2. Ainsi, dans le cas

---

<sup>1</sup> Ces procédures ont surtout concerné en pratique l'analyse factorielle classique ou analyse des psychologues, dont l'analyse en composantes principales est un cas particulier - cas des variances spécifiques égales ou nulles (cf. § 3.2.9).

<sup>2</sup> Cf. également, parmi de très nombreuses publications sur ce thème, les articles de synthèse de Anastassakos et d'Aubigny (1984), Francisco et Finch (1980) et la revue faite par Jolliffe (1986).

<sup>3</sup> C'est le cas notamment en analyse des correspondances multiples quand les modalités d'une variable nominale supplémentaire possède des valeurs-test très significatives sur un ou plusieurs axes.

d'une analyse discriminante à partir des facteurs d'une analyse, on peut sélectionner le nombre de facteurs à partir des pourcentages de bien classés (ou pourcentages de succès) calculés sur des échantillons-tests ou par validation croisée (cf. l'exemple numérique du paragraphe 3.3.6.b, ainsi que la figure 3.3 - 7 de la section 3.3).

C'est donc ici le pouvoir de prédiction sur la variable externe qui permet de choisir la dimension de l'espace des prédicteurs<sup>1</sup>. Les procédures externes jouent en fait un rôle important dans la méthodologie de l'analyse des données.

### c – Critères de choix statistiques, résultats asymptotiques

Les travaux relatifs aux études distributionnelles des valeurs propres et des vecteurs propres ainsi qu'aux comportements asymptotiques de ces éléments sont considérables, mais peu de résultats sont vraiment utilisables par le praticien. Sauf mention contraire, tous les résultats ci-dessous supposent que les observations suivent une loi multinormale de matrice des covariances théoriques  $\Sigma$ .

#### - Critère pour l'analyse en composantes principales

Dans sa publication déjà citée (§ 4.1.1) donnant l'expression de la densité des valeurs propres d'une matrice de Wishart, Girshick (1939) calcule les variances et covariances asymptotiques (quand le nombre d'observations  $n$  tend vers l'infini) des valeurs propres et vecteurs propres de la matrice des covariances expérimentales  $S$ , ceci dans le cas où la matrice des covariances théoriques  $\Sigma$  a toutes ses valeurs propres distinctes.

Il donne également les variances et covariances théoriques des valeurs propres de la matrice des corrélations expérimentales, lorsque la matrice de corrélation théorique  $R$  a également toutes ses valeurs propres distinctes.

Bartlett (1950, 1951) propose une méthode pour tester l'égalité de  $p - q$  valeurs propres des matrices  $\Sigma$  ou  $R$ . Lawley (1956) approfondit le cas des  $p - q$  plus petites valeurs propres de  $\Sigma$ .

Anderson (1963) a généralisé ces résultats, en déterminant les lois limites des valeurs propres sans nécessairement supposer que les valeurs théoriques correspondantes sont distinctes.

Il démontre en particulier, pour tester l'égalité des  $r$  plus petites valeurs propres  $\hat{\lambda}_\alpha$  de la matrice des covariances expérimentales  $S$ , que la statistique :

---

<sup>1</sup> Plus généralement, ce type de procédure permet de sélectionner un sous-espace qui n'est pas nécessairement engendré par des axes consécutifs.

$$X^2 = nr \log \frac{\left(\frac{1}{r}\right) \sum_{\alpha=p-r+1}^{\alpha=p} \hat{\lambda}_\alpha}{\left(\prod_{\alpha=p-r+1}^{\alpha=p} \hat{\lambda}_\alpha\right)^{1/r}}$$

( $nr$  fois le logarithme du rapport de la moyenne arithmétique des  $r$  plus petites valeurs propres à leur moyenne géométrique) est asymptotiquement distribué comme un  $\chi^2$  à  $\left[\frac{r(r+1)}{2} - 1\right]$  degrés de liberté.

Les intervalles de confiance asymptotiques d'Anderson utilisés en pratique pour les valeurs propres remontent en fait aux travaux de Girshick.

Si les valeurs propres théoriques  $\lambda_\alpha$  de  $\Sigma$  sont distinctes, les valeurs propres  $\hat{\lambda}_\alpha$  de la matrice des covariances empiriques  $S$  suivent asymptotiquement des lois normales d'espérance  $\lambda_\alpha$  et de variance  $2\lambda_\alpha^2/(n-1)$  où  $n$  est la taille de l'échantillon.

On en déduit les intervalles de confiance approchés au seuil 95% :

$$\lambda_\alpha \in \left[ \hat{\lambda}_\alpha \left(1 - 1.96\sqrt{2/(n-1)}\right) ; \hat{\lambda}_\alpha \left(1 + 1.96\sqrt{2/(n-1)}\right) \right]$$

L'ampleur de l'intervalle donne une indication sur la stabilité de la valeur propre vis-à-vis des fluctuations dues à l'échantillonnage supposé laplacien. L'empiétement des intervalles de deux valeurs propres consécutives suggérera donc l'égalité de ces valeurs propres. Les axes correspondants sont alors définis à une rotation près. Ainsi l'utilisateur pourra éviter d'interpréter un axe instable selon ce critère.

On peut trouver des généralisations de ces résultats asymptotiques au cas non laplacien (Waternaux, 1976; Davis, 1977), mais leur utilisation n'est guère pratique<sup>1</sup>.

Les intervalles de confiance d'Anderson concernent en fait aussi bien les valeurs propres des matrices des covariances que des matrices de corrélations. Les simulations entreprises montrent que les intervalles de confiance obtenus sont en général "prudent" : le pourcentage de couverture de la vraie valeur est le plus souvent supérieur au seuil de confiance annoncé (Morineau, 1983).

Dans tous les cas, la nature asymptotique des résultats et l'hypothèse sous-jacente de normalité<sup>2</sup> font considérer les résultats comme indicatifs.

<sup>1</sup> On trouvera une revue de résultats asymptotiques relatifs à l'analyse en composantes principales dans Muirhead (1982), Anderson (1958 - seconde édition de 1984), Flury (1988), Pousse (1992).

<sup>2</sup> Muirhead (1982) a montré que l'hypothèse d'existence des quatre premiers moments pour la loi théorique de l'échantillon suffisait pour valider ces intervalles.

## - Critère pour l'analyse des correspondances

Dans le cas de l'analyse des correspondances des vraies tables de contingence (analyse des correspondances simples), la loi des valeurs propres ne permet que de juger la signification du premier axe, puisque les lois conditionnelles des autres valeurs propres ne sont pas connues. Une procédure approchée proposée par Malinvaud (1987) peut être utilisée pour déterminer le rang à partir duquel les valeurs propres ne sont plus significativement différentes entre elles.

Revenons au "modèle" que représente la formule de reconstitution approchée avec  $m$  facteurs (cf. formule [1.1 - 7] du § 1.1.5.b) des fréquences relatives  $f_{ij}$  du tableau de contingence  $\mathbf{K}$  d'ordre  $(n, p)$  de terme général  $k_{ij}$  et d'effectif total  $k$  ( $f_{ij} = k_{ij} / k$ ).

$$f_{ij} \approx g_{ij} = \alpha_i \beta_j \left( 1 + \sum_{h=1}^m \sqrt{\lambda_h} \varphi_h(i) \psi_h(j) \right) \quad [4.1 - 1]$$

Les restrictions suivantes sont imposées aux différents paramètres du modèle (moyenne nulle, variance unité, et orthogonalité des facteurs) :

- (a)  $\sum_{i=1}^n \alpha_i \varphi_h(i) = \sum_{j=1}^p \beta_j \psi_h(j) = 0$ , pour  $h \leq m$
- (b)  $\sum_{i=1}^n \alpha_i \varphi_h^2(i) = \sum_{j=1}^p \beta_j \psi_h^2(j) = 1$ , pour  $h \leq m$
- (c)  $\sum_{i=1}^n \alpha_i \varphi_h(i) \varphi_{h'}(i) = \sum_{j=1}^p \beta_j \psi_h(j) \psi_{h'}(j) = 0$ , pour  $h \neq h'$
- (d)  $\lambda_h \geq 0$ , pour  $h \leq m$

On note que  $g_{ij} = f_{ij}$  si  $m = p - 1$  (modèle dit saturé). Il ressort de la première ligne de contraintes que  $\alpha_i \beta_j = f_i \cdot f_j$ .

Le tableau reconstitué dans le cas où  $m = 0$  correspond à l'indépendance entre les lignes et les colonnes du tableau. Pour savoir si cette hypothèse est rejetée, on calcule la statistique du  $\chi^2$  usuelle (à  $(n-1)(p-1)$  degrés de liberté) :

$$\chi^2 = k \sum_{j=1}^p \sum_{i=1}^n \left( \frac{f_{ij} - f_i \cdot f_j}{f_i \cdot f_j} \right)^2$$

Plus généralement, pour  $m \geq 1$  fixé, l'ajustement du modèle [4.1 - 1] en rendant minimal le critère sous les contraintes précédentes :

$$u = \sum_{j=1}^p \sum_{i=1}^n \left( \frac{f_{ij} - g_{ij}}{f_i \cdot f_j} \right)^2$$



fournit la reconstitution  $\hat{g}_{ij}$  de  $f_{ij}$  à partir de  $m$  facteurs.

La statistique  $X^2 = k \sum_{j=1}^p \sum_{i=1}^n \left( \frac{f_{ij} - \hat{g}_{ij}}{\hat{g}_{ij}} \right)^2$ , mesurant l'écart entre le modèle saturé  $f_{ij}$  et le modèle estimé  $\hat{g}_{ij}$ , doit suivre une loi du  $\chi^2$  dont le nombre de degrés de liberté  $d(m) = (n - m - 1)(p - m - 1)$  s'obtient en retranchant de la dimension de l'espace le nombre de paramètres indépendants <sup>1</sup>.

Chaque nouveau facteur (de rang  $m$ ) demande en effet l'estimation de  $n + p + 1$  nouveaux paramètres ( $\varphi_h, \psi_h$  et  $\lambda_h$ ), liés par 2 contraintes de centrage (a), 2 contraintes de normalisation (b),  $(m-1)$  contraintes d'orthogonalité (c), d'où :

$$d(m) = (n - 1)(p - 1) - m(n + p - 2 - m)$$

soit

$$d(m) = (n - m - 1)(p - m - 1)$$

Utilisée avec prudence <sup>2</sup>, cette statistique constitue une aide à l'interprétation préférable aux méthodes empiriques et plus utile que les études de distribution des valeurs propres sous l'hypothèse d'indépendance.

<sup>1</sup> C'est le nombre de paramètres estimés duquel on retranche le nombre de contraintes qui peuvent être écrites sous la forme de fonctions admettant des dérivées partielles du premier ordre continues. (cf. Cramer, 1946 ; Rao, 1973).

<sup>2</sup> Il faut cependant éviter de procéder à des approximations dans le calcul de la statistique  $X^2$ . Il peut en effet être tentant de remplacer le dénominateur  $\hat{g}_{ij}$  par  $f_i \cdot f_j$ , et d'approcher alors  $X^2$  par la somme des valeurs propres de rang supérieur à  $m$ . Cette approximation conduit à comparer la  $m^{\text{ième}}$  valeur propre  $\lambda_m$  à un  $\chi^2$  à  $(n+p-2m-1)$  degrés de liberté. Ce résultat, proposé par plusieurs auteurs, a été réfuté par Lancaster (1963), qui a montré notamment que la plus grande valeur propre sous l'hypothèse d'indépendance a une espérance supérieure à  $(n+p-3)$ .

## Stabilité des axes, des formes, des classes

Cette section, dévolue aux calculs de stabilité et aux méthodes de validation par rééchantillonnage, comporte quatre paragraphes.

Le premier est consacré aux méthodes empiriques de validation qui consistent à modifier certains éléments du tableau initial en fonction des informations externes connues *a priori* sur les mesures, comptages ou codages qui composent ce tableau et à effectuer des calculs de stabilité des résultats. Le deuxième présente les principales méthodes de rééchantillonnage (Jackknife, Bootstrap, validation croisée) et leur application à l'étude de la stabilité des formes. Le troisième paragraphe, notamment en application du paragraphe précédent, décrit les zones de garde ou de confiance que l'on peut tracer autour des points dans les espaces de visualisation. Enfin le quatrième examine le cas de la classification en passant brièvement en revue les travaux relatifs au nombre et à la signification des classes.

### 4.2.1 Méthodes de validation empiriques : calculs de stabilité et de sensibilité

Les calculs de stabilité et de sensibilité sont probablement les procédures de validation les plus probantes. L'essentiel des opérations consiste en une vérification de la stabilité des configurations après diverses perturbations apportées au tableau initial.

#### a – Calculs de stabilité et de sensibilité

D'un point de vue théorique, Escofier et Leroux (1972), Escofier (1979) ont traité de la stabilité des facteurs en analyse en composantes principales et en analyse des correspondances. Ces auteurs étudient les variations maximales des facteurs et des valeurs propres lorsque l'on apporte des modifications bien déterminées aux données : suppression ou ajout d'éléments au tableau de données, influence du regroupement de plusieurs éléments ou de petites modifications des valeurs du tableau, influence du choix de la métrique et de la pondération.

Les sous-espaces correspondant au haut du spectre sont les plus stables vis-à-vis des éventuelles perturbations de la matrice à diagonaliser (cf. Wilkinson, 1965 ; Kato, 1966). De plus cette matrice elle-même (par exemple la matrice des corrélations expérimentales en analyse en composantes principales normée) est moins sensible aux fluctuations d'échantillonnage que les

moments d'ordre 1 (moyennes ou pourcentages)<sup>1</sup>. Ces perturbations ne doivent pas affecter l'orientation des axes ni les configurations si on les suppose stables, et la structure mise en évidence sera alors significative.

Des résultats assez forts existent dans le cas de perturbations symétriques de matrices symétriques, comme par exemple le théorème de Wielandt-Hoffman (cf. Wilkinson, *op.cit.*) qui énonce que si  $\mathbf{A}$ ,  $\mathbf{B}$ ,  $\mathbf{C}$ , sont des matrices symétriques ( $p, p$ ) ayant respectivement pour valeurs propres classées par ordre décroissant  $\alpha_i$ ,  $\beta_i$ ,  $\gamma_i$  et telles que  $\mathbf{C} = \mathbf{A} + \mathbf{B}$  ( $\mathbf{B}$  représente une perturbation additive,  $\mathbf{C}$  est la matrice perturbée), alors :

$$\sum_{i=1}^p (\gamma_i - \alpha_i)^2 \leq \sum_{i=1}^p \beta_i^2$$

Un autre théorème classique très utilisé dans les travaux précités énonce que, avec les mêmes notations, pour tout  $i$  tel que  $1 \leq i \leq p$  :

$$\alpha_i + \beta_p \leq \gamma_i \leq \alpha_i + \beta_1$$

Montrons brièvement, utilisant une formulation empruntée à Gifi (1990), comment des variations de la matrice symétrique à diagonaliser  $\mathbf{A}$ , supposées ici infinitésimales, influencent les éléments propres.

La relation  $\mathbf{A} \mathbf{u}_r = \lambda_r \mathbf{u}_r$  se note, pour l'ensemble du spectre :

$$\mathbf{A} \mathbf{U} = \mathbf{U} \Lambda \quad [4.2 - 1]$$

avec, rappelons-le, les contraintes  $\mathbf{U}'\mathbf{U} = \mathbf{I}$ .

La matrice  $\mathbf{A}$ , et par conséquent  $\mathbf{U}$  et  $\Lambda$ , sont supposés dépendre continûment d'un paramètre  $s$ . La dérivation des relations [4.2 - 1] par rapport à  $s$  donne les systèmes :

$$\begin{aligned} \frac{\partial \mathbf{A}}{\partial s} \mathbf{U} + \mathbf{A} \frac{\partial \mathbf{U}}{\partial s} &= \frac{\partial \mathbf{U}}{\partial s} \Lambda + \mathbf{U} \frac{\partial \Lambda}{\partial s} \\ \text{et } \frac{\partial \mathbf{U}'}{\partial s} \mathbf{U} + \mathbf{U}' \frac{\partial \mathbf{U}}{\partial s} &= 0 \end{aligned} \quad [4.2 - 2]$$

Prémultipliant les deux membres de la première relation par  $\mathbf{U}'$ , il vient après simplification (mettant à profit les relations précédentes) :

$$\mathbf{U}' \frac{\partial \mathbf{A}}{\partial s} \mathbf{U} = \frac{\partial \Lambda}{\partial s} + \left\{ \mathbf{U}' \frac{\partial \mathbf{U}}{\partial s} \Lambda - \Lambda \mathbf{U}' \frac{\partial \mathbf{U}}{\partial s} \right\}$$

Or la matrice entre accolades a ses éléments diagonaux nuls (comme toute matrice de la forme  $(\mathbf{B}\Lambda - \Lambda\mathbf{B})$ , avec  $\Lambda$  diagonale), d'où l'expression de  $\frac{\partial \Lambda}{\partial s}$  :

<sup>1</sup> Les travaux de Tanaka (1984) concernent également l'analyse des correspondances (connue également au Japon sous le nom de *méthode de quantification n°3* de Hayashi). Sur l'analyse en composantes principales, on mentionnera les travaux de Krzanowski (1984), Critchley (1985), Benasseni (1986a, 1986b).

$$\frac{\partial \Lambda}{\partial s} = \text{diag} \left( \mathbf{U}' \frac{\partial \mathbf{A}}{\partial s} \mathbf{U} \right)$$

(où le symbole *diag* signifie "diagonale de...").

Posant :

$$\mathbf{Q} = \mathbf{U}' \frac{\partial \mathbf{A}}{\partial s} \mathbf{U}$$

d'où

$$q_{rt} = \sum_{i,j} u_{ir} u_{jt} \frac{\partial a_{ij}}{\partial s}$$

on peut écrire :

$$\frac{\partial \lambda_r}{\partial s} = q_{rr} = \sum_{i,j} u_{ir} u_{jr} \frac{\partial a_{ij}}{\partial s} \quad [4.2 - 3]$$

Pour les vecteurs propres, le calcul, plus complexe, conduit à :

$$\frac{\partial u_{jr}}{\partial s} = \sum_{t \neq r} u_{jt} \frac{q_{rt}}{\lambda_r - \lambda_t} \quad [4.2 - 4]$$

Ces formules nous montrent donc, d'une part que la partie principale de la variation des valeurs propres ne dépend pas des variations des vecteurs propres (la variance le long d'un axe dépend plus, par exemple, de l'adjonction ou du retrait d'un élément que de petites variations d'angle de l'axe), d'autre part que les variations des composantes d'un vecteur propre dépendent des écarts entre la valeur propre correspondante et les autres valeurs propres, c'est-à-dire de l'isolement de cette valeur propre, résultat également intuitif et récurrent dans tous les calculs de perturbation.

Ainsi, beaucoup des résultats de Escofier et Leroux (*op. cit.*) se fondent sur un théorème que ces auteurs établissent à partir de résultats de Davis et Kahan (1970), qui s'énonce, avec les mêmes notations que pour le théorème de Wielandt-Hoffman :

Soit deux sous-espaces invariants de  $\mathbf{A}$  et de  $\mathbf{C}$  ( $\mathbf{C} = \mathbf{A} + \mathbf{B}$ ) correspondant à des valeurs propres de mêmes rangs  $r, r+1, \dots, r+t$ .

Si  $\theta$  est le plus grand angle canonique entre ces sous-espaces, on a la majoration :

$$\sin 2\theta \leq \frac{\beta_1 - \beta_n}{\varepsilon}, \quad \text{avec } \theta \leq \frac{\pi}{4} \quad \text{si } \beta_1 - \beta_n < \varepsilon$$

avec :  $\varepsilon = \inf \{ \alpha_{s-1} - \alpha_s, \alpha_{s-r} - \alpha_{s+r+1} \}$  si  $s \neq 1$

$$\varepsilon = \inf \{ \alpha_{r+1} - \alpha_{r+2} \} \quad \text{si } s = 1.$$

Ce sont donc les écarts entre les valeurs propres qui "bordent" le sous-espace qui définissent la stabilité de ce sous-espace. Dans le cas du sous-espace engendré par les premiers facteurs (cas  $s = 1$ ), c'est l'écart entre la dernière valeur propre correspondant à ce sous-espace et la valeur propre

immédiatement consécutive qui compte. L'angle entre le sous-espace du tableau initial **A** et le sous-espace homologue du tableau perturbé **C** sera d'autant plus petit que cet écart entre valeurs propres est grand.

## **b – Epreuves empiriques de stabilité**

En pratique, les méthodes empiriques de validation qui permettent un premier contrôle de la qualité des résultats et de leur stabilité font partie intégrante du processus d'analyse des données.

Quels sont les différents éléments qui peuvent conditionner la qualité et la stabilité des résultats d'une analyse factorielle?

Nous en distinguons quatre :

- 1 - le choix et le poids des variables,
- 2 - le codage des variables,
- 3 - les erreurs de mesure,
- 4 - les poids des individus, les fluctuations d'échantillonnage <sup>1</sup>.

Les quatre sources de perturbation donnent lieu à des modifications du tableau initial et permettent de vérifier la permanence de la configuration initiale. Elles sont implicitement pratiquées dans la mesure où l'usage des méthodes factorielles, dans un but exploratoire, nécessite non pas une analyse mais une série d'analyses : à chaque étape, le tableau de données sera modifié par le choix des variables ou d'individus (avec ajout ou retrait de certains éléments), les corrections d'éventuelles erreurs, le recodage des données, etc..

Cette démarche proche de la "structuration en escalade" décrite par Mallows et Tukey (1982) permet une connaissance progressive du phénomène et constitue en soi une procédure de validation des résultats. Un exemple simple d'instabilité rappelé par Holmès (1985), est celui du point aberrant (marginal ou résultant d'une erreur) qui attire de façon excessive le plan principal et dont le retrait de l'analyse change l'orientation du plan.

### *- Le choix et le poids des variables*

Le problème se pose lorsque le statisticien a la possibilité d'échantillonner dans l'espace des variables, ce qui n'est pas toujours le cas. Les critères d'homogénéité et d'exhaustivité ne fournissent qu'un cadre général.

On pourra effectuer des "ponctions aléatoires" dans l'ensemble des variables, afin d'éprouver la sensibilité des résultats vis-à-vis de la composition de cet ensemble.

---

<sup>1</sup> Les trois premiers points correspondent à ce que Greenacre (1984) désigne par *stabilité interne* (l'univers est constitué par le tableau analysé, sans référence à une population plus large). Le quatrième point répond plutôt aux demandes de *stabilité externe* (visant à étendre les faits structuraux observés à partir du tableau analysé à une population plus générale).

Le problème du poids des variables se pose surtout en analyse en composantes principales (ou en analyse des correspondances s'il s'agit de tableaux de notes ou de mesures, et non de comptages).

Pour mettre en évidence une éventuelle invariance par rapport au système de pondération des variables, on procédera par exemple à la transformation suivante : l'analyse initiale étant faite sur les variables réduites (écarts-types unités), on dilatera ces écarts-types entre 1 et 2 (par un tirage pseudo-aléatoire) et on effectuera de nouveau l'analyse non-normée sur la matrice des covariances obtenues.

#### - *Le codage des variables*

Le codage apparaîtra comme source de perturbation éventuelle des résultats dans le cas des notes, des échelles ou des classements (par exemple en analyse des rangs ou des préférences). Il est alors important de vérifier que les configurations obtenues résistent à des changements de variables monotones très déformants (logarithme, exponentiel, etc.), afin de s'assurer que l'ordre des notes est plus important que les *propriétés métriques* particulières à l'échelle utilisée.

Il est intéressant alors de mettre en évidence un *codage minimal*, c'est-à-dire le codage le plus fruste susceptible de conserver les configurations observées.

Citons deux exemples : une analyse factorielle réalisée sur un tableau de dépenses individuelles de consommation donne lieu à une certaine typologie des postes de consommation. Cette analyse, refaite en codant simplement par "1" les dépenses strictement positives quels que soient leurs montants, produit alors une typologie des postes très voisine de la précédente (Jousselin, 1972). On conçoit que l'interprétation de la première analyse soit modifiée par ce résultat, qui souligne l'importance de l'accès à certains types de consommation, indépendamment de l'intensité de ces consommations.

On a obtenu un résultat analogue à propos de la typologie des activités réalisée à partir de budget-temps ; cette typologie n'a pas été bouleversée lorsque les durées positives des activités ont été remplacées par des "1", les durées nulles étant toujours codées par des "0". La simple mention d'une activité (lecture, promenade, soins aux enfants, etc.) jouait donc un rôle prédominant.

Nous avons brièvement envisagé le codage des variables dans une optique de validation des résultats. Mais le codage est une étape fondamentale dans le processus d'analyse des données.

Comme l'analyse des données elle-même, cette transformation a pour raison d'être l'augmentation de la *valeur pratique de l'information* (cf. § 4.1.5 b). Il s'agit de rendre celle-ci utilisable par l'algorithme et interprétable par l'utilisateur (cf. § 1.4.3 i pour quelques références concernant les diverses procédures de codage dans le cadre de l'analyse des correspondances multiples).

### - Les erreurs de mesure

L'ordre de grandeur de ces erreurs, ainsi que leur distribution approximative dans la population, doivent être spécifiés par l'utilisateur en fonction de sa propre connaissance du domaine étudié.

Par exemple dans le cas classique des réponses ordonnées du type : "pas du tout d'accord; pas d'accord; assez d'accord; tout à fait d'accord", on peut supposer que l'individu enquêté a une chance sur deux d'avoir exprimé exactement ce qu'il ressentait, une chance sur quatre (sauf aux extrémités) de répondre à une modalité immédiatement contiguë.

Les programmes de calcul permettront en général de simuler une grande variété de situations dont la traduction analytique serait inextricable. De ce fait les hypothèses que l'on soumet à l'épreuve d'un test peuvent être beaucoup mieux adaptées aux situations réelles et aux préoccupations des utilisateurs que les hypothèses classiques donnant lieu à une formulation analytique. En revanche la mise en œuvre de ces validations exige un certain travail de programmation (qui peut être facilité par l'utilisation d'un langage de simulation approprié).

### - Les poids des individus, les fluctuations d'échantillonnage

Deux types de calculs de stabilité peuvent être exécutés comme dans le cas du choix et du poids des variables cité ci-dessus : modifications des pondérations des individus; ponctions ou fractionnements de l'échantillon. Ces deux opérations doivent permettre d'apprécier la stabilité des résultats et l'on parle alors du "*poids bootstrap*" en référence à la méthode bootstrap présentée au paragraphe suivant.

Toutefois, les typologies obtenues par analyse factorielle n'exigent pas une représentativité de l'échantillon aussi stricte<sup>1</sup> que les estimations de pourcentages ou de moments d'ordre 1 (moyennes, fréquences). Cette relative stabilité vis-à-vis de la représentativité de l'échantillon est un fait d'expérience, étayé par les considérations sur la stabilité du paragraphe 4.2.1.a ci-dessus.

Dans les enquêtes par sondage, lorsque l'échantillon n'est pas représentatif et privilégie par exemple une sous-population de la population mère, chaque individu de l'échantillon est alors affecté d'un "coefficient de redressement" qui permet d'ajuster les moyennes ou les marges sur des valeurs connues dans la population parente<sup>2</sup>. Il n'est pas rare que les

<sup>1</sup> Bien entendu, un échantillon où certains aspects de la population parente sont absents, ne pourra pas fournir de résultats "extrapolables", même si les configurations obtenues sont stables.

<sup>2</sup> Ces redressements de tableaux à partir de leur marge se font en général à partir d'algorithmes itératifs (iterative proportional fitting) proposés à l'origine par Deming et Stephan (1940). Pour une vision historique et générale, cf. Thionet (1976) et d'autres articles de ce numéro spécial d'Annales de l'INSEE consacré aux ajustements de tableaux. Pour des développements récents sur ce thème, voir Deville et Särndal (1992).

typologies obtenues fassent preuve d'une bonne stabilité et qu'elles soient les mêmes que l'échantillon soit "redressé", ou que l'analyse soit faite sur les données brutes.

Mais les méthodes privilégiées pour étudier la stabilité des résultats vis-à-vis de fluctuations d'échantillonnage sont les techniques de rééchantillonnage examinées au paragraphe 4.2.2.

#### 4.2.2 Méthodes de rééchantillonnage (Jackknife, Bootstrap, Validation croisée)

Ce sont des méthodes de calculs intensifs qui reposent sur des techniques de simulations d'échantillons à partir d'un seul échantillon. Rendues possibles par la puissance de calcul des ordinateurs, ces techniques se substituent dans certains cas aux procédures plus classiques reposant sur des hypothèses contraignantes. Elles sont les seules procédures possibles lorsque la complexité analytique du problème ne permet pas d'inférence classique.

Elles consisteront pour nous à répéter des analyses sur les différents échantillons simulés puis à étudier les fluctuations des résultats obtenus, (valeurs propres, facteurs ou tout autre paramètre statistique à estimer). Pour cela, on évalue la variabilité réelle d'un paramètre par le truchement de sa variabilité pour l'ensemble de ces séries de données.

Plusieurs méthodes de validation permettent d'engendrer, de façon différente, les échantillons artificiels. Nous citerons les trois plus connues : Le Jackknife, le Bootstrap, la validation croisée.

##### a – Technique de Jackknife

Cette technique, proposée par Quenouille (1949) et reprise par Tukey<sup>1</sup> (1958) et décrite de façon complète par Miller (1974), consiste à supprimer de l'échantillon de base une seule observation ce qui rejoint le principe d'échantillonnage de la validation croisée. Si  $n$  est la taille de l'échantillon, on construit ainsi  $n$  échantillons de taille  $n - 1$ .

Le Jackknife s'utilise surtout pour calculer l'estimation du biais et de la variance d'estimateurs classiques, alors que la validation croisée sert principalement à calculer des erreurs de prédiction.

Un paramètre statistique à estimer  $\theta$  vaut  $\hat{\theta} = s(x)$  sur l'échantillon initial ou de base  $x = (x_1, x_2, \dots, x_n)$ . Il est calculé pour chacun des nouveaux échantillons obtenus et l'on peut ainsi évaluer sa variabilité.

---

<sup>1</sup> C'est Tukey qui lui a donné le nom de Jackknife (couteau à tout faire) que Tomassone *et al.* (voir diverses références) traduisent par Eustache (notion équivalente en Français), de même que ces auteurs traduisent Bootstrap (utilisé dans une expression désignant le fait de se hisser en tirant sur ses propre lacets de soulier) par Cyrano (allusion au voyage dans la lune).



Notons  $\hat{\theta}_{(i)} = s(\mathbf{x}_{(i)})$  l'estimateur obtenu sur l'échantillon  $\mathbf{x}_{(i)}$  qui n'est autre que l'échantillon de base  $\mathbf{x}$  privé de l'observation  $x_i$ .

L'estimateur jackknife  $\tilde{\theta}$  est donné par :

$$\tilde{\theta} = n\hat{\theta} - (n-1)\hat{\theta}_{(\cdot)}$$

où  $\hat{\theta}_{(\cdot)}$  est la moyenne des  $\hat{\theta}_{(i)}$  c'est-à-dire :

$$\hat{\theta}_{(\cdot)} = \frac{1}{n} \sum_i \hat{\theta}_{(i)}$$

Le biais jackknife vaut :

$$b = \hat{\theta} - \tilde{\theta} = (n-1)(\hat{\theta}_{(\cdot)} - \hat{\theta})$$

L'estimation jackknife de la variance vaut :

$$\text{Var}(\tilde{\theta}) = \frac{(n-1)}{n} \sum_i (\hat{\theta}_{(i)} - \hat{\theta}_{(\cdot)})^2$$

Le coefficient  $(n-1)$  peut surprendre dans le calcul du biais et dans celui de la variance. En fait, les  $n$  nouveaux échantillons  $\mathbf{x}_{(i)}$  ont deux à deux  $n-2$  valeurs en commun : leurs moyennes sont donc anormalement proches ce qui justifie, au moins intuitivement dans cette brève présentation, le coefficient  $(n-1)$ .

Cette technique est performante pour réduire le biais d'un estimateur et est, de ce fait, bien adaptée aux calculs de statistiques biaisées.

En particulier, dans le cas où l'estimation  $\hat{\theta} = s(\mathbf{x})$  est biaisée en  $1/n$ , c'est à dire si :

$$E(\hat{\theta}) = \theta + \frac{a}{n} + o\left(\frac{1}{n}\right)$$

on a alors :

$$E(\hat{\theta}_{(i)}) = \theta + \frac{a}{n-1} + o\left(\frac{1}{n}\right)$$

et donc :

$$E(\tilde{\theta}) = E(n\hat{\theta} - (n-1)\hat{\theta}_{(\cdot)}) = \theta + o\left(\frac{1}{n^2}\right)$$

La partie principale du biais est donc supprimée. Dans ce cas (favorable), l'estimation Jackknife est moins biaisée que l'estimation usuelle.

Bien qu'il exige moins de calculs que le bootstrap (voir ci-dessous), le Jackknife semble, cependant, moins souple et parfois moins fiable.

Il est surtout moins bien adapté pour calculer des intervalles de confiance et estimer les paramètres d'une analyse multidimensionnelle. Il est en échec s'il s'agit d'estimer une quantité qui n'est pas une fonction continue des données (comme la médiane) ou une quantité dont l'espérance théorique dépend de  $n$ , comme les valeurs propres en analyse en composantes principales ou en analyse des correspondances.

## b – Technique de Bootstrap

Cette technique, introduite par Efron (1979), consiste à simuler  $m$  (généralement supérieur à 30) échantillons de même taille  $n$  que l'échantillon initial. Ils sont obtenus par tirage au hasard avec remise parmi les  $n$  individus observés au départ, ceux-ci ayant tous la même probabilité  $1/n$  d'être choisis. Certains individus auront de ce fait un poids élevé (2, 3,...) alors que d'autres seront absents (poids nul).

Cette méthode est employée pour analyser la variabilité de paramètres statistiques simples en produisant des intervalles de confiance de ces paramètres. Elle peut aussi être appliquée à de nombreux problèmes pour lesquels on ne peut pas estimer analytiquement la variabilité d'un paramètre. Ceci est le cas pour les caractéristiques des méthodes multidimensionnelles où les hypothèses de multinormalité sont rarement vérifiées.

Le bootstrap n'est rien d'autre qu'une technique de simulation particulière, fondée sur la distribution empirique de l'échantillon de base. Efron et Tibshirani (1993) réservent le nom de *non-parametric bootstrap* à ce type de simulation, et qualifie de *parametric bootstrap* les simulations qui mettent en jeu une distribution théorique et des paramètres calculés à partir de l'échantillon (simulations classiques)<sup>1</sup>. On constate que le Jackknife est déterministe et fait intervenir de façon symétrique l'échantillon (sans nécessiter de procédure de tirage pseudo-aléatoire), contrairement au bootstrap<sup>2</sup>.

Prenons l'exemple de l'estimation du coefficient de corrélation  $r$  entre deux variables. Le principe consiste à calculer le coefficient de corrélation pour chaque échantillon simulé (pour lequel on effectue un tirage avec remise des *couples* d'observations). On établit alors la distribution des fréquences du coefficient de corrélation (on porte en ordonnée le nombre d'échantillons ayant une même valeur de  $r$ , laquelle est représentée en abscisse). Puis on calcule la probabilité pour que le coefficient de corrélation d'un échantillon soit compris dans différentes fourchettes de valeurs définissant ainsi les intervalles de confiance. On obtient une estimation de la précision de la valeur de  $r$  obtenue sur l'échantillon de base sans faire l'hypothèse d'une distribution normale des données. Les bornes de l'intervalle de confiance peuvent être estimées directement par les quantiles de la distribution simulée.

Pour estimer les valeurs propres, les taux d'inertie et les coordonnées factorielles issus d'une analyse en composantes principales, par exemple, le

<sup>1</sup> Sur les techniques de génération de variables pseudo-aléatoires, cf. par exemple Newman et Odell (1971), Ripley (1983).

<sup>2</sup> Il existe des variantes du Jackknife (comme le *delete-d Jackknife*) qui préconisent des ponctions plus importantes dans l'échantillon et nécessitent soit un découpage arbitraire, soit des tirages pseudo aléatoires, car le nombre d'échantillons possibles devient alors dirimant (cf. Efron et Tibshirani, *op.cit.*, pour une revue des différentes méthodes de rééchantillonnage, et une comparaison bootstrap-Jackknife).

principe est le même que pour le coefficient de corrélation ; on effectue sur chaque échantillon simulé, une analyse en composantes principales puis on établit une distribution de fréquences pour chacune des composantes.

La méthode de bootstrap donne dans la plupart des cas une bonne image de la précision statistique de l'estimation sur un échantillon.

Les recherches théoriques menées par Efron en particulier montrent que, pour de nombreux paramètres statistiques, l'intervalle de confiance correspondant à la distribution simulée par bootstrap et celui correspondant à la distribution réelle sont généralement de même amplitude.

Un exemple classique d'échec de l'estimation par bootstrap est l'estimation des bornes d'un intervalle pour une loi uniforme dans cet intervalle. Il est clair en effet que dans ce cas, l'estimation donnée par les valeurs extrêmes ne sera pas améliorée par des tirages à l'intérieur de l'échantillon de base<sup>1</sup>.

### c – Validation croisée

La validation croisée, dont le principe et les principales références ont été donnés à la section 3.3 (paragraphe 3.3.5 dévolue aux règles d'affectation en analyse discriminante) est plutôt utilisée lorsque l'échantillon de base est de taille petite ou moyenne et permet d'évaluer les procédures d'ajustement et de reclassement comme la régression et l'analyse discriminante.

Rappelons qu'il s'agit à l'origine de séparer l'échantillon de base en deux blocs de tailles pouvant être inégales. Un des blocs constitue l'échantillon d'apprentissage sur lequel est formulé le modèle et sont élaborées les règles de décision ou d'affectation ; l'autre compose l'échantillon-test sur lequel sont appliquées les règles et estimées les performances du modèle.

Afin d'utiliser un échantillon d'apprentissage qui soit le plus grand possible, la validation croisée a adopté un principe proche de celui du Jackknife pour constituer les blocs. Il consiste à construire l'échantillon-d'apprentissage en retirant un seul individu (ou un groupe) de l'échantillon de base. La prédiction est faite sur un seul individu à la fois, et par rotation sur tous les individus, ce qui produit finalement un échantillon d'épreuves aussi important que l'échantillon initial. Par exemple dans le cadre de l'analyse discriminante, on estime par la validation croisée le taux d'erreur de classement<sup>2</sup>.

Pour des raisons d'économie de calcul, dans le cas de grands échantillons, il est possible de retirer de l'échantillon d'apprentissage  $k$  individus ( $k > 1$ )

<sup>1</sup> Pour une revue critique de l'utilisation du bootstrap (avec discussions), cf. Young (1994).

<sup>2</sup> Ce taux d'erreur est distinct, mais en général peu différent, de ce que serait une estimation jackknife du taux d'erreur dans l'échantillon d'apprentissage. Dans le cadre de ces problèmes de classement, la validation croisée, plus simple dans son principe et plus intuitive, est systématiquement utilisée. Pour une mise au point et des confrontations de méthodes de calcul de taux d'erreur, voir Hand (1986, 1987).

plutôt qu'un seul, ce qui divise approximativement par  $k$  le volume total de calcul.

Complément indispensable des méthodes de prédiction et de classement, les méthodes de validation croisée sont moins utilisées que le bootstrap pour valider les représentations des méthodes plus descriptives.

### 4.2.3 Zones de confiance, zones de garde ; nombre d'axes

Les résultats fournis par les méthodes factorielles ne sont pas des assertions, mais des représentations, c'est-à-dire des objets complexes, auxquels s'appliquent mal les différentes techniques de mesure d'information usuelles en statistique.

Comment valider une forme observée dans un plan factoriel ?

- Par des procédures externes, analogues à celles mentionnées pour le choix du nombre d'axes : connaissances a priori, positionnement de variables supplémentaires.
- Par des calculs de stabilité adaptés (exploration d'un voisinage des données construit à partir des erreurs de mesure ou de réponses).
- Par des calculs de zones de confiance pour les positions des points-lignes et des points-colonnes. Ces calculs peuvent être analytiques, fondé sur des hypothèse probabilistes, ou au contraire, fondés sur les techniques de rééchantillonnage exposées au paragraphe précédent.

On commencera par présenter le cadre de l'utilisation des simulations pour le calcul de ces zones de confiance. Dans ce cadre, le bootstrap, qui constitue une méthode de simulation non paramétrique d'une grande souplesse, jouera un rôle de premier plan.

#### a – Zones de confiance établies par bootstrap

Le bootstrap est un outil privilégié pour étudier la stabilité des formes. Une application à l'exemple d'analyse des correspondances (section 1.3) nous montrera la simplicité et l'efficacité de la méthode.

##### - Présentation à partir d'un exemple

Reprenons le sous-tableau du tableau 1.3 - 10 correspondant aux seules lignes actives.

Une simulation bootstrap classique consiste à tirer avec remise les  $k = 12\ 888$  contacts-média (chacun d'entre eux correspondant à une case  $(i,j)$  du tableau 4.2 - 1). Cela revient à faire autant de tirages selon une loi multinomiale dont les probabilités de tirage sont :  $p_{ij} = k_{ij} / k$ . On peut vérifier (empiriquement) qu'il est équivalent, au niveau des résultats de la simulation, d'utiliser l'approximation normale de la loi multinomiale,

c'est-à-dire de générer  $k_{ij}^*$ , variable normale de moyenne  $k_{ij}$  et de variance  $k_{ij} \frac{1-k_{ij}}{k}$  (la valeur ainsi générée sera arrondie à l'entier supérieur)<sup>1</sup>.

**Tableau 4.2 - 1**  
Tables de contingence croisant 6 types de contacts-média (colonnes)  
avec 8 professions (lignes) [partie active du tableau 1.3 - 10].

Professions	Radio	Tél.	Quot.N.	Quot	R. P.Mag.	P.TV
1 - Agriculteur	96	118	2	71	50	17
2 - Petit patron	122	136	11	76	49	41
3 - Prof.Cad.S.	193	184	74	63	103	79
4 - Prof.interm	360	365	63	145	141	184
5 - Employé	511	593	57	217	172	306
6 - Ouvrier qual	385	457	42	174	104	220
7 - Ouvrier n-q	156	185	8	69	42	85
8 - Inactif	1474	1931	181	852	642	782

Le tableau 4.2 - 2 donne un exemple de deux tableaux générés de cette façon. Les numéros en début de ligne seront ceux qui figureront sur le plan factoriel de la figure 4.2 - 1. On va, pour cet exemple, générer 30 réplifications<sup>2</sup>, chiffre largement suffisant, on va le voir, pour donner une bonne idée de la stabilité des résultats.

**Tableau 4.2 - 2**  
Exemple de deux réplifications des valeurs du tableau 4.2 - 1

1	109	120	1	78	48	20
2	126	142	8	76	53	30
3	196	181	80	77	109	72
4	384	365	60	133	138	203
5	514	596	59	228	172	316
6	378	467	33	171	100	223
7	169	188	8	79	38	81
8	1519	1961	158	893	632	764
1	83	138	3	79	62	19
2	142	142	8	82	50	26
3	198	163	63	68	114	85
4	359	367	73	155	132	196
5	503	561	56	266	173	294
6	395	432	25	171	104	220
7	149	179	16	74	50	83
8	1488	1919	182	852	611	794

<sup>1</sup> Les trois cases d'effectif faible (<12) pour lesquelles une telle approximation est discutable ont en fait une influence quasi nulle sur les résultats. Cela ne serait pas le cas si une colonne (ou une ligne) entière était formée d'effectifs faibles. Pour des programmes de génération de variables pseudo-aléatoires normales, cf. Neave (1973), Brent (1974).

<sup>2</sup> On utilisera le terme d'origine anglaise *réplification* pour désigner une simulation d'échantillon.

Notons qu'une analyse des correspondances faite sur un seul tableau répliqué suffit à donner une forte présomption de stabilité. L'observation, au sens près des axes, du même *pattern* (de la même forme) signifie que la structure observée a résisté à la perturbation constituée par la simulation. Il est en effet extrêmement improbable de retrouver par hasard un agencement complexe de points.

C'est là une différence fondamentale avec la statistique uni-dimensionnelle, pour laquelle une répllication isolée n'est pas utilisable. Cependant, dans la plupart des cas, la structure est partiellement déformée et l'on souhaite pouvoir isoler ses éventuelles parties stables. C'est alors que la répétition des réplifications est utile, pour limiter la subjectivité dans les appréciations.

Il existe plusieurs façon de mettre à profit ces 30 réplifications pour construire des intervalles de confiance.

Procéder à 30 analyses indépendantes est exclu, car les axes correspondant à des valeurs propres voisines peuvent changer de rang ou subir des rotations. De plus, ces axes sont définis au signe près, et donc les tentatives de superposition des structures peuvent être laborieuses<sup>1</sup>.

Il reste comme possibilité :

- a - analyser les juxtapositions de tableaux de contingence, en lignes (comme esquissé dans le tableau 4.2 - 2), de façon à étudier la variabilité des lignes, et en colonnes (pour positionner les différentes colonnes simulées).
- b - projeter en éléments supplémentaires les lignes (et les colonnes) simulées) dans les plans factoriels issus de l'analyse de la table de contingence initiale (qui est pour ce modèle, rappelons-le, l'espérance des matrices simulées).
- c - calculer un tableau de contingence moyen, et projeter les lignes ou les colonnes comme en *b*.

Les trois procédures *a*, *b* et *c* donnent en fait des résultats extrêmement voisins dans le cas de l'exemple traité, et plus généralement dans tous les cas où il existe effectivement une structure stable.

En effet, dans ce cas, il y a une forte redondance dans les réplifications, et la propriété d'équivalence distributionnelle nous assure que les distances calculées entre les colonnes sur les tableaux juxtaposés en ligne (par exemple) sont voisines des distances calculées entre les colonnes du tableau moyen (obtenu en agrégeant les lignes homologues, par exemple les deux lignes 1, les deux lignes 2, etc. du tableau 4.2 - 2). Or ce tableau moyen converge, lorsque le nombre de réplifications augmente, vers le tableau initial qui est l'espérance mathématique des différents tableaux pseudo-aléatoires.

---

<sup>1</sup> Les méthodes d'analyse procrustéennes (cf. § 3.8.2) ont précisément pour objectif de détecter des structures superposables après déplacement et dilatation.

La figure 4.2 - 1 représente le premier plan factoriel, ou plan  $(F_1, F_2)$ , de l'analyse des correspondances de la table à 6 colonnes et 248 lignes obtenues en juxtaposant à la table originale les 30 tables simulées suivant le modèle précédent ( $248 = 8 + 30 \times 8$ ).

On voit que les enveloppes convexes des points correspondant à des lignes homologues (catégories socio-professionnelles) sont bien séparées, à l'exception des catégories 5 et 6 (employés et ouvriers qualifiés).

Bien entendu, la même procédure peut être appliquée aux points-colonnes (contact-média).

La forme, mais aussi la taille des enveloppes convexes apportent une information supplémentaire par rapport à la figure originale 1.3 - 23 du chapitre 1. Ainsi, on peut affirmer que les ouvriers non-qualifiés (symbole 7) ont un comportement en continuité, mais cependant résolument distinct de celui des ouvriers qualifiés (6). Même observation en ce qui concerne l'absence de solution de continuité entre agriculteurs (1) et petits patrons (2).

La figure 4.2 - 2 représente le second plan factoriel, ou plan  $(F_2, F_3)$ , de la même table. On retrouve sur l'axe horizontal les distinctions observées sur l'axe vertical précédent, mais la confusion est totale sur l'axe  $q$  (vertical). Seule la classe 8 (inactifs) occupe une position typée sur l'axe 3, en s'opposant aux autres classes.

En conclusion, en même temps qu'un enrichissement de l'information sur le plan  $(F_1, F_2)$ , on a un critère pour choisir le nombre d'axes de représentation, limité ici à 2.

Notons que le processus s'applique également aux variables supplémentaires. Ainsi, pour prendre un exemple, les lignes du tableau de contingence 1.3 - 10 (sexe, âge, niveau d'éducation) de la section 1.3 peuvent être répliquées en utilisant un schéma multinomial similaire à celui utilisé pour les variables actives, les projections de ces lignes simulées dans les plans factoriels définissant alors des zones de confiance<sup>1</sup>.

## **b – Autres types de simulation bootstrap**

On vient de voir une application du bootstrap à la validation des représentations issues de l'analyse des tables de contingence simples.

### *- Cas de l'analyse des correspondances multiples*

Dans le cas de l'analyse des correspondances multiples, une réplique bootstrap est obtenue en tirant avec remise les individus, lignes du tableau de données  $R$  ou de façon équivalente, lignes du tableau disjonctif associé  $Z$ .

<sup>1</sup> En fait, nous utilisons ici l'expression zone de confiance en parlant simplement des enveloppes convexes des projections des valeurs répliquées. Les enveloppes convexes, étudiées par Efron (1965) peuvent être "pelées" progressivement de façon à obtenir des estimations non-paramétriques de zones de confiance (cf. Barnett, 1976; Green, 1981; et Holmes, 1985, qui publie également des exemples et les programmes de calcul correspondant).

Figure 4.2-1

Zones de confiance "bootstrap"  
pour les lignes actives de  
l'exemple de la section 1.2  
Plan (F<sub>1</sub>, F<sub>2</sub>)

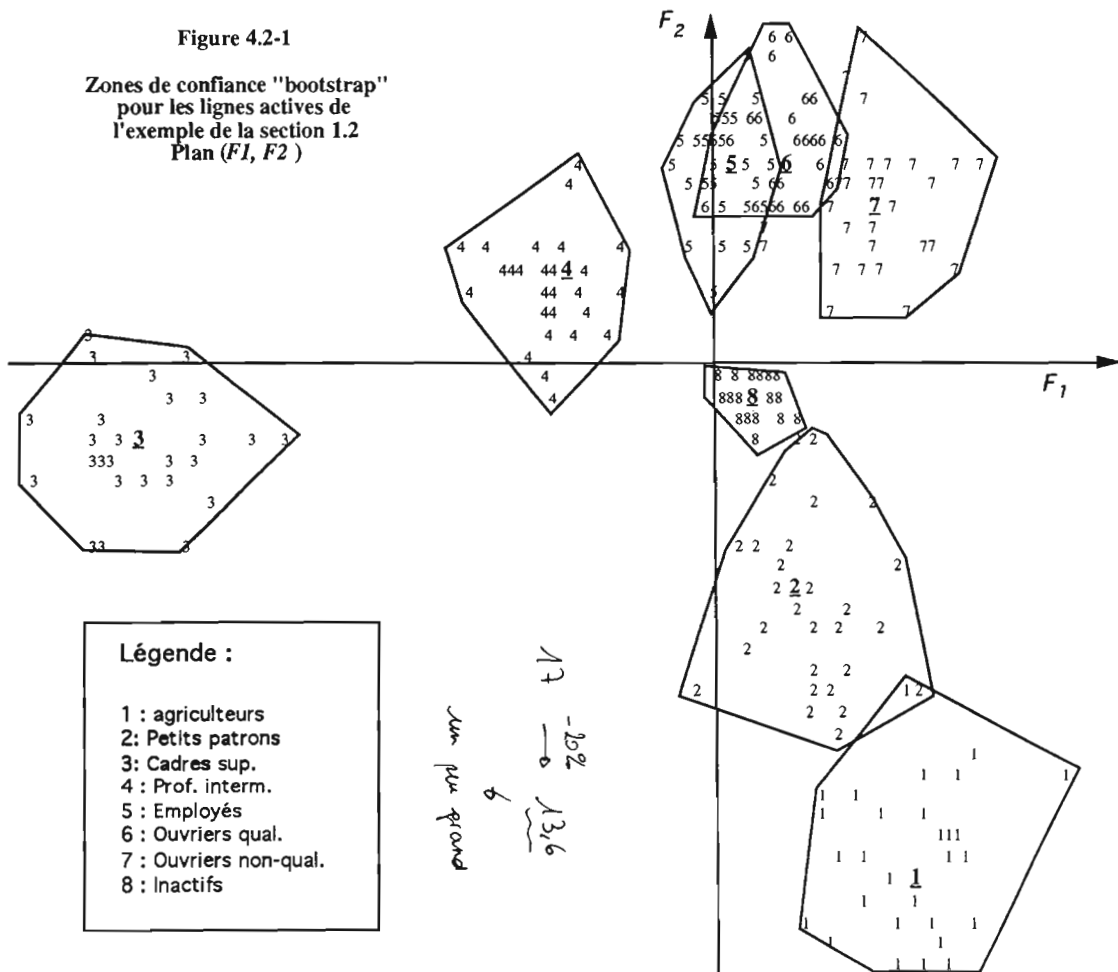
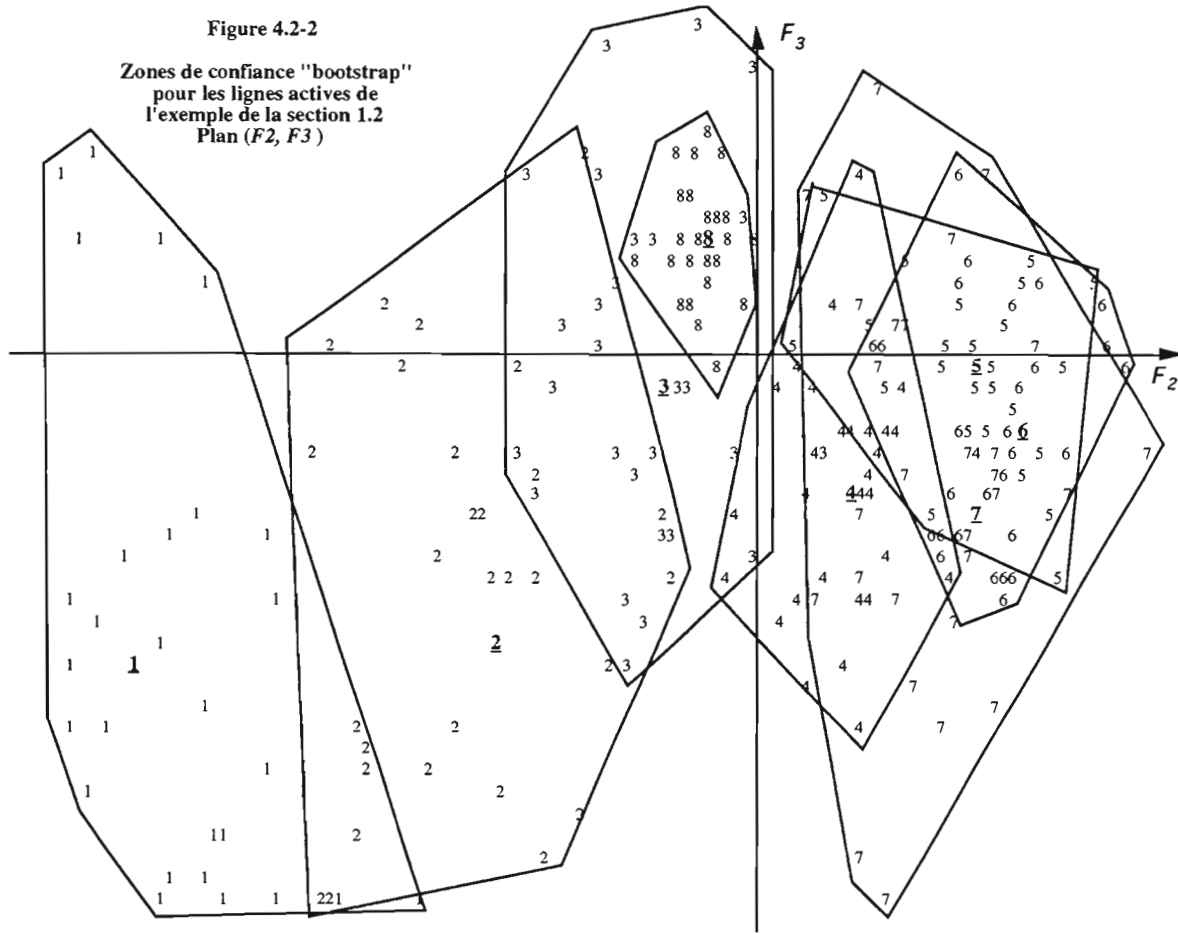




Figure 4.2-2

Zones de confiance "bootstrap"  
 pour les lignes actives de  
 l'exemple de la section 1.2  
 Plan ( $F_2$ ,  $F_3$ )



Chaque réplication permet de construire un tableau de Burt, dont les lignes sont projetées en éléments supplémentaires dans les plans factoriels issus de l'analyse du tableau de Burt initial. Les zones de confiance obtenues sont d'autant plus utiles ici pour choisir la dimension de l'espace de représentation que les valeurs propres et les taux d'inertie sont, on l'a vu, d'une interprétation difficile<sup>1</sup>.

#### - Cas de l'analyse en composantes principales

C'est le domaine d'application qui a donné lieu au plus grand nombre de travaux, utilisant notamment des méthodes de rééchantillonnage antérieures au bootstrap, comme la validation croisée et ses variantes. C'est ainsi que Wold (1978), puis Eastment et Krzanowski (1982), Krzanowski (1987) proposent des méthodes de validation croisée pour déterminer la dimension de l'espace de représentation. Ces auteurs utilisent la théorie de la perturbation pour alléger les calculs (l'omission d'un ou plusieurs éléments, qui est à la base de la validation croisée, est considérée comme une perturbation du tableau de données (cf. § 4.2.1), et une formule approchée permet d'éviter de refaire une analyse complète). Besse et Ferré (1993) ont montré que la réitération de ces approximations revenait en fait à utiliser le critère classique de la part de variance expliquée par les axes.

Si l'on excepte les travaux de Gifi (1981) qui concernent plus spécifiquement l'analyse des correspondances (le principe du bootstrap est en fait sensiblement le même pour toutes les méthodes factorielles) les premiers travaux d'application du bootstrap à la validité des résultats en analyse en composantes principale sont ceux de Diaconis et Efron (1983), Stauffer *et al.* (1985), Holmes (1985), Beran et Srivastava (1985), Daudin *et al.* (1988), Holmes (1989), qui construisent des intervalles de confiance pour les valeurs propres et les composantes, ou étudient les propriétés asymptotiques des intervalles ou des estimations obtenus.

L'algorithme qui nous paraît le mieux adapté pour les intervalles de confiance est analogue à celui préconisé pour l'analyse des correspondances multiples : une réplication consiste en un tirage avec remise des  $n$  individus (vecteurs-observations), suivi du positionnement des  $p$  nouvelles variables ainsi obtenues en variables supplémentaires sur les  $q$  premiers axes factoriels de l'analyse de base. Après  $r$  répliques, on obtient, pour chacune des  $p$  variables, un nuage de  $r$  points dont l'enveloppe convexe (éventuellement "pelée") constituera la zone de confiance empirique cherchée de la variable<sup>2</sup>. On peut de la même façon construire

<sup>1</sup> Les premières application du bootstrap pour évaluer la stabilité et pour construire des zones de confiance à partir d'une analyse des correspondances multiples (*homogeneity analysis* selon la terminologie de ces auteurs) sont celles de Gifi (1981), Meulman (1982), puis Greenacre (1984).

<sup>2</sup> Il existe plusieurs variantes possibles quant au choix de l'espace factoriel commun. Holmes (1985) applique une méthode d'analyse conjointe de tableaux (méthode STATIS, cf. L'Hermier des Plantes, 1976; Lavit, 1988) au tableau initial et à l'ensemble de ses répliques.

des régions de confiance pour les éventuelles modalités de variables nominales supplémentaires<sup>1</sup>.

### c – Zones de garde en analyse des correspondances

La notion de zone de garde n'est pas une application des méthodes de rééchantillonnage, mais est cependant mentionnée ici comme une procédure rapide d'appréciation de la position des points. Cette notion s'applique surtout aux éléments supplémentaires ou à ceux qui ont une faible contribution à l'inertie sur les axes considérés, et seulement dans le cas de l'analyse des correspondances des tables de contingences .

Un point a une position "significative" si son éloignement au centre de gravité n'est pas dû au hasard. Rappelons que la distance du  $\chi^2$  d'un point-profil  $i$  au centre de gravité s'écrit :

$$d^2(G,i) = \sum_j \frac{1}{f_{i,j}} (f_{ij} - f_{i,j})^2 = \frac{1}{f_{i,j}} \sum_j \frac{(f_{ij} - f_{i,j})^2}{f_{i,j}}$$

Considérons l'hypothèse nulle  $H_0$  selon laquelle un point  $i$  ne diffère du centre de gravité  $G$  du nuage que par des fluctuations aléatoires.

Alors  $k f_{i,j} d^2(G,i) = k \sum_j \frac{(f_{ij} - f_{i,j})^2}{f_{i,j}}$  suit approximativement un  $\chi^2_{(p-1)}$ .

La projection orthogonale de ce  $\chi^2_{(p-1)}$  sur un sous-espace à  $q$  dimensions fixé à l'avance ( $q < p - 1$ ) est un  $\chi^2$  à  $q$  degrés de liberté noté  $\chi^2_q$ . Il est important de spécifier que le sous-espace est fixé à l'avance car le raisonnement de s'applique pas à un sous-espace qui serait calculé après la réalisation du  $\chi^2$  dans l'espace à  $p-1$  dimensions et donc ne s'applique pas aux variables actives ayant une influence sur les  $q$  premiers axes.

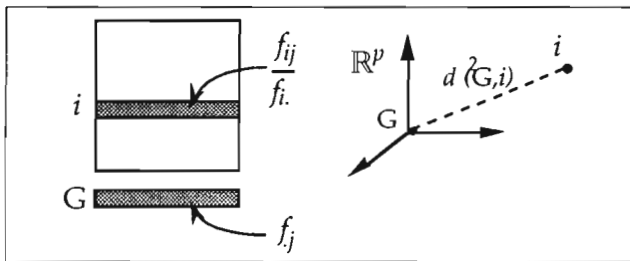


Figure 4.2 - 2  
Représentation de la distance d'un point au centre de gravité

<sup>1</sup> Pour chaque réplification, on positionne les centres de gravité des individus appartenant à chacune des modalités de la variable nominale.

Dans un espace factoriel à deux dimensions, la distance du point  $i$  projeté au centre de gravité s'écrit :

$$r^2(G, i) = \psi_{i1}^2 + \psi_{i2}^2$$

Sous l'hypothèse nulle  $H_0$  et sous la condition que la contribution du point  $i$  est faible ou nulle (indépendance de  $i$  et du plan factoriel), la quantité  $k f_i \cdot r^2(G, i)$  est une réalisation d'un  $\chi^2$  à 2 degrés de liberté.

Au seuil de 5%, on a :  $\text{Prob}[\chi_2^2 \geq 5.99] = 0.05$  et l'on comparera  $k f_i \cdot r^2(G, i)$  à 5.99, soit encore  $r(G, i)$  à  $r_i = \sqrt{\frac{5.99}{k f_i}}$ . On peut donc calculer le rayon des cercles correspondant qui ne dépendent que de l'effectif  $k_i = k f_i$  de la catégorie concernée.

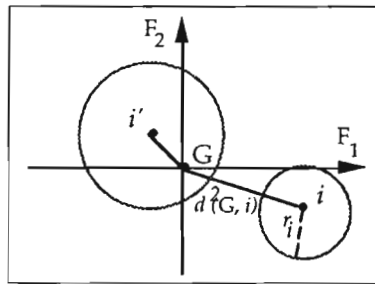


Figure 4.2 - 3  
Cercles de garde tracés autour des points

Pour juger la "position significative" d'un point  $i$  dans un plan factoriel, on calculera par conséquent le rayon  $r_i$  du cercle de garde relatif à un point  $i$ . Ces cercles de garde doivent être tracés autour de l'origine des axes qui est le centre de gravité du nuage. On ne rejettera pas l'hypothèse d'indépendance (donc de non pertinence du point  $i$ ) si ce cercle contient le point  $i$ .

S'il existe plusieurs points dont la position est à éprouver, il est plus simple de centrer les cercles correspondants sur les points eux-mêmes et de regarder s'ils contiennent ou non le centre de gravité  $G$ <sup>1</sup>.

Plus généralement, dans un sous-espace à  $q$  dimensions, on comparera la quantité :

$$k f_i \cdot r^2(G, i) = k f_i \cdot \sum_{\alpha=1}^q \psi_{i\alpha}^2$$

aux fractiles d'un  $\chi_q^2$  à  $q$  degrés de liberté.

<sup>1</sup> Il ne s'agit cependant pas de régions de confiance des points, mais d'un test simple de signification statistique de la position de chaque point par rapport à l'origine. On ne peut pas en particulier comparer entre elles les positions relatives de deux points.

### d - Autres régions de confiances

Bien que les techniques de *multidimensional scaling*<sup>1</sup> ne soient pas traitées dans cet ouvrage, il faut mentionner dans ce paragraphe, pour leur intérêt méthodologique, les travaux de Ramsay (1978) (zones de confiances fondées sur la distribution des distances entre individus pour la méthode dite *MULTISCALE*) et de Weinberg *et al.* (1984) (zone de confiances fondés sur le Jackknife et le bootstrap pour la méthode *INDSCAL*).

Plus proches de nos préoccupations, on mentionnera également les ellipses de confiance proposées par Saporta et Hatabian (1986), qui s'appliquent à toute catégorie de variable nominale supplémentaire (qu'il s'agisse d'analyse des correspondances ou d'analyse en composantes principales). Ces ellipses se calculent à partir de la matrices des covariances des coordonnées factorielles des points appartenant à la catégorie.

A côté des travaux précités sur les régions de confiance bootstrap, Gifi (1981, 1990) a également proposé des ellipsoïdes de confiance fondés sur la méthode *delta* (cf. par exemple, Rao, 1973 ; Efron, 1982). Cette méthode généralise au cas multidimensionnel le résultat élémentaire suivant :

Proposons-nous de calculer la variance et la loi asymptotique d'une (bonne) fonction  $g(X)$  d'une variable aléatoire  $X$  de moyenne  $\mu$  et de variance  $\sigma^2$ . A partir du développement de Taylor de  $g$  autour de  $\mu$  :

$$g(t) \approx g(\mu) + (t - \mu) g'(\mu)$$

on déduit immédiatement :

$$\text{var} [g(X)] \approx g'(\mu)^2 \sigma^2$$

Plus généralement, la méthode *delta* est fondée sur le résultat suivant :

Si l'on a une fonction  $\mathbf{y}_n = \Phi(\mathbf{x}_n)$  d'une suite  $\mathbf{x}_n$  de vecteurs aléatoires tels que  $\sqrt{n}(\mathbf{x}_n - \mu)$  est asymptotiquement normal de moyenne nulle et de matrice des covariances  $\Sigma$  ( $\Phi$  est supposée différentiable en  $\mu$ ), alors  $\sqrt{n}(\mathbf{y}_n - \Phi(\mu))$  est aussi asymptotiquement normal de moyenne nulle et de matrice des covariances  $\mathbf{V}(\sqrt{n} \mathbf{y}_n) = \partial\Phi(\mu) \Sigma \partial\Phi(\mu)$ , où  $\partial\Phi(\mu)$  est la matrice des dérivées partielles de  $\Phi$  au point  $\mu$ .

Si les composantes de  $\mathbf{y}_n$  sont les coordonnées d'un point sur deux axes factoriels et les composantes de  $\mathbf{x}_n$  sont les éléments du tableau de données ( $\mathbf{x}_n$  est par exemple une table de contingence d'effectif total  $n$ , dont la loi asymptotique est une loi normale - cf. § 4.1.2), la méthode *delta* permet

<sup>1</sup> Ensemble de techniques de représentations de systèmes de distances entre points, développé autour des *Bell Laboratories* et de la revue *Psychometrika*, avec, à l'origine, des contributions de R. Shepard, L. Guttman, J. Kruskal, D. Carrol (cf. par exemple Kruskal et Wish, 1978 ; Schiffman *et al.*, 1981). Cf. également l'article de synthèse sur les méthodes et les programmes de Drouet d'Aubigny (1993).

d'estimer la matrice des covariances de  $y_n$ , et donc de construire des zones de confiance ellipsoïdales autour du point correspondant.

Des formules analogues à la formule [4.2 - 4] (§ 4.2.1) permettent d'estimer  $\partial\Phi(\mu)$ . Dans le calcul de  $V(\sqrt{n} y_n)$ , les valeurs théoriques sont remplacées par leurs estimations empiriques<sup>1</sup>.

#### 4.2.4 Nombre de classes et validation des classifications

Dès les premières tentatives de classification s'est posé le problème du nombre de classes à retenir en vue d'une utilisation particulière. Déjà, sous cette formulation pragmatique, le problème est moins ambitieux que celui de savoir combien de classes existent réellement dans le corpus de données soumis à l'analyse.

La classification peut en effet être utilisée simplement pour explorer les données, généralisant au cas multidimensionnel l'histogramme qui permet de schématiser une distribution numérique unidimensionnelle. Il s'agit en fait de l'utilisation la plus courante dans le cas des traitements de fichiers d'enquêtes. On peut aussi espérer découvrir des classes existantes, dans les cas les plus favorables. Les questions sont aussi simples que les réponses sont complexes : Existe-t-il des classes ? Si oui, combien ?

On évoquera brièvement quelques travaux réalisés à propos de l'existence et de la détermination du nombre des classes. La méthodologie de la validation est analogue à celle déjà rencontrée à propos des méthodes factorielles.

##### - *Cadre inférentiel général*

Il sera possible de tester des hypothèses nulles (analogues de l'hypothèse d'indépendance pour les méthodes factorielles) qui sera selon les cas une hypothèse d'homogénéité ou d'uniformité spatiale de la distribution multidimensionnelle des observations à classer. Toutefois, comme dans le cas des méthodes factorielles, ce type de test, tout en fournissant des repères et un cadre conceptuel intéressant, sera de peu d'utilité pratique, car l'hypothèse d'absence de structure, trop sévère, sera la plupart du temps rejetée.

##### - *Validation empirique, calculs de stabilité*

Des procédures empiriques, en général variables selon les domaines d'application ou la nature du tableau des données, seront assez largement utilisées.

---

<sup>1</sup> Comme les zones de confiance bootstrap, avec lesquelles la compatibilité empirique semble bonne, les zones déterminées par la méthode delta peuvent concerner les variables actives, ce qui leur confère un avantage important sur les autres techniques évoquées dans ce paragraphe.

Enfin, des calculs de stabilité, utilisant des méthodes de simulation ou de rééchantillonnage, permettront d'éprouver la qualité de résultats et de porter une appréciation sur la réalité des classes produites par les algorithmes.

#### *- Importance des critères externes*

Le rôle des critères externes (connaissances a priori, identification ou caractérisation des classes à partir de variables supplémentaires) sera souvent primordial dans la pratique. Ainsi, une classe mal différenciée, mais identifiée par une catégorie de variable nominale supplémentaire deviendra, dans bien des cas, digne d'être retenue.

Il existe cependant une différence fondamentale avec les méthodes factorielles : il n'y a pas en classification l'équivalent du théorème d'Eckart et Young (décomposition aux valeurs singulières), et donc pas de paramètres aussi intrinsèques que les valeurs propres<sup>1</sup>. Il existe en revanche une riche flore d'algorithmes dont l'utilisation simultanée sur un même tableau constitue d'ailleurs une épreuve pragmatique de stabilité de structures observées.

Commençons par mentionner quelques travaux de synthèse sur le sujet. Une contribution récente de Bock (1994) sur les problèmes et l'avenir des méthodes de classification comprend une brève mais dense revue des problèmes de validation. D'autres revues intéressantes sont celles de Gordon (1987) (limitée à la classification hiérarchique), de Hartigan (1985), de Bock (1985), de Perruchet (1983), de Dubes et Jain (1979). Enfin on trouvera plus bas plusieurs références de contributions consacrées à des comparaisons de méthodes.

### **a – L'hypothèse d'absence de structure, les modèles**

Il existe de nombreux travaux sur ce thème et, à de rares exceptions près, ils ne concernent que les méthodes de classification utilisées de façon isolée. Dans cet ouvrage où nous considérons les méthodes factorielles et les méthodes de classification comme complémentaires (et devant être utilisées simultanément), on peut donc préconiser sans hésiter, au moins dans un premier temps, les tests d'indépendance ou de sphéricité déjà évoqués à propos des méthodes factorielles. Il est en effet extrêmement improbable que des variations de densité à l'intérieur d'un nuage de points ne se répercutent pas sur une ou plusieurs valeurs propres d'une analyse en axes principaux.

On peut objecter qu'un ellipsoïde peut être allongé, mais parfaitement homogène. Dans ce cas, une coupure en deux de son grand axe produit deux classes qui, même si elles ne sont pas séparées par une zone de faible densité, ne peuvent être considérées comme le fruit du hasard. Il s'agit en

---

<sup>1</sup> On a vu au chapitre 2, § 2.5.3, les relations qui pouvaient exister dans certains cas entre valeurs propres et indices de niveaux relatifs à une même table de contingence.

fait de la meilleure coupure en deux classes de l'échantillon. On voit qu'il faudrait préciser ce que l'on entend par classe. En fait, il y a presque autant de définitions des classes que de critères de classification utilisés pour les obtenir.

### - Modèles de mélanges

Le modèle théorique de base le plus répandu est le modèle des mélanges de distributions. L'observation  $x_i$  ( $i \leq n$ ) est alors une réalisation d'une variable aléatoire  $x$  de densité  $f(x)$  :

$$f(x) = \sum_{k=1}^q p_k f_k(x)$$

avec pour tout  $k$ ,

$$0 < p_k < 1 \quad \text{et} \quad \sum_{k=1}^q p_k = 1$$

Dans cette formule,  $f_k(x)$  est la densité de la classe  $k$  (dont la forme doit être spécifiée ; par exemple : densité d'une loi normale de moyenne  $\mu_k$  et de matrice des covariances  $\Sigma_k$ ). On note que le nombre de classes  $q$  est supposé connu.

Dans ces conditions, l'hypothèse d'absence de structure peut être celle de l'identité des diverses composantes  $f_k(x)$  de la densité  $f(x)$ .

De nombreux travaux ont été publiés sur l'estimation des mélanges de densités, dont on trouvera une synthèse dans Celeux (1992). Parmi les premiers travaux sur ce thème, citons ceux de Day (1969), de Bock (1974, 1977)<sup>1</sup>.

### - Modèles de partitions fixes

Ces modèles de référence supposent l'existence d'une partition inconnue en

$q$  classes ( $I_1, I_2, \dots, I_q$ ) d'effectifs respectifs ( $n_1, n_2, \dots, n_q$ ) avec  $\sum_{k=1}^q n_k = n$ .

A chacune des  $q$  classes  $I_k$  est associée une densité  $f_k(x)$ . Dans le cas où les densités  $f_k(x)$  sont celles de lois normales sphériques de même matrice des covariances  $\sigma^2 \mathbf{I}$  et de moyennes  $\mu_k$ , la partition qui réalise le maximum de vraisemblance est celle qui minimise le critère :

$$cr(q) = \sum_{k=1}^q \sum_{i \in I_k} \|x_i - \bar{x}_k\|^2 = \sum_{k=1}^q \sum_{i \in I_k} d^2(i, C_k)$$

<sup>1</sup> Cette formalisation donne lieu à beaucoup de travaux théoriques intéressants (cf. l'ouvrage de Everitt et Hand, 1981), mais peu d'entre eux débouchent sur des procédures utilisables en pratique pour valider les classifications ou déterminer le nombre de classes.



où l'on a noté, comme au § 2.1.2,  $C_k$  le centre de gravité de la classe  $I_k$ , de composantes  $\bar{x}_k$ . On reconnaît le critère utilisé dans l'agrégation autour de centres mobiles. La partition optimale exacte est actuellement impossible à déterminer, mais la méthode des centres mobiles, on l'a vu, conduit rapidement à un optimum local. Ce critère permet donc, dans le cadre fourni par ce modèle, d'évaluer la qualité d'une partition.

#### - Autre modèles

Une autre modélisation directe de l'hypothèse nulle est l'hypothèse d'homogénéité spatiale, développée par Dubes et Zeng (1987). Ces auteurs s'inspirent des tests de répartition spatiale aléatoire et de processus de Poisson généralisés (cf. par exemple Ripley, 1981) pour explorer, par des simulations extensives, les possibilités de ces épreuves de validation dans le domaine de la classification.

Aux critères qui permettraient de détecter l'existence d'une partition, on peut préférer les critères plus pragmatiques et modestes qui permettraient de comparer deux partitions (ou éventuellement d'améliorer une partition).

Parmi les critères les plus utilisés, citons le critère dit "critère F", quotient de la variance totale inter-classes par la variance totale intra-classes (traces des matrices  $E$  et  $D$ , matrices déjà rencontrées pour calculer les fonctions linéaires discriminantes au § 3.3.2), le critère dit "critère de Wilks", quotient des déterminants des deux matrices des covariances précédentes<sup>1</sup>.

### b – Combien de classes retenir ?

On présentera dans ce paragraphe les méthodes visant à déterminer, par des procédures empiriques (souvent inspirées par les modèles évoqués plus haut) le nombre de classes, sans faire intervenir d'information externe. On examinera tout d'abord le cas de la classification mixte qui a été exposée au chapitre 2, section 2.3.

Cette stratégie de classification adaptée à l'analyse exploratoire de grands tableaux (plusieurs milliers d'individus, plusieurs centaines de variables ou modalités) comporte en effet des possibilités de contrôle et de validation dans son processus même de calcul.

#### - Cas de la classification mixte

On reprendra les étapes du § 2.4.1.b, décrivant l'enchaînement analyse factorielle-classification mixte.

- 1- La première étape est une analyse en axes principaux, qui permet d'éprouver l'hypothèse d'absence éventuelle de structure, et donne une idée de l'éventuelle concentration du nuage de points à classer dans un

---

<sup>1</sup> On a pu établir la loi asymptotique des maxima de ces deux critères (maxima calculés sur toutes les partitions possibles) sous l'hypothèse nulle de distributions uniformes ou unimodales (cf. Bock, 1989).

sous-espace. Cette étape produit un système de coordonnées euclidiennes que les variables de départ soient numériques (analyse en composantes principales), fréquentielles (analyse des correspondances) ou nominales (analyse des correspondances multiples). On peut alors choisir de garder tous les axes correspondant à des valeurs propres non-nulles, ou de tronquer le support de façon à réaliser un filtrage. La possibilité de moduler le nombre d'axes permettra d'éprouver la stabilité des résultats de la classification qui va suivre.

2- La seconde étape est la classification mixte proprement dite.

2-1- Lors de la première phase (partition préliminaire par les centres mobiles destinées à réduire la dimension du problème), la possibilité de calculer des groupements stables (ou formes fortes, cf. Diday 1972) constitue une première épreuve de validation, fondée sur l'initialisation aléatoire de la méthode des centres mobiles.

2-2- La seconde phase (classification hiérarchique sur données agrégées utilisant le critère de Ward généralisé, bien adaptée à la classification de données regroupées) produit un dendrogramme et un histogramme des indices de niveau (schématisés sur les figures 2.3 - 1 et 2.3 - 2, par exemple) qui permettent d'apprécier les sauts importants de l'indice et donc de proposer, sur inspection visuelle, une coupure de l'arbre hiérarchique, à laquelle correspondra le nombre de classes retenu<sup>1</sup>. Si le critère d'inspection visuelle a été retenu ici, c'est par absence de consensus sur les autres critères, nombreux, qui ont été proposés dans la littérature<sup>2</sup>.

### - *Cas général*

Nous sommes vraiment ici, plus encore que dans le cas de la validation des méthodes factorielles, dans le domaine de la statistique expérimentale. Même si les développements théoriques sont parfois importants, il reste indispensable de tester empiriquement l'adéquation des résultats à la réalité, par simulation et bootstrap ou/et par essai sur des jeux de données-test homologués.

On mentionnera les premiers travaux de simulation de Gower et Banfield (1975), qui étudient empiriquement, à partir de plusieurs critères, la

---

<sup>1</sup> La consolidation de la coupure par réaffectation des individus (centres mobiles) donne également une information importante sur la qualité des résultats. Si l'évolution de la variance inter-classes (par exemple) est trop importante au cours de la consolidation, cela met en question la qualité de la coupure de l'arbre, qui se révèle loin d'un optimum local. Cela doit inciter à la prudence dans le maniement de la partition obtenue.

<sup>2</sup> Mollière (1986, 1989) propose également dans le cadre d'une stratégie d'agrégation mixte d'utiliser le coefficient CCC (Cubic Clustering Criterion) proposé par Sarle (1983) qui est une fonction de  $R^2$  (rapport de la variance interclasses à la variance totale) déterminée empiriquement. Ce coefficient CCC a été considéré comme satisfaisant à l'issue des simulations de Milligan et Cooper (1985).

distorsion entre la métrique initiale et l'ultramétrique produite par agrégation hiérarchique. Matusita et Ohsumi (1980) proposent un critère dit *d'affinité* pour comparer plusieurs partitions dans le cadre d'un algorithme à centres mobiles. Milligan et Cooper (1985) ont étudié et comparé plus de 30 tests et critères par simulation. Wong (1985), Jain et Moreau (1987) utilisent systématiquement le bootstrap pour étudier la stabilité des résultats et en déduire le nombre de classes stables. Hardy (1994) compare 7 critères appliqués à des résultats de classifications issues de six méthodes différentes, chaque couple classification-critère étant appliqué à 4 jeux de données artificiels différents choisis en raison de leurs aptitudes à représenter des situations typiques distinctes. Rasson et Kubushishi (1994) proposent un nouveau test (*Gap test*), fondé sur des processus de Poisson stationnaires, qui utilise les éventuelles zones vides entre classes. Testé sur des jeux de données simulées ou classiques, il est efficace pour reconnaître les classes isolées.

### c – Les critères externes

Comme le souligne Bock (1994), il ne faudrait pas exagérer la pertinence et l'importance de la notion de nombre de classes d'une classification, car une classification n'est jamais une fin en soi. C'est beaucoup plus souvent d'une dissection dont on a besoin, selon la terminologie de Kendall (1966) qui considère qu'un découpage de la réalité multidimensionnelle est toujours utile, même si les classes ne sont pas bien séparées, même si tous les individus ne sont pas classés.

Que signifie alors un critère global de qualité, qui pourrait nous faire rejeter des traits structuraux importants ? Et quels modèles théoriques pourraient rendre compte d'une situation aussi complexe ?

William et Lance (1965) pensent qu'une classification "ne peut pas être vraie ou fausse, ni probable ou improbable, mais seulement profitable ou non profitable". Cette notion de profitabilité ne peut qu'être externe au tableau de données. Elle est liée au contexte et aux objectifs de la recherche ou de l'étude, aux méta-données (meta-data), c'est-à-dire à l'information sur l'information.

Les procédures de description automatique des classes (cf. section 2.3) à partir des variables actives ayant créé la partition, mais aussi à partir de toute l'information externe disponible (ayant le statut de variables supplémentaires, numériques ou nominales) sont des procédures de validation potentielles. Elles nous disent que telle portion connexe de l'espace engendré par les variables actives présente de l'intérêt vis-à-vis d'autres informations présentes dans la base de données. De même que sur un histogramme unidimensionnel, on peut identifier certaines zones à partir d'information extérieure sur les individus.

# Bibliographie

- Agrawala A.K. (Ed.) (1977) - *Machine Recognition of Patterns*. IEEE Press, New York.
- Agresti A. (1990) - *Categorical Data Analysis*. J. Wiley, Chichester.
- Agresti A. (1992) - A survey of exact inference for contingency tables. *Statistical Science*, 7, 1, p 131-177.
- Aitkin M. A. (1979) - A simultaneous test procedure for contingency tables. *Appl. Statist.*, 28, p 233-242.
- Aitchison J. (1983) - Principal component analysis of compositional data. *Biometrika*, 70, (1), p 57-65.
- Aitchison J., Aitken C. G. G. (1976) - Multivariate binary discrimination by the kernel method. *Biometrika*, 63, p 413-420.
- Akaike H. (1973) - Information theory and an extension of the maximum likelihood principle. In : *Second Internat. Symp. on Information Theory*, Petrov B.N., Czaki F., eds., Akademiai Kiado, Budapest, p 267-281.
- Aluja Banet T., Lebart L. (1984) - Local and partial principal component analysis and correspondence analysis. In : *COMPSTAT, Proceedings in Computational Statistics*, Physica Verlag, Vienna, p 113-118.
- Amari S. (1990) - Mathematical foundation of neurocomputing. *Proc of the IEEE*, 78, n°9.
- Anastassakos I., D'Aubigny G. (1984) - L'utilisation de tests de sphéricité pour la recherche de la dimension de l'espace latent en analyse factorielle classique et en analyse en composantes principales. *Revue Statist. Appl.*, 32, (2), p 45-57.
- Anderberg M.R. (1973) - *Cluster Analysis for Applications*. Academic Press, New York.
- Anderson J.A. (1982) - Logistic Discrimination. in : *Handbook of Statistics*, 2, Krishnaiah P.R. and Kanal L. (Eds) North Holland, Amsterdam, p 169-191.
- Anderson T. W. (1951) - The asymptotic distribution of certain characteristic roots and vectors. *Proc. of the 2nd Berkeley Symp. on Math. Statist. and Prob.*, p 103-130, Univ. of California Press.
- Anderson T.W. (1958) - *An Introduction to Multivariate Statistical Analysis* (Second edition : 1984). J. Wiley, New York.
- Anderson T. W. (1963) - Asymptotic theory for principal component analysis. *Ann. Math. Statist.*, 34, p 122-148.
- Anderson T. W., Rubin H. (1956) - Statistical Inference in factor analysis. *Proc. of the 3rd Berkeley Symp. on Math. Statist.*, 5, p 111-150.
- Andrews D. F. (1972) - Plots of High-dimensional data. *Biometrics*, 28, p 125-136.
- Arabie P. (1978) - Constructing blockmodels : how and why. *J. of Math. Psychology*, 17, (1), p 21-63.
- Arabie P. (1991) - Was Euclid an unnecessarily sophisticated psychologist?. *Psychometrika*, 56, p 567-587.
- Art D., Gnanadesikan R, Kettenring J.R. (1982) - Data based metrics for cluster analysis. *Utilitas Mathematica*, 21 A, p 75-99.
- ASU, (Lebart L., ed.) (1992) - *La qualité de l'information dans les enquêtes*. Dunod, Paris.
- Atkinson A.C. (1981) - Likelihood ratios, posterior odds and information criteria. *J. Econometrics*, 16, p 15-20.

- Atkinson A.C. (1985) - *Plots, Transformation and Regression : an Introduction to Graphical Methods of Diagnostic Regression Analysis*. Clarendon Press, Oxford.
- Babeau A., Lebart L. (1984) - Les conditions de vie et aspirations des Français. *Futuribles*, 1, p 37-53.
- Bailey R.A. (1981) - A unified approach to Design of Experiments. *J. Royal Statist. Soc. (A)*, 144(2), p 214-233.
- Balbi S. (1994) - *L'Analisi Multidimensionale dei dati negli anni'90*. Dipartimento di Matematica e Statistica. (Univ. Federico II), Rocco Curto Editore, Napoli.
- Baldi P., Hornik K. (1989) - Neural networks and principal component analysis : learning from examples without local minima. *Neural Networks*, 2, p 52-58.
- Ball G.H., Hall D.J. (1965) - *ISODATA, A Novel Method of Data Analysis and Pattern Classification*. AD 699616, Stanford Research Institute, Menlo Park, California.
- Ball G.H., Hall D.J. (1967) A clustering technique for summarizing multivariate data. *Behavioral Sciences*, 12, p 153-155.
- Ballif J.-F. (1986) - *Analyse multivariée : un modèle descriptif général*. Univ. de Lausanne, Peter Lang, Berne.
- Bardos M. (1984) - Le risque de défaillance d'entreprise. *Cahiers Economiques et Monétaires*. 19, p 1-190.
- Bardos M. (1989) - Trois méthodes d'analyse discriminante. *Cahiers Economiques et Monétaires*. 33, p 151-190.
- Barnett V. (1976) - The ordering of multivariate data. *J. Royal Statist. Soc. (A)*, 139, p 318-354.
- Bartlett M.S. (1950) - Tests of significance in factor analysis. *British J. Psych. (Stat. Section)*, 3, p 77-85.
- Bartlett M.S. (1951) - The effect of standardization on  $\chi^2$  approximation in factor analysis (with an appendix by W. Lederman). *Biometrika*, 38, p 337-344.
- Beltrami E. (1873) - Sulle funzioni bilineari. *Giorn. Math. Battaglin*. 11, p 98-106.
- Benali H., Escofier B. (1987) - Stabilité de l'analyse factorielle des correspondances multiples en cas de données manquantes et modalités à faibles effectifs. *Revue Statist. Appl.*, 35, n°1, p 41-52.
- Benali H., Escofier B. (1990) - Analyse factorielle lissée et analyse des différences locales. *Revue Statist. Appl.* 38, 2, p 55-76.
- Benasseni J. (1986a) - Stabilité de l'analyse en composantes principales par rapport à une perturbation des données. *Revue Statist. Appl.*, 35, 3, p 49-64.
- Benasseni J. (1986b) - Stabilité en ACP par rapport aux erreurs de mesure. In : *Data Analysis and Informatics*, 4, Diday E. et al. (eds), North-Holland, Amsterdam, p 523-533.
- Benzécri J.-P. (1969a) - Statistical analysis as a tool to make patterns emerge from clouds. In : *Methodology of Pattern Recognition* (S.Watanabe, Ed.) Academic Press, p 35-74.
- Benzécri J.-P. (1969b) - Approximation stochastique dans une algèbre normée non commutative. *Bull. Soc. Math. France*, 97, p 225-241.
- Benzécri J.-P. (1973) - *L'Analyse des Données*. Tome 1: *La Taxinomie*. Tome 2: *L'Analyse des Correspondances* (2<sup>de</sup> éd. 1976). Dunod, Paris.
- Benzécri J.-P. (1974) - La place de l'a priori. In : *Organum- Encyclopaedia Universalis*. Paris.
- Benzécri J.-P. (1977 a) - Analyse discriminante et analyse factorielle. *Les Cahiers de l'Analyse des Données*, 4, p 369-406.
- Benzécri J.-P. (1977 b) - Choix des unités et des poids dans un tableau en vue d'une analyse des correspondances. *Cahiers de l'Analyse des Données*, 2, p 333-352.

- Benzécri J.-P. (1979) - Sur le calcul des taux d'inertie dans l'analyse d'un questionnaire. *Cahiers de l'Analyse des Données*, 4, p 377-378 .
- Benzécri J.-P. (1982 a) - *Histoire et préhistoire de l'analyse des données*. Dunod, Paris.
- Benzécri J.-P. (1982 b) - Sur la généralisation du tableau de Burt et son analyse par bandes. *Cahiers de l'Analyse des Données*, 7, p 33-43.
- Benzécri, J.-P. (1982 c) - Construction d'une classification ascendante hiérarchique par la recherche en chaîne de voisins réciproques. *Cahiers d'Analyse des Données*, 7, p 209-218.
- Benzécri, J.-P. (1983) - Analyse d'inertie intraclasse par l'analyse d'un tableau de correspondance. *Les Cahiers d'Analyse des Données*, 8, p 351-358.
- Benzécri J.-P. (1992) - *Correspondence Analysis Handbook*. Marcel Dekker, New York.
- Benzécri J.-P., Cazes P. (1978) - Problème sur la classification. *Les Cahiers de l'Analyse des Données*, 3, 1, p 95-101.
- Benzécri J.-P., Jambu M. (1976) - Agrégation suivant le saut minimum et arbre de longueur minimum. *Les Cahiers de l'Analyse des Données*, 1, p 441-452.
- Benzécri, J.-P., Lebeaux M.-O., and Jambu M. (1980) - Aides à l'interprétation en classification automatique. *Les Cahiers de l'Analyse des Données*, 5, p 101-123.
- Beran R., Srivastava M.S. (1985) - Bootstrap test and confidence region for functions of a covariance matrix. *Ann. of Statist.*, 13, p 95-115.
- Berge C. (1963) - *Théorie des graphes et ses applications*. Dunod, Paris.
- Berge C. (1973) - *Graphs and Hypergraphs*. North Holland, Amsterdam.
- Berk R.H. (1972) - Consistency and asymptotic normality of MLE's for exponential models. *Ann. Math. Statist.*, 43, p 193-204.
- Berry W. D. (1984) - *Non recursive causal models*. Sage, Beverly Hills.
- Bertin J. (1973) - Article : "Graphique (représentation -)". In : *Encyclopaedia Universalis*.
- Bertrand P., Diday E. (1990) - Une généralisation des arbres hiérarchiques : les représentations pyramidales. *Revue Statist. Appl.*, 38, (3), p 53-78.
- Besley D. A., Kuh E., Welsh R. E. (1980) - *Regression Diagnostics: Identifying Influential Data and Sources of Colinearity*, J. Wiley, New York.
- Besse P., Ferré L. (1993) - Sur l'usage de la validation croisée en analyse en composantes principales. *Revue Statist. Appl.*, 41, (1), p 71-76.
- Birch M. W. (1963) - Maximum likelihood in three-way contingency tables. *J. Royal Statist. Soc. (B)*, 25, p 220-233.
- Bishop Y., Fienberg S., Holland P. (1975) - *Discrete Multivariate Analysis*. MIT Press, Cambridge, Mass.
- Bock H. H. (1974) - *Automatische Klassifikation. Theoretische und praktische Methoden zur Gruppierung und Strukturierung van Daten (Cluster Analysis)*. Vandenhoeck & Ruprecht, Göttingen.
- Bock H. H. (1977) - On tests concerning the existence of a classification. In : *First International Symposium on Data Analysis and Informatics*. INRIA, Rocquencourt, p 449-464.
- Bock H. H. (1979) - Simultaneous clustering of objects and variables. In : *Analyse des données et informatique*. European C.C. Courses, INRIA, p 187-203.
- Bock H. H. (1985) - On some significance tests in cluster analysis. *J. of Classification*, 2, p 77-108.
- Bock H. H. (1989) - Probabilistic aspects in cluster analysis. In : *Conceptual and numerical analysis of data*. Opitz O. (ed.), Springer-Verlag, Berlin, Heidelberg.
- Bock H. H. (1994) - Classification and clustering : Problems for the future. In : *New Approaches in Classification and Data Analysis*, Diday E. et al. (eds), Springer Verlag, Berlin, p 3-24.

- Boeswillwald E. (1992) - L'expérience du CESP en matière de qualité des mesures d'audience. In : *La qualité de l'information dans les enquêtes*, (ASU), Dunod, Paris, p 313-341.
- Bourlard H., Kamp Y. (1988) - Auto-association by multi-layers perceptrons and singular value decomposition. *Biological Cybernetics*, 59, p 291-294.
- Bouroche J.-M., Tenenhaus M. (1970) - Quelques méthodes de segmentation. *RAIRO*, 5, 2, p 29-42.
- Bouroche J.-M., Saporta G. (1980) - *L'analyse des données*. coll."Que sais-je", n°1854, PUF, Paris .
- Bourret P., Reggia J., Samuelides M. (1991) - *Réseaux Neuronaux*. Teknea, Toulouse.
- Box G. E. P., Cox D. R., (1982) - An analysis of transformations revisited, rebutted. *J. Amer. Statist. Assoc.*, 77, p 209-210.
- Breiman L., Friedman J. H., Ohlsen R. A., Stone C. J. (1984) - *Classification and Regression Trees*. Wadsworth, Belmont.
- Brent R.P. (1974) - A gaussian pseudo-random number generator. *Com. ACM*, 17, p 704-706.
- Brillouin L. (1959) - *La science et la théorie de l'information*. Masson, Paris.
- Bruynooghe M. (1978) - Classification ascendante hiérarchique des grands ensembles de données : un algorithme rapide fondé sur la construction des voisinages réductibles. *Les Cahiers de l'Analyse des Données*, 3, p 7-33.
- Burt C. (1950) - The factorial analysis of qualitative data. *British J. of Statist. psychol.* 3, 3, p 166-185.
- Burtschy B., Lebart L. (1991) - Contiguity analysis and projection pursuit. In : *Applied Stochastic Models and Data Analysis*, R. Gutierrez and M.J.M. Valderrama, Eds, World scientific, Singapore, p 117-128.
- Cacoullos T. (Ed.) (1973) - *Discriminant Analysis and Applications*. Academic Press, New York.
- Caillez F., Pagès J.P. (1976) - *Introduction à l'Analyse des Données*. S.M.A.S.H., Paris.
- Callant C. M. (1991) - *Technique de Lissage et de Régularisation en Analyse Discriminante*. Thèse. Université Paris IX, Dauphine, (Publ. INRIA TU177), Paris.
- Caraux G. (1984) - Réorganisation et représentation visuelle d'une matrice de données numériques : un algorithme itératif. *Revue de Statist. Appl.*, 32, p 5-24.
- Carlier A. (1985) Analyse des évolutions sur tables de contingences, quelques aspects opérationnels. In : *Data Analysis and Informatics*, Diday E. et al. (eds), North Holland, Amsterdam, p 421-428.
- Carlier A., Lavit C., Pagès M., Pernin M.-O., Turlot J.-C. (1988) - A comparative review of methods which handles a set of indexed data tables. In : *Multiway Data Analysis*, Coppi R., Bolasco S. (eds), North Holland, Amsterdam, p 85-102.
- Carrol J. D. (1968) - Generalization of canonical correlation to three or more sets of variables. *Proc. Amer. Psychological Assoc.* p 227-228.
- Carrol J. D., Chang J. J. (1970) - Analysis of individual differences in multidimensional scaling via an n-way generalization of Eckart-Young' decomposition. *Psychometrika*, 35, p 283-319.
- Carrol J. D., Pruzansky S., and Green P. F. (1977) - Estimation of the parameters of Lazarsfeld's Latent Class Model by application of canonical decomposition CANDECOP to multi-way contingency tables. *AT&T Bell Laboratories*, unpublished paper.
- Casin P., Turlot J.-C. (1986) - Une présentation de l'analyse canonique généralisée dans l'espace des individus. *Revue Statist. Appl.* 35, (3), p 65-75.
- Cattell R.B. (1966) - The scree test for the number of factors. *Mult. Behavioral Research*, 1, p 245-276.

- Caussinus H. (1992) - Projections révélatrices. In : *Modèles pour l'analyse des données multidimensionnelles*. J.J. Dreesbeke, B. Fichet, P.Tassi, eds, Economica, Paris.
- Caussinus H., Ruiz A. (1990) - Interesting projections of multidimensional data by means of generalized principal component analysis. In : *COMPSTAT 90*, Physica Verlag, Heidelberg, p 121-126.
- Cazes P. (1977) - Etude des propriétés extrémales des sous-facteurs issus d'un sous-tableau d'un tableau de Burt. *Les Cahiers de l'Analyse des Données*, 2, p 143-160.
- Cazes P. (1980) - Analyse de certains tableaux rectangulaires décomposés en blocs. *Les Cahiers de l'Analyse des Données*, 5, p 145-161, et p 387-403.
- Cazes P. (1981) - Note sur les éléments supplémentaires en analyse des correspondances. *Les Cahiers de l'Analyse des Données*, 1, p 9-23; 2, p 133-154.
- Cazes P. (1982) - Analyse de certains tableaux rectangulaires décomposés en blocs : Codage simultané de variables qualitatives et quantitatives. *Les Cahiers de l'Analyse des Données*, 6, p 9-18.
- Cazes P. (1984) - Correspondance hiérarchiques et ensembles associés. *Cahiers du B.U.R.O.*, n° 43-44, Université Pierre et Marie Curie, p 43-142.
- Cazes P. (1986 a) - Une généralisation des correspondances multiples et des correspondances hiérarchiques. *Cahiers du B.U.R.O.*, 46-47, Université Pierre et Marie Curie, p 37-64.
- Cazes P. (1986 b) - Correspondance entre deux ensembles et partition de ces deux ensembles. *Les Cahiers de l'Analyse des Données*, 11, p 335-340.
- Cazes P. (1990) - Codage d'une variable continue en vue de l'analyse des correspondances. *Revue Statist. Appl.*, 38, 3, p 35-51.
- Cazes P., Chessel D., Doledec S. (1988) - L'analyse des correspondances interne d'un tableau partitionné : son usage en hydrobiologie. *Revue Statist. Appl.* 36, (1), p 39-54.
- Cazes P., Moreau J. (1991) - Contingency table in which the rows and columns have a graph structure. In : E.Diday, Y.Lechevallier (Eds) *Symbolic-Numeric Data Analysis and Learning*, Nova Science Publishers. New York, p 271-280.
- Celeux G. (ed) (1990) - *Analyse discriminante sur variables continues*. INRIA, Rocquencourt.
- Celeux G. (1992) - Résultats asymptotiques et validation en classification. In : *Modèles pour l'analyse des données multidimensionnelles*. J.J. Dreesbeke, B. Fichet, P.Tassi, eds, Economica, Paris.
- Celeux G., Diday E., Govaert G., Lechevallier Y., Ralambondrainy H. (1989) - *Classification automatique des données: environnement statistique et informatique*. Dunod, Paris.
- Celeux G., Hébrail G., Mkhadri A., Suchard M. (1991). Reduction of a large scale and ill-conditioned statistical problem on textual data. In : *Applied Stochastic Models and Data Analysis, Proceedings of the 5th Symposium*. In : ASMDA, Gutierrez R. and Valderrama M.J. Eds, World Scientific, p 129-137.
- Celeux G., Nakache J.-P. (eds) (1994) - *Analyse discriminante sur variables qualitatives*. Polytechnica, Paris.
- Chabanon C., Dubuisson B. (1991) - Méthodes non probabilistes. In : *Analyse discriminante sur variables continues*, Celeux G. (ed.), INRIA, Paris.
- Chandon J.-L., Pinson S. (1981) - *Analyse typologique : Théorie et applications*. Masson, Paris.
- Chateau F. (1994) - Probabilités a priori inégales dans la règle des k plus proches voisins. *Actes des XXVIèmes Journées de Statistiques* (Neuchâtel), p 195-198.
- Chatterjee S., Price B. (1991) - *Regression Analysis by Examples*. J. Wiley, New York.
- Cheng B., Titterton D.M. (1994) - Neural networks: a review from a statistical perspective. *Statistical Science*, 9, n°1, p 2-54.



- Chernoff H. (1973) - The use of faces to represent points in  $k$ -dimensional space graphically. *J. Amer. Statist. Assoc.*, **68**, p 361-368.
- Chessel D., Lebreton J.-D., Yoccoz N. (1987) - Propriétés de l'analyse canonique de correspondances ; une illustration en hydrobiologie. *Revue de Statist. Appl.*, **35**, (4), p 55-72.
- Choudary Hanumara R. Thompson W.A. (1968) - Percentage points of the extreme roots of a Wishart matrix. *Biometrika*, **55**, p 505-512.
- Christensen R. (1990) - *Log-Linear Models*. Springer-Verlag, New York.
- Clemm D.S., Krishnaiah P.R., Waikar V.B. (1973) - Tables of the extreme roots of a Wishart matrix. *J. of Statist. Comput. and Simul.* **2**, p 65-92.
- Cliff N. (1966) - Orthogonal rotation to congruence. *Psychometrika*, **31**, p 33-42.
- Cliff A.D. and Ord J.K. (1981) - *Spatial Processes : Models and Applications*. Pion, London.
- Cochran W.G., Cox G.M. (1957) - *Experimental Design* (2nd ed.). J. Wiley, New York.
- Cohen A. (1980) - On the graphical display of the significant components in two-ways contingency tables. *Comm. in Statistics, Theory.Meth.*, **A9** (10), p 1025-1041.
- Cohen J. (1967) - *Statistical Power Analysis for the Behavioral Sciences*. Academic Press, New York.
- Cook R.D., Weisberg S. (1982) - *Residuals and Influence in Regression*. Chapman and Hall, London.
- Cook R.D., Weisberg S. (1994) - *An Introduction to Regression Graphics*. J. Wiley, New York.
- Coppi R., Bolasco S. (eds) (1989) - *The Analysis of Multiway Data Matrices*. North Holland, Amsterdam.
- Cormack R.M. (1971) - A review of classification. *J. of Royal Statist. Society, Serie A*, **134**, Part. 3, p 321-367.
- Cornfield J. (1962) - Joint dependence of risk of coronary heart disease on serum cholesterol and systolic blood pressure: a discriminant function approach. *Fed. Amer. Socs. Exper. Biol. Proc. Suppl.*, **11**, p 58-61.
- Corsten L. C. A. (1976) - Matrix approximation, a key to application of multivariate methods. In : *Proc. 9th Int. Biometric Conf.*, **1**, p 61-77, Raleigh, North Carolina.
- Cottrell M., Fort J.-C. (1987) - Etude d'un algorithme d'auto-organisation. *Ann. de l'Inst. Henri Poincaré*, **23**, p 1-20.
- Cox D.R. (1958) - *Planning of Experiments*. J. Wiley, New York.
- Cox D. R. (1972) - *Analyse des données binaires*. Dunod, Paris.
- Cox D. R. (1977).- The role of significance tests. *Scandinavian Journal of Statist.*, **4**, p 49-70.
- Craddock J.M., Flood C.R. (1970) - The distribution of the  $\chi^2$  statistic in small contingency tables. *Appl. Statist.*, **19**, p 173-181.
- Cramer H. (1946) - *Mathematical Methods of Statistics*. Princeton University Press, Princeton.
- Critchley F. (1985) - Influence in principal component analysis. *Biometrika*, **72**, p 626-636.
- Dagnelie P. (1981) - *Principes d'expérimentation*. Les Presse Agronomiques de Gembloux, Gembloux.
- Darmois G. (1957) - *Statistique et applications*. Armand Colin, Paris.
- Darroch J. N., Lauritzen S. L., Speed T. P. (1980) - Markov field and log-linear interaction models for contingency tables. *Ann. of Statist.*, **8**, 522-539.
- Daudin J.-J., Duby C., Trécourt P. (1988) - Stability of principal components studied by the bootstrap method. *Statistics*, **19**, p 241-258.

- Daudin J.-J., Trécourt P. (1980) - Analyse factorielle des correspondances et modèle log-linéaire : Comparaison des deux méthodes sur un exemple. *Revue Statist. Appl.* 28, n° 1, p 5-24.
- Davis A. W. (1977) - Asymptotic theory for principal component analysis : the non-normal case. *Australian J. of Statist.*, 19, p 206-212.
- Davis C., Kahan W. M. (1970) - The rotation of eigenvectors by a perturbation. *Journal of SIAM (Numerical Analysis)*, 7, p 1-46.
- Day N. E. (1969) Estimating the component of a mixture of normal distribution. *Biometrika*, 56, p 463-474.
- Delecroix M. (1983) - *Histogrammes et estimation de densité*. P.U.F., Paris.
- Deming W. E., Stephan F. F. (1940) - On a least squares adjustment of a sampled frequency table when the expected marginal total are known. *Ann. Math. Statist.*, 11, p 427-444.
- Dempster A.P. (1971) - An overview of multivariate data analysis. *J. Mult. Analysis*, 1, p 316-346.
- Deroo M., Dussaix A.-M. (1980) - *Pratique et analyse des enquêtes par sondage*. P.U.F., Paris.
- Devaud J.-M. (1985) - Discrimination et description sur variables qualitatives : un exemple comparatif sur données réelles. *Revue Statist. Appl.* 33, n° 2, p 5-18.
- Devijver P., Kittler J. (1982) - *Pattern Recognition : A Statistical Approach*. Prentice Hall, New York.
- Deville J.-C., Malinvaud E. (1983) - Data analysis in official socio-economic statistics. *J. Royal Statist. Soc.*, A, 146, part 4.
- Deville J.-C., Särndal C.-E. (1992) - Calibration estimator in Survey Sampling. *J.A.S.A.*, 87, 418, p 376-382.
- Diaconis P., Efron B. (1983) - Computer intensive methods in statistics. *Scientific American*, 248, (May), p 116-130.
- Diday E. (1972) - Optimisation en classification automatique et reconnaissance des formes. *Revue Française de Recherche Opérationnelle*, 3, p 61-96.
- Diday E. (1974) - Classification automatique séquentielle pour grands tableaux. *Revue Fr. Inf. Rech. Opér.* 9, (Mars 1975), p 1-29.
- Diday E. (1992) - From data to knowledge : Probabilist objects for a symbolic data analysis. In : *Computational Statistics*, Dodge Y., Whittaker J. (Eds), Physica Verlag, Heidelberg, p 193-214.
- Diday E. (1971) - La méthode des nuées dynamiques. *Revue Statist. Appl.* 19, n° 2, p 19-34.
- Diday E., Lemaire J.L., Pouget J., Testu F. (1982) - *Eléments d'Analyse des Données*. Dunod, Paris.
- Dobson A. (1983) - *An Introduction to Statistical Modelling*. Chapman and Hall, New York.
- Dodge Y. (ed.) (1987) - *Statistical Data Analysis Based on the  $L_1$ -Norm and Related Methods*. North Holland, Amsterdam.
- Domenges D., Volle M. (1979) - Analyse factorielle sphérique: Une exploration. *Annales de l'INSEE*, n° 35.
- Draper N. R., Smith H. (1981) - *Applied Regression Analysis (2nd ed.)*. J. Wiley, New York.
- Droesbeke J.-J., Fichet B., Tassi P. (ed.) (1987) - *Les sondages*. Economica, Paris.
- Droesbeke J.-J., Fichet B., Tassi P. (ed.) (1992) - *Modèles pour l'analyse des données multidimensionnelle*. Economica, Paris.

- Droesbeke J.-J., Tassi P. (1990) - *Histoire de la statistique*. Que-sais-je? PUF, Paris.
- Drouet d'Aubigny G. (1993) - Analyse des proximités et programmes de codage multidimensionnel. *La Revue de Modulad*, INRIA, Rocquencourt, **12**, p 1-32.
- Dubes R., Jain A. K. (1979) - Validity studies in clustering methodology. *Pattern Recognition*, **11**, p 235-254.
- Dubes R. C., Zeng G. (1987) - A test for spatial homogeneity in cluster analysis. *J. of Classification*, **4**, p 33-56.
- Dubuisson B. (1990) - *Diagnostic et reconnaissance des formes*. Hermès, Paris.
- Duda R.O., Hart P.E. (1973) - *Pattern Classification and Scene Analysis*. J. Wiley, New York.
- Dugué D. (1958) - *Traité de statistique théorique et appliquée*. Masson, Paris.
- Eastment H. T., Krzanowski W., J. (1982) - Cross validatory choice of the number of components of a principal component analysis. *Technometrics*, **24**, p 73-77.
- Eckart C., Young G. (1936) - The approximation of one matrix by another of lower rank. *Psychometrika*, **1**, p 211-218.
- Eckart C., Young G. (1939) - A principal axis transformation for non-hermitian matrices. *Bull. Amer. Math. Assoc.*, **45**, p 118-121.
- Edwards A. W. F., Cavalli-Sforza L. L. (1965) - A method for cluster analysis. *Biometrics*, **21**, p 362-375.
- Efron B. (1965) - The convex hull of a random set of points. *Biometrika*, **52**, p 331-343.
- Efron B. (1979) - Bootstraps methods : another look at the Jackknife. *Ann. Statist.*, **7**, p 1-26.
- Efron B. (1982) - *The Jackknife, the Bootstrap et other Resampling Plans*. SIAM, Philadelphia.
- Efron B., Tibshirani R. J. (1993) - *An Introduction to the Bootstrap*. Chapman and Hall, New York.
- Engelman L., Hartigan J. A. (1969) - Percentage points of a test for clusters. *J. Amer. Statist. Assoc.*, **64**, p 1647-1648.
- Enyukov I. S. (1988) - Detecting structure by mean of projection pursuit. *COMPSTAT Proceedings*, Physica Verlag, Heidelberg, p 47-58.
- Escofier B. [Cordier B.] (1965) - *l'Analyse des correspondances*. Thèse, Faculté des Sciences de Rennes ; publiée en 1969 dans les *Cahiers du Bureau Universitaire de Recherche Opérationnelle*, n°13.
- Escofier B. (1978) - Analyse factorielle et distances répondant au principe d'équivalence distributionnelle. *Revue de Statist. Appl.*, **26**, p 29-37.
- Escofier B. (1979 a) - *Stabilité et approximation en analyse factorielle*. Thèse d'Etat, Université Pierre et Marie Curie, Paris.
- Escofier B. (1979 b) - Traitement simultané de variables qualitatives et quantitatives. *Les Cahiers de l'Analyse des Données*, **4**, (2), p 137-146.
- Escofier B. (1979 c) - Une représentation des variables dans l'analyse des correspondances multiples. *Revue de Statist. Appl.*, **27**, p 37-47.
- Escofier B. (1984) - Analyse factorielle en référence à un modèle: application à l'analyse des tableaux d'échanges. *Revue de Statist. Appl.*, **32**, 25-36.
- Escofier B. (1987) - Analyse des correspondances multiples conditionnelles. In : *Data Analysis and Informatics*, Diday E. (ed.), North Holland, Amsterdam, p 13-22.
- Escofier B. (1989) - Multiple correspondence analysis and neighboring relation. In : *Data Analysis, Learning Symbolic and Numeric Knowledge*. Diday E. (ed.), Nova Science Publishers, New York, p 55-62.
- Escofier B., Leroux B. (1972) - Etude de trois problèmes de stabilité en analyse factorielle. *Publication de l'Institut Statistique de l'Université de Paris*, **11**, p 1-48

- Escoufier B., Pagès J. (1983) - Méthode pour l'analyse de plusieurs groupes de variables. Application à la caractérisation des vins rouges du Val de Loire. *Revue Statist. Appl.* 31, p 43-59.
- Escoufier B., Pagès J. (1984) - Analyses factorielles multiples. *Cahiers du BURO*, 2, ISUP, Paris.
- Escoufier B., Pagès J. (1988) - *Analyses factorielles multiples*. Dunod, Paris.
- Escoufier Y. (1970) - Echantillonnage dans une population de variables aléatoires réelles. *Publication de l'Institut Statistique de l'Université de Paris*, 19, Fasc 4, p 1-47.
- Escoufier Y. (1980) - L'analyse conjointes de plusieurs matrices de données. In: *Biométrie et Temps*, Jolivet et al. (eds), p 59-76.
- Escoufier Y. (1982) - L'Analyse des correspondances simples et multiples. *Metron*, 1-2, p 53-78.
- Escoufier Y. (1985 a) - Objectifs et procédures de l'analyse conjointe de plusieurs tableaux. *Statist. et Anal. des Données*. 10, (1), p 1-10.
- Escoufier Y. (1985 b) - L'Analyse des correspondances, ses propriétés, ses extensions. *Bull. of the Int. Statist. Inst.*, 4, p 28-2.
- Escoufier, Y. (1988) - Beyond correspondence analysis. In: *Classification and Related Methods of Data Analysis*, H.H.Bock, Ed., North Holland, p 505-514.
- Everitt B. S., Hand D. J. (1981) - *Finite Mixture Distributions*. Chapman and Hall, London.
- Falguerolles (de) A., Jmel S. (1993) - Un modèle graphique pour la sélection de variables qualitatives. *Revue de Statist. Appl.* 41, p 23-41.
- Falissard B. (1995) - Déploiement d'une matrice de corrélation sur la sphère unité de  $R^3$ . *Revue de Statist. Appl.*, 43, (2) p 35-48.
- Faraj A. (1993) - Analyse de contiguïté: une analyse discriminante généralisée à plusieurs variables qualitatives. *Revue Statist. Appl.* 41, (3), p 73-84.
- Farebrother R. W. (1987) - The historical development of the  $L_1$  and  $L_\infty$  estimation procedures. in: *Statistical Data Analysis Based on the  $L_1$ -Norm and Related Methods*, Dodge Y. ed., North Holland, Amsterdam, p 37-64.
- Fichet B. (1987) - The role played by  $L_1$  in data analysis. in: *Statistical Data Analysis Based on the  $L_1$ -Norm and Related Methods*, Dodge Y. ed., North Holland, Amsterdam, p 185-194.
- Fichet B. (1988) -  $L_p$  space in Data Analysis. In: *Classification and Related Methods of Data analysis*. Boch H. H. (ed.), North-Holland, Amsterdam, p 439-444.
- Fienberg S.E. (1980) - *The Analysis of Cross-classified Categorical Data*. MIT Press, Cambridge, Mass.
- Fine J. (1992) - Modèles graphiques d'associations. In: ASU, (Droesbeke J.-J., Fichet B., Tassi P., ed.), *Modèles pour l'analyse des données multidimensionnelle*, Economica, Paris.
- Fine J. (1993) - Problèmes d'indétermination en analyse en facteurs et analyse en composantes principales optimale. *Revue de Statist. Appl.*, 41, (4), p 45-72.
- Fisher R. A. (1915) - Frequency distribution of the value of the correlation coefficient in samples from an indefinitely large population. *Biometrika*, 10, p 507-521.
- Fisher R. A. (1935) - *The Design of Experiments*. Oliver and Boyd, Edinburgh.
- Fisher R. A. (1936) - The use of multiple measurements in taxonomic problems. *Ann. of Eugenics*, 7, p 179-188.
- Fisher R.A. (1939) - The sampling distribution of some statistics obtained from non linear equations. *Ann. Eugen.*, 7, p 179-188.
- Fisher R.A. (1940) - The precision of discriminant functions. *Ann. Eugen.*, 10, p 422-429.

- Fisher R.A., Yates F. (1949) - *Statistical Tables for Biological, Agricultural and Medical Research*. Hafner Publishing Company.
- Fisher W.D. (1958) - On grouping for maximum homogeneity. *J. of Amer. Statist. Assoc.*, **53**, p 789-798.
- Fix E., Hodges J. L. (1951) - Discriminatory analysis - nonparametric discrimination : consistency properties. Report of the U.S.A.F. School of Aviation Medicine. In : Agrawala (1977).
- Florek K. (1951) - Sur la liaison et la division des points d'un ensemble fini. *Colloq. Math.*, **2**, p 282-285.
- Flury B. (1988) - Common principal components and related multivariate models. J. Wiley, New York.
- Forgy E. W. (1965) - Cluster analysis of multivariate data : efficiency versus interpretability of classifications. *Biometric Society Meetings*, Riverside, California (Abstract in : *Biometrics* **21**, 3, p 768).
- Fortin M. (1975) - Sur un algorithme pour l'analyse des données et la reconnaissance des formes. *Revue de Statist. appl.*, **23**, p 37-46.
- Fourgeaud C., Lenclud B. (1978) - *Econométrie*. P.U.F., Paris.
- Francisco C. A., Finch M. D. (1980) - A comparison of methods used for determining the number of factors to retain in factor analysis. *Technometrics*, **22**, p 105-110.
- Friedman J. H. (1987) - Exploratory projection pursuit. *J. of Amer. Statist. Assoc.*, **82**, (397), p 249-266.
- Friedman J. H. (1989) - Regularized discriminant analysis. *J. of Amer. Statist. Assoc.*, **84**, p 165-175.
- Friedman J. H., and Tukey J.W. (1974) - A Projection pursuit algorithm for exploratory data analysis. *IEEE Transactions on Computers*, Ser. C, **23**, p 881-889.
- Fukunaga K. (1972) - *Statistical Pattern Recognition*. Academic Press, Boston.
- Furnas G. W., Deerwester S., Dumais S. T., Landauer T. K., Harshman R. A., Streeter L. A., Lochbaum K. E. (1988) - Information retrieval using a singular value decomposition model of latent semantic structure, *Proceedings of the 14th Int. ACM Conf. on Res. and Dev. In : Information Retrieval*, p 465-480.
- Furnival G. M. (1971) - All possible regressions with less computation, *Technometrics*, **13**, p 403-408.
- Furnival G. M., Wilson R.W. (1974) - Regressions by leaps and bounds, *Technometrics*, **16**, p 499-511.
- Gabriel K.R. (1969) - Simultaneous test procedures: some theory of multiple comparisons. *Ann. Math. Statist.*, **40**, 1, p 224-250.
- Gabriel K.R. (1971) - The biplot graphic display of matrices with application to principal component analysis. *Biometrika*, **58**, 3, p 453-467.
- Gallego F. J. (1982) - Codage flou en analyse des correspondances, *Les Cahiers de l'Analyse des Données*, **7**, n° 4, p 413-430.
- Gallinari P., Thiria S., Fogelman-Soulié F. (1988) - Multilayer perceptrons and data analysis, *International Conference on neural Networks, IEEE*, **1**, p 391-399.
- Garnett J.-C. (1919) - General ability, cleverness and purpose. *British J. of Psych.*, **9**, p 345-366.
- Geary R.C. (1954) - The contiguity ratio and statistical mapping. *The Incorporated Statistician*, **5**, 3, p 115-145.
- Geisser S. (1975) - The Predictive sample reuse method with applications. *J. of Amer. Statist. Assoc.* **70**, p 320-328.
- Gifi A. (1981) - *Non Linear Multivariate Analysis*, Department of Data theory, University of Leiden.
- Gifi A. (1990) - *Non Linear Multivariate Analysis*, J. Wiley, Chichester.

- Gilula, Z. (1986) - Grouping and association in contingency tables: an exploratory canonical correlation approach, *J. of Amer. Statist. Assoc.*, **81**, p 773-779.
- Gilula Z., Ritov Y. (1990) - Inferential ordinal correspondence analysis : motivation, derivation and limitations. *Inter. Statist. Review*, **58**, p 99-108.
- Girshick M.A. (1939) - On the sampling theory of roots of determinantal equations. *Ann. Math. Statist.*, **1**, 10, p 203-224.
- Gnanadesikan R. (1989) - Discriminant analysis and clustering, panel of experts. *Statistical Science*, 1989, 4, n°1, p 34-69.
- Gnanadesikan R., Kettenring J.R., Landwehr J.M. (1982) - Projection plots for displaying clusters, In : *Statistics and Probability, Essays in Honor of C.R. Rao*, G. Kallianpur, P.R. Krishnaiah, J.K.Ghosh, eds, North-Holland.
- Goldstein M., Dillon W. R. (1978) - *Discrete Discriminant Analysis*, J. Wiley, Chichester.
- Good P. (1994) - *Permutation Test - A practical Guide to Resampling Method for Testing Hypotheses*. Springer Verlag, New York.
- Goodman L.A. (1970) - The multivariate analysis of qualitative data: interaction among multiple classifications. *J. of Amer. Statist. Assoc.*, **65**, p 226-256.
- Goodman L.A. (1986) - Some useful extensions of the usual correspondence analysis approach and the usual log-linear approach in the analysis of contingency tables, *International Statist. Review*, **54**, p 243-270.
- Goodman L.A. (1991) - Measures, models, and graphical displays in the analysis of cross-classified data (with Discussion), *J. of Amer. Statist. Assoc.*, **86**, 416, p 1085-1138.
- Goodman L.A., Kruskal W.H. (1954) - Measures of association for cross classification. *J. of Amer. Statist. Assoc.*, **49**, p 732-764.
- Gordon A. D. (1979) - On the assessment and comparison of classification. In : *Analyse des données et informatique. Cours de la C.E.E.*, Tomassone R. (ed.), INRIA, Rocquencourt. p 149-160.
- Gordon A. D. (1981) - *Classification : Methods for the Exploratory Analysis of Multivariate Data*. Chapman and Hall, London.
- Gordon A.D. (1987) - A review of hierarchical classification, *J.R.Statist.Soc.*, A, **150**, Part2, p 119-137.
- Gordon A.D., Finden C.R. (1985) - Classification of spatially located data. *Comp. Statist. Quarterly*, **2**, p 315-328.
- Gourlay A.R., Watson G.A. (1973) - *Computational Methods for Matrix Eigen Problems*. J. Wiley, New York.
- Govaert G. (1977) - Algorithmes de classification d'un tableau de contingence. In : *Premières Journées Internationales Analyse des Données et Informatique (Versailles 1977)* INRIA, p 487-500.
- Govaert, G. (1984) - Classification simultanée de tableaux binaires.- In : *Data Analysis and Informatics*, 4, E. Diday et al., Eds, North Holland, p 223-236.
- Gower J. C. (1966) - Some distance properties of latent and vector methods used in multivariate analysis. *Biometrika*, **53**, p 325-328.
- Gower J. C. (1975) - Generalized Procrustes Analysis. *Psychometrika*, **40**, (1), p 33-51.
- Gower J. C. (1984) - Procrustes analysis. In : *Handbook of Applicable Mathematics*, 6, Lloyd E.H. (ed.), J. Wiley, Chichester, p 397-405.
- Gower J. C., Banfield C. F. (1975) - Goodness-of-fit criteria in cluster analysis and their empirical distributions. In : *Proceeding of the 8th Intern. Biometric Conf.*, Corsten L. C. A., Postelnicu T., (eds), p 347-361.
- Gower J. C., Harding A. (1988) - Nonlinear biplot. *Biometrika*, **75**, p 445-455.
- Gower J. C., Ross G. (1969) - Minimum spanning trees and single linkage cluster analysis. *Appl. Statistics*, **18**, p 54-64.

- Green P. J. (1981) - Peeling bivariate data. In: *Interpreting multivariate data*, Barnett V. (ed.), J. Wiley, Chichester, p 3-20.
- Greenacre M. (1984) - *Theory and Applications of Correspondence Analysis*. Academic press, London.
- Greenacre M. (1988) - Clustering the rows and columns of a contingency table, *J. of Classification*, 5, p 39-51.
- Greenacre M. (1993) - *Correspondence Analysis in Practice*. Academic Press, London.
- Greenacre M., Blasius J. (Eds) (1994) - *Correspondence Analysis in the Social Sciences*. Academic Press, London.
- Grelet Y. (1993) - Préparation des tableaux pour l'analyse des données : le codage des variables. In : *Traitement statistique des enquêtes*, Grangé D., Lebart L. (eds), Dunod, Paris.
- Grizzle J. E., Starner C. F., Koch G. G. (1969) - Analysis of categorical data by linear models. *Biometrics*, 25, p 489-504.
- Grosbras, J.-M. (1986) - *Méthodes statistiques des sondages*. Economica, Paris.
- Guéguen A., Nakache J.-P. (1988) - Méthode de discrimination basée sur la construction d'un arbre de décision binaire. *Revue de Statist. Appl.*, 36, (1), p 19-38.
- Guttman L. (1941) - The quantification of a class of attributes: a theory and method of a scale construction. In : *The prediction of personal adjustment* (Horst P., ed.) p 251-264, SSCR New York.
- Guttman L. (1954) - Some necessary conditions for common factor analysis. *Psychometrika*, 19, p 149-161.
- Haberman S. J. (1974) - *The Analysis of Frequency Data*. University of Chicago University Press, Chicago.
- Hand D. J. (1981) - *Discrimination and Classification*. J. Wiley, New York.
- Hand D. J. (1982) - *Kernel Discriminant Analysis*. J. Wiley, New York.
- Hand D., J. (1986) - Recent advances in error-rate estimation. *Pattern Recogn. lett.*, 4, p 335-346.
- Hand D. J. (1987) - A shrunken leaving-one-out estimator of error rate. *Comput. Math. Applic.*, 14, (3), p 161-167.
- Hand D. J. (1992) - Microdata, macrodata, metadata. In : *Computational Statistics*, Dodge Y., Whittaker J. (Eds), Physica Verlag, Heidelberg, 2, p 325-340.
- Hardy A. (1994) - An examination of procedures for determining the number of clusters in a data set. In : *New Approaches in Classification and Data Analysis*, Diday E. et al. (eds), Springer Verlag, Berlin, p 178-185.
- Harman H.H. (1967) - *Modern Factor Analysis* (2nd ed.). Chicago University Press, Chicago.
- Harshman R. A. (1970) - Foundation of the PARAFAC procedure : Models and conditions for an explanatory multi-modal factor analysis. *UCLA working paper in Phonetics*, 16, UCLA, Los Angeles.
- Harter H.L. (1974-1975) - The method of least squares and some alternatives. *Internat. Statist. Review*, Part 1 and 2: 42, p 147-174, p 235-264; Part 3 to 5: 43, p 1-44, p 125-190, p 269-278.
- Hartigan J. A. (1972) - Direct clustering of a data matrix, *J. of Amer. Statist. Assoc.*, 67, p 123-129.
- Hartigan J. A. (1975) *Clustering Algorithms*. J. Wiley, New York.
- Hartigan J. A. (1985) - Statistical theory in clustering. *J. of Classification*, 2, 63-76.
- Harvatopoulos Y., Livian Y.-F., Sarnin P. (1989) - *L'art de l'enquête*, Eyrolles, Paris.
- Hayashi C. (1956) - Theory and examples of quantification. (II) *Proc. of the Institute of Statist. Math.* 4 (2), p 19-30.

- Hayashi C., Hayashi F. (1982) - A new algorithm to solve PARAFAC model. *Behaviormetrika*, **14**, p 27-48.
- Heiser W. J. (1986) - Undesired nonlinearities in nonlinear multivariate analysis. In : *Data Analysis and Informatics IV*, Diday E. et al. (eds), North Holland, Amsterdam, p 455-469.
- Hertz J., Krogh A., Palmer R.G. (1991) - *Introduction to the Theory of Neural Computation*. Addison- Wesley, Reading, (Mass.).
- Highleyman W.H. (1962) - The design and analysis of pattern recognition experiments. *Bell Syst. Tech. Journal.*, **41**, p 723-744.
- Hill M.O. (1974) - Correspondence analysis: a neglected multivariate method. *Appl. Statist.* **3**, p 340-354.
- Hirschfeld H.D. (1935) - A Connection between correlation and contingency. *Proc. Camb. Phil. Soc.* **31**, p 520-524.
- Holmes S. (1985) - *Outils Informatiques pour l'évaluation de la pertinence d'un résultat en analyse des données*. Thèse USTL, Montpellier.
- Holmes S. (1989) - Using the bootstrap and the RV coefficient in the multivariate context. in : *Data Analysis, Learning Symbolic and Numeric Knowledge*, E. Diday (ed.), Nova Science, New York, p 119-132.
- Hornik K. (1994) - Neural networks : more than "statistics for amateurs". In : *New approaches in Classification and Data Analysis*. Diday E; et al. (eds), Springer Verlag, Berlin.
- Horst P. (1961) - Relation among m sets of measures. *Psychometrika*, **26**, p 129-149.
- Horst P. (1965) - *Factor Analysis of Data Matrices*. Holt, Rinehart, Winston, New York.
- Hosmer D. W., Lemeshow S. (1989) - *Applied Logistic Regression*, J. Wiley, New York.
- Hotelling H. (1933) - Analysis of a complex of statistical variables into principal components. *J. Educ. Psy.* **24**, p 417-441, p 498-520.
- Hotelling H. (1936) - Relation between two sets of variables. *Biometrika*, **28**, p 129-149.
- Householder A.S. (1953) - *Principles of Numerical Analysis*. Mc Graw-Hill, New York.
- Hsu P. L. (1939) - On the distribution of the roots of certain determinantal equations. *Ann. Eugen.* **9**, p 250-258.
- Huber P.J. (1981) - *Robust Statistics*. J. Wiley, New York.
- Huber P.J. (1987) - The place of the  $L_1$ -Norm in robust estimation, in: *Statistical Data Analysis Based on the  $L_1$ -Norm and Related Methods*, Dodge Y. ed., North Holland, Amsterdam, p 23-34.
- Hudon G. (1990) - Une comparaison des résultats de modèles log-linéaires et de généralisations de l'analyse des correspondances. *Revue de Statist. Appl.*, **38**, (2), p 43-53.
- Hurley J. R. , Cattell R. B. (1962) - The Procrustes program : Producing direct rotation to test an hypothesized factor structure. *Behavioural Science*, **7**, p 258-262.
- Jain A. K., Moreau J. V. (1987) Bootstrap technique in cluster analysis. *Pattern Recognition*, **20**, p 547-568.
- Jambu M. (1991) - *Exploration statistique et informatique des données*. Dunod, Paris.
- Jambu M., Lebeaux M.O. (1978) - *Classification automatique pour l'analyse des données*. Dunod, Paris.
- Jeffreys H. (1946) - An Invariant form for the prior probability in estimation problems. *Proc. Roy. Soc. (A)*, **186**, p 453-461.
- Johnson S. C. (1967) - Hierarchical clustering schemes. *Psychometrika*, **32**, p 241-254.
- Jolliffe I. (1986) - *Principal Component Analysis*. Springer-Verlag, New York.



- Jones M.C., and Sibson R. (1987) - What is projection pursuit (with discussion). *J. of Royal Statist. Society, A*, **150**, p 1-36.
- Jordan C. (1874) - Mémoire sur les formes bilinéaires. *J. Math. Pures et Appliquées*, **19**, p 35-54.
- Joreskog K. (1963) - *Statistical Estimation in Factor Analysis : a New Technique and its Foundation*. Almqvist & Wiksell, Uppsala.
- Joreskog K., Sörbom D. (1979) - *Advances in Factor Analysis and Structural Equation Models*. Abt, Cambridge (MA).
- Jousselin B. (1972) - Les choix de consommation et les budgets des ménages. *Consommation*, Dunod, **1**, p 41-72.
- Kaiser H. F. (1961) - A note on Guttman's lower bound for the number of common factors. *Brit. J. Statist. Psychol.*, **14**, p 1-2.
- Kato T. (1966) - *Perturbation Theory for Linear Operators*. Springer, New York.
- Kaufman L., Rousseeuw P. J. (1986) - Clustering large data sets (with discussion). *Pattern recognition in practice II* (E.S. Gelsema and L.N. Kanal, eds), North-Holland, Amsterdam, p 425-437.
- Kaufman L., Rousseeuw P. J. (1990) - *Finding Groups in Data*. J. Wiley, New York.
- Kazmierczak J.-B. (1985) - Analyse logarithmique : deux exemples d'application. *Revue de Statist. Appl.*, **33**, (1), p 13-24.
- Kendall M. G. (1962) - *Rank Correlation Methods*. Griffin, London.
- Kendall M. G. (1966) - Discrimination and classification. In : *Proc. Symp. Mult. Analysis*. Dayton, Ohio, (Krishnaiah P. R. (ed.), Academic Press, New York, p 165-185.
- Kendall M.G., Stuart A. (1961) - *The Advanced Theory of Statistics*. Charles Griffin, London.
- Kettenring R. J. (1971) - Canonical analysis of several sets of variables. *Biometrika*, **58**, (3), p 433-450.
- Kharchaf I., Rousseau R. (1988, 1989) Reconnaissance de la structure de blocs d'un tableau de correspondance par la classification ascendante hiérarchique: parties 1 et 2, *Les Cahiers de l'Analyse des Données*, **13**, p 439-443; et : **14**, p 257-266.
- Kiers H. A. L. (1989) - *Three-way Methods for the Analysis of Quantitative and Qualitative Two-way Data*. DSWO Press, Leiden.
- Kohonen T. (1989) - *Self-Organization and Associative Memory*. Springer Verlag, Berlin.
- Krishnaiah P.R., Chang T. C. (1971) - On the exact distribution of the extreme roots of the Wishart and MANOVA matrix. *J. of Multivariate Anal.*, **1**, (1), p 108-116.
- Krishnaiah P.R., Kanal L. (Eds) (1982) - *Handbook of Statistics (2)*. North Holland, Amsterdam.
- Kroonenberg P. (1983) - *Three-Mode Principal Component Analysis*. DSWO Press, Leiden.
- Kroonenberg P. M., de Leeuw J. (1980) - Principal component analysis of three-mode data by means of alternating least-square algorithms. *Psychometrika*, **45**, p 69-97.
- Kruskal J. B. (1956) - On the shortest spanning subtree of a graph and the traveling salesman problem. *Proc. Amer. Math. Soc.*, **7**, p 48-50.
- Kruskal J. B., Wish M. (1978) - *Multidimensional Scaling*. Sage University Paper, **11**, Sage, Beverly Hills.
- Krzanowski W. J. (1984) - Sensitivity of principal components. *J. Royal Statist. Soc.*, **46**,(3), p 558-563.
- Krzanowski W. J. (1987) - Cross-validation choice in principal component analysis. *Biometrics*, **43**, p 575-584.
- Kshirsagar A.M. (1972) - *Multivariate Analysis*. Marcel Dekker, New York.

- Kullback S. (1959) - *Information Theory and Statistics*. J. Wiley, New York.
- Kullback S., Leibler R.A. (1951) - Information and sufficiency. *Ann. Math. Statist.*, **22**, p 79-86.
- L'Hermier des Plantes H. (1976) - *STATIS : Structuration de tableaux à trois indices de statistique*. Thèse (3c), USTL, Montpellier.
- Lachenbruch P.A., Goldstein M. (1979) - Discriminant Analysis. *Biometrics*, **35**, p 68-85.
- Lachenbruch P.A., Mickey M.R. (1968) - Estimation of error rate in discriminant analysis. *Technometrics*, **10**, p 1-11.
- Lafosse R. (1985) - *Analyses Procrustéennes de deux Tableaux*. Thèse, Univ. Paul Sabatier, Toulouse.
- Lancaster H. O. (1963) - Canonical correlation and partition of  $\chi^2$ . *Quart. J. Math.*, **14**, p 220-224.
- Lancaster H. O. (1969) - *The Chi-squared Distribution*. J. Wiley, New York.
- Lance G. N., Williams W. T. (1967) - A general theory of classification sorting strategies. *Computer J.*, **9**, p 373-380.
- Laplace P.S. (1793) - Sur quelques points du système du monde. *Mémoires de l'Académie Royale des Sciences de Paris*, p 1-87; Réédition: *Oeuvres*, (1895), **11**, Gauthier-Villars, Paris, p 477-558.
- Lauro N. C, D'Ambra L. (1984) - L'Analyse non-symétrique des Correspondances. In : *Data Analysis and Informatics*, III, Diday et al. Ed., North-Holland, p 433-446.
- Lauro N. C., Decarli A. (1982) - Correspondence analysis and log-linear models. In : multiway contingency tables study. *Metron*, **1-2**, p 213-234.
- Lavit C. (1988) - *Analyse Conjointe de Tableaux Quantitatifs*. Masson, Paris.
- Lawley D. N. (1956) - Tests of significance for the latent roots of covariance and correlation matrices. *Biometrika*, **43**, p 128-136.
- Lawley D. N., Maxwell A. E. (1963) - *Factor Analysis as a Statistical Method*. Methuen, London.
- Le Calvé G. (1987) -  $L_1$  -embeddings of a data structure (I,D). in: *Statistical Data Analysis Based on the  $L_1$ -Norm and Related Methods*, Dodge Y. ed., North Holland, Amsterdam, p 195-202.
- Le Foll Y. (1982) - Pondération des distances en analyse factorielle. *Statist. et Anal. des Données*, **7**, p 13-31.
- Le Foll Y., Burtschy B. (1983) - Représentations optimales des matrices imports-exports. *Revue de Statist. Appl.*, **31**, (3), p 57-72.
- Lebart L. (1969 a) - L'Analyse statistique de la contiguïté. *Publications de l'ISUP*, XVIII- p 81 - 112.
- Lebart L. (1969 b) - Introduction à l'analyse des données : Analyse des correspondances et validité des résultats. *Consommation*, Dunod. **4**, p 65-87.
- Lebart L. (1974) - On the Benzécri's method for finding eigenvectors by stochastic approximation. *Proceedings in Comp. Statist.*, COMPSTAT, Physica verlag, Vienna, p 202-211.
- Lebart L. (1975 a) - L'orientation du dépouillement de certaines enquêtes par l'analyse des correspondances multiples. *Consommation*, **2**, p 73-96. Dunod.
- Lebart L. (1975 b) - *Validité des résultats en analyse des données*. Rapport Credoc-Cordes. Credoc, Paris.
- Lebart L. (1976) - The significance of eigenvalues issued from correspondence analysis. *Proceedings in Comp. Statist.*, COMPSTAT, Physica verlag, Vienna, p 38-45.

- Lebart L. (1982) - Exploratory analysis of large sparse matrices, with application to textual data. *COMPSTAT*, Physica Verlag, Vienna, p 67-76.
- Lebart L. (1986) - Qui pense quoi ? Evolution et structure des opinions en France de 1978 à 1984. *Consommation Revue de Socio-Economie*, Dunod, 4, p 3-22.
- Lebart L. (1987 a) - Some recent advances in data analysis practice. In : *New Perspective in Theoretical and Applied Statistics*. M.L. Puri and al., Eds. J. Wiley, New York.
- Lebart L. (1987 b) - Conditions de vie et aspirations des Français, Evolution et structure des opinions de 1978 à 1984. *Futuribles*, 1, p 25-56.
- Lebart L. (1988) - Contribution of classification to the processing of longitudinal socio-economic surveys. In : *Classification and Related Methods of Data Analysis*, H. Bock Ed., North Holland, p 113-120.
- Lebart L. (1992) - Discrimination through the regularized nearest cluster method. *COMPSTAT; Proceedings of the 10th Symposium on Computational Statistics*, Physica Verlag, Vienna, p 103-118.
- Lebart L., Fénelon J.P. (1971) - *Statistique et informatique appliquées*. Dunod, Paris.
- Lebart L., Houzel Y. (1980) - Le système d'enquête sur les aspirations des Français. *Consommation Revue de Socio-Economie*, Dunod, 1, p 3-25.
- Lebart L., Mirkin B. (1993) - Correspondence analysis and classification. In : *Multivariate Analysis: Future Directions 2*, C. M. Cuadras and C.M.Rao, Eds., North Holland, Amsterdam, p 341-357.
- Lebart L., Morineau A., Fénelon J.P. (1981) - *Traitement des Données Statistiques*. Dunod, Paris.
- Lebart L., Morineau A., Lambert T., Pleuvret P. (1991) - *SPAD.N version 2 Système Portable pour l'Analyse des Données*. CISIA, 1 avenue Herbillon 94160, Saint-Mandé.
- Lebart L., Morineau A., Tabard N. (1977) - *Techniques de la description statistique*. Dunod, Paris.
- Lebart L., Morineau A., Warwick K. (1984) - *Multivariate Descriptive Statistical Analysis*. J. Wiley, New York.
- Lebart L., Salem A. (1994) - *Statistique textuelle*. Dunod, Paris.
- Lebart L., Tabard N. (1973) - *Recherches sur la description automatique des données socio-economiques*. Rapport CORDES-CREDOC, C.R n°13/1971.
- Lebreton J.-D., Chessel D., Prodon R., Yoccoz N. (1988) - L'analyse des relations espèces-milieu par analyse canonique des correspondances. *Acta Oecologica, Oecol. Gener.*, 9, (1), p 53-67.
- Leclerc A. (1975) - L'analyse des correspondances sur juxtaposition de tableaux de contingence. *Revue Statist. Appl.*, 23, p 5-16.
- Leclerc A. (1976) - Une étude de la relation entre une variable qualitative et un groupe de variables qualitatives. *Int. Statist. Review*, 44, p 221-248.
- Leclerc A., Chevalier A., Luce D., Blanc M. (1985) - Analyse des correspondances et modèle logistique : possibilités et intérêt d'approches complémentaires. *Revue Statist. Appl.*, 33, p 25-38.
- Lelu A. (1991) - From data analysis to neural networks : new prospects for efficient browsing through databases. *Journal of Information Science*, 17, p 1-12.
- Lelu A. (1994) - Clusters and factors : neural algorithm for a novel representation of highly multidimensional data sets. In : *New Approaches in Classification and Data Analysis*, Diday et al. (eds), Springer Verlag, Berlin, p 241-248.
- Lerman I., C. (1970) - *Les Bases de la Classification Automatique*. Gauthier-Villars, Paris.
- Lerman I., C. (1981) - *Classification et analyse ordinale des données*. Dunod, Paris.

- Li G., and Chen Z. (1985) - Projection-pursuit approach to robust dispersion matrices and principal components: primary theory and monte carlo. *J. of Amer. Statist. Assoc.* 80, p 759-766.
- Ling R. F. (1973) - A probability theory of cluster analysis. *J. Amer. Statist. Assoc.*, 68, p 159-164.
- MacQueen J. B. (1967) - Some methods for classification and analysis of multivariate observations. *Proc. Symp. Math. Statist. and Probability (5th)*, Berkeley, 1, p 281-297, Univ. of Calif. Press, Berkeley.
- Mahalanobis P.C. (1936) - On the generalized distance in statistics. *Proc. Nat. Inst. Sci. India*, 12, p 49-55.
- Malinvaud E. (1964) - *Méthodes statistiques de l'économétrie*. (3ème éd. 1978), Dunod, Paris.
- Malinvaud E. (1987) - Data analysis in applied socio-economic statistics with special consideration of correspondence analysis. *Marketing Science Conference Proceedings*, HEC-ISA, Jouy en Josas.
- Mallows C.L. (1973) - Some comments on  $C_p$ . *Technometrics*, 15, p 661-675.
- Mallows C.L., Tukey J.W. (1982) - An overview of technique of data analysis emphasizing its exploratory aspects. In : *Some Recent Advances in Statistics* (J. Tiago de Oliveira Ed.) Academic Press, p 11-172.
- Marcotorchino F. (1987) - Block seriation problems: a unified approach, *Applied Stochastic Models and Data Analysis*, 3, p 73-93.
- Marsaglia G., Bray T.A. (1964) - A convenient method for generating normal variables. *SIAM Rev.* 6, p 260-264.
- Martin R.S., Reinsch C., Wilkinson J.H. (1968) - Householder's tridiagonalisation of a symmetric matrix. *Num. Math.* 11, p 181-195.
- Martin R.S., Wilkinson J.H. (1968) - Implicit QL algorithm. *Num. Math.* 12, p 377-383.
- Masson M. (1974) - Analyse non linéaire de données. *C.R. Acad. Sc.*, 278 (11 mars).
- Masson M. (1980) - *Méthodologie générale du traitement statistique de l'information de masse*. Nathan, Paris.
- Matheron G. (1963) - Principles of geostatistics. *Economic Geology*, 58, p 1246-1266.
- Matheron G. (1965) - *Les variables régionalisées et leur estimation*. Masson, Paris.
- Matusita K. (1955) - Decision rules based on the distance, for problems of fit, two samples, and estimation. *Ann. of Math. Statist.* 26, 4, p 631-640.
- Matusita K., Ohsumi N. (1980) - A criterion for choosing the number of clusters in cluster analysis. In : *Recent Development in Statistical Inference and Data Analysis*, Matusita K. (ed.) North-Holland, Amsterdam, p 203-213.
- McCullagh P., Nelder J.A. (1989) - *Generalized Linear Models*. Chapman and Hall, London.
- McLachlan G.J. (1992) - *Discriminant Analysis and Statistical Pattern Recognition*. J. Wiley, New York.
- McQuitty L.L. (1966) - Single and multiple classification by reciprocal pairs and rank order type. *Educational Psychology Measurements*. 26, p 253-265.
- Mehta C. R., Patel N., R. (1991) - *Statistique non-paramétrique exacte, Introduction à StatXact*. CISIA, Saint Mandé.
- Mehta M.L. (1960) - On the statistical properties of the level spacing in nuclear spectra. *Nucl. Phys.* 18, p 395-419.
- Mehta M.L. (1967) *Random Matrices and the Statistical Theory of Energy Levels*. Academic Press, New York.
- Meot A., Chessel D., Sabatier R. (1993) - Opérateur de voisinage et analyse des données spatio-temporelles. In *Biométrie et environnement*, Lebreton J.-D., Asselain B., (eds), Masson, Paris, p 45-71.

- Meulman J. (1982) - *Homogeneity Analysis of Incomplete Data*. DSWO Press, Leiden.
- Meyer R. (1994) - An eigenvector algorithm to fit  $L_p$ -distance matrices. In : *New Approches in Classification and Data Analysis*, Diday E. et al. (eds), Springer Verlag, Berlin, p 502-509.
- Michelat G. et Simon M. (1985) - Les sans-réponses aux questions politiques, *Revue Pouvoirs*, **33**, PUF, Paris.
- Milgram M. (1993) - *Reconnaissance des formes, méthodes numériques et connexionnistes*. Armand Colin, Paris.
- Miller R. G. (1966) - *Simultaneous Statistical Inference*. Mac Graw Hill, New York.
- Müller R. G. (1974) - The Jakknife-a review. *Biometrika*, **61**, p 1-15.
- Milligan G. W., Cooper M. C. (1985) - An examination of procedures for determining the number of cluster in a data set. *Psychometrika*, **50**, p 159-179.
- Mirkin, B.G. (1990) - A sequential fitting procedure for linear data analysis models, *J. of Classification*, **7**, p 167-195.
- Mollière J.-L. (1986) - What's the real number of clusters? In : *Classification as Tool of Research*, Gaul W., Schader M. (eds), North-Holland, Amsterdam, p 311-320.
- Mollière J.-L. (1989) - Stratégie de classification pour de grands ensembles de données. *La Revue de Modulad (INRIA)*, **3**, p 31-69.
- Mom A. (1988) - *Méthodologie statistique de la classification des réseaux de transport*. Thèse, U.S.T.L., Montpellier.
- Mood A.M. (1951) - On the distribution of the characteristic roots of normal second moment matrices. *Ann. Math. Statist.* **22**, p 266-273.
- Moran P. A. P. (1948) - The interpretation of statistical maps. *J. Royal Statist. Soc., B*, **10**, p 243-251.
- Moran P.A.P. (1954) - Notes on continuous stochastic phenomena, *Biometrika*, **37**, p 17-23.
- Moreau J. (1992) - *Analyse de données structurées par des graphes. Cas de l'analyse des correspondance*. Thèse, E.P.F.L., Lausanne.
- Morgan J. M., Messenger R. C. (1973) - *THAID : a sequential search program for the analysis of nominal scale dependent variables*. Institute for Social Research, University of Michigan, Ann Arbor.
- Morgenthaler S., Tukey J.W. (1989) - The next future of data analysis. In *Data Analysis, Learning Numeric and Symbolic Knowledge*, 1989, Diday ed., Novascience, New York, p 1-12.
- Morineau A. (1983) - Etude de stabilité en analyse en composantes principales. *Bull. Techn. du Centre de Statist. et d'Infor. Appl.*, **1**, p 9-12.
- Morineau A. (1984) - Note sur la caractérisation statistique d'une classe et les valeurs-tests, *Bull. Techn. du Centre de Statist. et d'Infor. Appl.*, **2**, p 20-27.
- Morineau A. (1992) - L'Analyse de données et les tests de cohérence dans les données d'enquête. In : *La Qualité de l'Information dans les Enquêtes*, ASU (ed.), Dunod, Paris.
- Morineau A., Lebart L. (1986) - Specific clustering algorithms for large data sets and implementation in SPAD Software. In : *Classification as a tool of research*, Gaul W., Schader M., Eds, North Holland, Amsterdam, p 321-330
- Morineau A., Nakache J.-P., Krzyzanowski C. (1995) - *Le modèle log-linéaire et ses applications. (La procédure Logli de SPAD.N)*, CISIA, Paris.
- Morineau A., Sammartino A.-E., Gettler-Summa M., Pardoux C. (1994) - Analyses des données et modélisation des séries temporelles. *Revue Statist. Appl.*, **42**, (4), p 61-81.
- Morrison. D.F. (1967) - *Multivariate Statistical Methods* (2nd edition : 1976). Mac Graw Hill, New York.

- Mosteller F., Tukey J. W. (1977) - *Data Analysis and Regression*. Addison Wesley, Reading, (Mass).
- Muirhead R. J. (1982) - *Aspects of Multivariate Statistical Theory*. J. Wiley, New York.
- Mulaik S. A. (1972) - *The Foundation of Factor Analysis*. McGraw Hill, New York.
- Murtagh F. (1985) - *Multidimensional Clustering Algorithms*. COMPSTAT Lectures 4, Physica Verlag, Vienna.
- Nakache J.P. (1973) - Influence du codage des données en analyse factorielle des correspondances. Etude d'un exemple pratique médical. *Revue Statist. Appl.*, **21**, (2).
- Nakache J.-P., Lorente P., Benzécri J.-P., Chastang J.-F. (1977) - Aspects Pronostics et thérapeutiques de l'infarctus myocardique aigu. *Les Cahiers de l'Analyse des Données*, **2**, p 415-434.
- Nakhlé F. (1976) - Sur l'analyse d'un tableau de notes dédoublées. *Les Cahiers de l'Analyse des Données*, **1**, p 243-257.
- Neave H.R. (1973) - On using Box-Muller transformation with multiplicative congruential pseudo-random number generators. *Appl. Statist.*, **22**, p 92-97.
- Nelder J.A., Wedderburn R.W.M. (1972) - Generalized linear models. *J. R. Statist. Soc.*, **A**, **135**, p 370-384.
- Newman T.G., Odell P.L. (1971) - *The Generation of Random Variates*. GRIFFIN's Statistical Methods and Courses, n°29, Griffin.
- Nijenhuis A., Wilf H.S. (1975) - *Combinatorial Algorithms*. Academic Press, New York.
- Nishisato S. (1980) - *Analysis of Categorical Data. Dual Scaling and its Application*. Univ. of Toronto Press.
- O'Neill M. E. (1978) - Distributional expansion from canonical correlation from contingency tables. *J. Roy. Statist. Soc.*, **B**, **40**, p 301-312.
- O'Neill M. E. (1981) - A note on the canonical correlation from contingency tables. *Austr. J. Statist.*, **23**, p 58-66.
- Ohsumi N. (1988) - Role of computer graphics in interpretation of clustering results. In: *Recent Developments in Clustering and Data Analysis*, Diday E. et al. (eds), Academic Press, Boston.
- Oja E. (1982) - A simplified neuron model as a principal components analyzer. *J. of Math. Biology*, **15**, p 267-273.
- Oja E. (1992) - Principal components, minor components, and linear neural networks. *Neural Networks*, **5**, p 927-935.
- Oja E., Karhunen J. (1981) On stochastic approximation of the eigenvectors and eigenvalues of the expectation of a random matrix. Report of the Helsinki University of Technology (Dept of Technical Physics). Otaniemi, Finland.
- Pagès J.-P., Escoufier Y., Cazes P. (1976) - Opérateurs et analyse de tableaux à plus de deux dimensions. *Cahiers du BURO*, ISUP, Paris, p 61-89
- Palm R., lemma A. F. (1995) - Quelques alternatives à la regression classique dans le cas de la colinéarité. *Revue Statist. Appl.*, **43**, (2), p 5-33.
- Parzen E. (1962) - On the estimation of a probability density function and mode. *Ann. of Math. Statist.*, **33**, p 1065-1076.
- Pearson K. (1901) - On lines and planes of closest fit to systems of points in space. *Phil. Mag.*, **2**, n°11, p 559-572.
- Perruchet C. (1983) - Une analyse bibliographique des épreuves de classifiabilité en analyse des données. *Statist. et Anal. des Données*, **8**, p 18-41.
- Pillai K.C.S. (1965) - On the distribution of the largest root of a matrix in multivariate analysis. *Biometrika*, **52**, p 405-414.
- Pillai K.C.S. (1967) - Upper percentage point of the largest root of a matrix in multivariate analysis. *Biometrika*, **54**, p 189-194.

- Pillai K.C.S., Chang T.C. (1970) - An approximation to the c.d.f. of the largest root of a covariance matrix. *Ann. of the Inst. of Statist. Math.*, p 115-124.
- Piron M. (1990) - *Structuration de l'information à plusieurs niveaux et analyse des données*. Thèse, Université Pierre et Marie Curie.
- Piron M. (1992) - *Analyse statistique d'un système d'échelles*. Réseau ADOC, doc. 4, ORSTOM, Bondy.
- Pousse A. (1992) - Résultats asymptotiques. In : *Modèles pour l'analyse des données multidimensionnelles*, Dreesbeke et al., eds, Economica, Paris.
- Prim R.C. (1957) - Shortest connection matrix network and some generalizations. *Bell System Techn. J.*, 36, p 1389-1401.
- Proriol J. (1991) MLP - Programme de réseau de neurone multicouche. *La Revue de MODULAD*, 8, INRIA, p 23-29.
- Quenouille M. (1949) - Approximate tests of correlation in time series. *J. Royal Statist. Soc.*, B, 11, p 18-44.
- Ramsay J.O. (1978) - Confidence region for multidimensional scaling analysis. *Psychometrika*, 43, p 145-160.
- Rao C.R. (1964) - The use and interpretation of principal component analysis in applied research. *Sankhya serie A*, 26, p 329-357.
- Rao C.R. (1973) - *Linear Statistical Inference and its Application*. (1st ed. : 1965) J. Wiley, New York.
- Rao C.R. (1989) - *Statistics and Truth*. International Cooperative Publishing House, Fairland.
- Rasson J.-P., Kubushishi T. (1994) - The gap test : an optimal method for determining the number of natural classes in cluster analysis. In : *New Approaches in Classification and Data Analysis*, Diday E. et al. (eds), Springer Verlag, Berlin, p 186-193.
- Reinert M. (1986) - Classification descendante hiérarchique : un algorithme pour le traitement des tableaux logiques de grandes dimensions. In *Data Analysis and Informatics*, 4, Diday et al. Ed., North-Holland, p 23-28.
- Richardson M., Kuder G. F. (1933) - Making a rating scale that measures. *Personnel Journal.*, 12, p 71-75.
- Ripley B. D. (1981) - *Spatial Statistics*. J. Wiley, New York.
- Ripley B. D. (1983) - Computer generation of random variables : a tutorial. *Inter. Statist. Review*, 51, p 301-319.
- Ripley B. D. (1993) - Statistical aspects of neural networks. In : *Networks and Chaos-Statistical and Probabilistic Aspects*, Barndorff-Nielsen O.E., Jensen J. L., Kendall W. S., (eds), Chapman and Hall, London, p 40-123..
- Ripley B. D. (1994) - Neural networks and related methods of classification. *J. R. Statist. Soc. B*, 56, n°3, p 409-456.
- Ritter H., Martinez T., Schulten K. (1992) - *Neural Computation and Self-Organizing Maps : An Introduction*. Addison Wesley, Reading.
- Robert Ch. (1992) - *L'analyse statistique bayésienne*. Economica, Paris.
- Robert P., Escoufier Y. (1976) - A unifying tool for linear multivariate methods : the Rv coefficient. *Applied Statistics*, 25, (3), p 257-265.
- Romedor J.M. (1973) - *Méthodes et Programmes d'Analyse Discriminante*. Dunod, Paris.
- Rosenblatt M. (1956) - Remarks on some nonparametric estimates of the density function. *Ann. of Math. Statist.*, 27, p 823-835.
- Rouanet H., Le Roux B. (1993) - *Analyse des données Multidimensionnelles*. Dunod, Paris.
- Roux M. (1985) - *Algorithmes de Classification*. Masson, Paris.
- Roux M. (1991) - Basic procedures in hierarchical cluster analysis. *Applied Multivariate Analysis in SAR and Environmental Studies* (J. Devillers and W. Karcher, eds), p 115-135 ECSC, EEC, EAEC, Brussels and Luxembourg.

- Roy S.N. (1939) -  $p$  - Statistics or some generalisations of analysis of variance appropriate to multivariate problems. *Sankhya*, **4**, p 381-396.
- Rumelhart D. E., Hinton G. E., Williams R. J. (1986) - Learning internal representation by back-propagating errors. *Nature*, **323**, p 533-536.
- Sabatier R. (1984) - Quelques généralisations de l'analyse en composantes principales de variables instrumentales. *Statist. et Anal. des Données*, **9**, (3), p 75-103.
- Sabatier R. (1987) - Analyse factorielle de données structurées et métriques. *Statist. et Anal. des Données*, **12**, (3), p 75-96.
- Sabatier R., Lebreton J.-D., Chessel D. (1989) - Principal component analysis with instrumental variables as a tool for modeling composition data. In : *Mutivay Data Analysis*, Coppi R., Bolasco S. (eds), Elsevier, Amsterdam.
- Saporta G. (1975 a) - *Liaisons entre plusieurs ensembles de variables et codages de données qualitatives*. Thèse 3°C., Université Paris VI.
- Saporta G. (1975 b) - Dépendance et codage de deux variables aléatoires. *Revue Statist. Appl.* **23**, p 43-63.
- Saporta G. (1977) - Une méthode et un programme d'analyse discriminante sur variables qualitatives. In : *Premières Journées Int. Analyse des Données et informatiques*, INRIA, Rocquencourt.
- Saporta G. (1990) - *Probabilités, analyse des données et statistiques*. Technip, Paris.
- Saporta G., Hatabian G. (1986) - Régions de confiance en analyse factorielle. In : *Data Analysis and Informatics*, **4**, Diday E. et al. (eds), North-Holland, Amsterdam, p 499-508.
- Sarle W. S. (1983) - *Cubic clustering criterion*. SAS Technical Report. A-108. SAS Institute Limited. Cary, NC.
- Schiffman S. S., Lance Reynolds M., Young F. W. (1981) - *Introduction to Multidimensional Scaling*. Academic Press, New York.
- Schönemann P. H. (1968) - On two-sided orthogonal procrustes problems. *Psychometrika*, **33**, p 19-33.
- Schönemann P. H., Carroll R. M. (1970) - Fitting one matrix to another under choice of a central dilation and a rigid motion. *Psychometrika*, **35**, p 245-255.
- Schriever B.F. (1983) - Scaling of order dependent categorical variables with correspondence analysis. *Inter. Statist. Review*, **51**, p 225-238.
- Searle S.E. (1971) - *Linear Models*. J. Wiley, New York.
- Seber G.A.F. (1977) - *Linear Regression Analysis*, J. Wiley, New York.
- Shepard R. N. (1974) - Representation of structure in similarity data : problems and prospects. *Psychometrika*, **39**, (4), p 373-421.
- Silverman B. W. (1986) - *Density Estimation for Statistics and Data Analysis*. Chapman and Hall, London.
- Sirat J. A. (1991) - A fast neural algorithm for principal component analysis and singular value decomposition. *Internat. J. of Neural Systems*, **2**, p 147-155.
- Sneath P. H. A. (1957) - The Application of computers to taxonomy. *J. General Microbiology*, **17**, p 201-226.
- Snee R.D. (1974) - Graphical displays of two-ways contingency tables. *Amer. Statistician* **28**, p 9-12.
- Sokal R. R., Sneath P. H. A. (1963) - *Principles of Numerical Taxonomy*, Freeman and co., San-Francisco.
- Sonquist J. A. and Morgan J. N. (1964) - *The Detection of Interaction Effects*. Institute for Social Research, University of Michigan, Ann Arbor.
- Spearman C. (1904) - General intelligence, objectively determined and measured. *Amer. Journal of Psychology*, **15**, p 201-293.



- Stauffer D. F., Garton E. O., Steinhorst R. K. (1985) - A comparison of principal component from real and random data. *Ecology*, 66, p 1693-1698.
- Steinberg D.M., Hunter W.G. (1984) - Experimental Design : Review and comments. *Technometrics*, 26(2), p 71-97.
- Stone M. (1974) - Cross-validatory choice and assessment of statistical predictions. *J. R. Statist. Soc. B* 36, p 111-147.
- Stone M. (1977) - An asymptotic equivalence of choice of model by cross-validation and Akaike's criterion. *J. Royal Statist. Soc. B*, 39, p 44-47.
- Sugiyama T. (1966) - On the distribution of the largest latent root and the corresponding latent vector for principal component analysis. *Ann. Math. Statist.* 37, p 995-1001.
- Sylvester J.J. (1889) - *Messenger of Mathematics* (cité par Eckart, Young, 1939). 19, n°42.
- Tabard N. (1972) - Consommation et statut social, analyse multidimensionnelle des budgets familiaux. *Consommation*, 2, p 41-63.
- Tanaka Y. (1984) - Sensitivity analysis in Hayashi's third method of quantification. *Behaviormetrika*, 16, p 31-44.
- Tenenhaus M. (1994) - *Méthodes statistiques en gestion*. Dunod, Paris.
- Tenenhaus M., Leroux Y., Guimart C., Gonzales P. L. (1993) - Modèle linéaire généralisé et analyse des correspondances. *Revue de Statist. Appl.*, 41, (2) p 59-86.
- Tenenhaus M., Young F. W. (1985) - An analysis and synthesis of multiple correspondence analysis, optimal scaling, dual scaling, homogeneity analysis and other methods for quantifying categorical multivariate data. *Psychometrika*, 50, p 91-119.
- Ter Braak C. J. F. (1986) - Canonical Correspondence Analysis. : a new eigenvector technique for multivariate direct gradient analysis. *Ecology*, 67, (5), p 1167-1179.
- Ter Braak C. J. F. (1987) - The analysis of vegetation-environment relationships by canonical correspondence analysis. *Vegetatio*, 69, p 69-77.
- Ter Braak C. J. F. (1988) - Partial Canonical Correspondence Analysis. In : *Classification and Related Methods of Data Analysis*, Bock H. H. (ed.) North Holland, Amsterdam, p 551-558.
- Theil H. (1971) - *Principles of Econometrics*. J. Wiley, New York.
- Thionet P. (1976) - Construction et reconstruction de tableaux statistiques. *Annales de l'INSEE*, 22-23, p 5-28.
- Thom R. (1974) - *Modèles mathématiques de la morphogénèse*. 10/18, Bourgois, Paris.
- Thorndike R.L. (1953) - Who belongs in the family. *Psychometrika*, 18, p 267-276.
- Thurstone L. L. (1947) - *Multiple Factor Analysis*. The University of Chicago Press, Chicago.
- Tomassone R., Danzart M., Daudin J.J., Masson J.P. (1988) - *Discrimination et classement*. Masson, Paris.
- Tomassone R., Dervin C., Masson J.-P. (1993) - *Biométrie, Modélisation de phénomènes biologiques*. Masson, Paris.
- Tomassone R., Lesquoy E., Millier C. (1983) - *La régression : nouveaux regards sur une ancienne méthode statistique*. Masson, Paris.
- Toussaint G.T. (1974) - Bibliography on estimation of misclassification. *IEEE Trans. Inform. Theory*, IT-20, p 472-479.
- Tucker L. R. (1958) - An inter-battery method of factor analysis. *Psychometrika*, 23, (2).
- Tucker L. R. (1964) - The extension of factor analysis to three-dimensional matrices. In : *Contribution to Mathematical Psychology*, Harris C. W. (ed.), Univ. of Wisconsin Press, Madison, p 109-127.

- Tucker L. R. (1966) - Some mathematical notes on three-mode factor analysis. *Psychometrika*, **31**, p 279-311.
- Tukey J. W. (1958) - Bias and confidence in not quite large samples. *Ann. Math. Statist.*, (Abstract), **29**, p 614.
- Tukey J. W. (1977) - *Exploratory Data Analysis*. Addison Wesley, Reading, Mass.
- Van Cutsem B. (ed.) (1994) - *Classification and Dissimilarity Analysis*. Springer-Verlag, New York.
- van Buuren S., and Heiser W.J. (1989) - Clustering N objects into k groups under optimal scaling of variables. *Psychometrika*, **54**, 4, p 699-706.
- van der Heijden, P. G. M. (1987) - *Correspondence Analysis of Longitudinal Categorical Data*. DSWO Press, Leiden.
- van der Heijden P. G. M., de Leeuw J. (1985) - Correspondence analysis used complementary to log-linear analysis. *Psychometrika*, **50**, p 429-447.
- van der Heijden P. G. M., de Falguerolles A., de Leeuw J. (1989) - A combined approach to contingency table analysis with correspondence analysis and log-linear analysis (with discussion). *Applied Statistics*, **38**, p 249-292.
- van Rijckevorsel J. (1987) - *The application of fuzzy coding and horseshoes in multiple correspondances analysis*. DSWO Press, Leiden.
- Wakimoto K., Taguri M. (1978) - Constellation graphical methods for representing multidimensional data. *Ann. of the Inst. of Statist. Math.*, **30**, (1), p 97-104.
- Ward J.H. (1963) - Hierarchical grouping to optimize an objective function. *J. of Amer. Statist. Assoc.*, **58**, p 236-244.
- Watermaux C. M. (1976) - Asymptotic distribution of the sample roots for a non-normal population. *Biometrika*, **63**, p 639-645.
- Weinberg S.L., Carrol J. D., Cohen H.S. (1984) - Confidence region for INDSCAL using the Jackknife and bootstrap techniques. *Psychometrika*, **49**, p 475-491.
- Werbos P. J. (1974) - *Beyond Regression : New Tools for Prediction and Analysis in the Behavioral Sciences*. Ph.D. Thesis, Harvard University.
- Werbos P. J. (1990) - Backpropagation through time : what it does and how to do it. *Proceedings of the IEEE*, **78**, (10), p 1550-1560.
- Wermuth N. (1976) - Analogies between multiplicative models in contingency tables and covariance selection. *Biometrics*, **32**, p 95-108.
- Wermuth N., Cox D. R. (1992) - Graphical models for dependencies and associations. In : *Computational Statistics* (Dodge Y., Whittaker J., eds), **1**, p 235-250, Physica Verlag, Heidelberg.
- Whittaker J. (1990) - *Graphical Models in Applied Multivariate Statistics*. J. Wiley, Chichester.
- Wilkinson J. H. (1965) - *The Algebraic Eigenvalue Problem*. Clarendon Press, Oxford.
- Wilkinson J. H., Reinsch C. (1971) - *Handbook for Automatic Computation*. **2**, Linear Algebra, Springer-Verlag.
- Williams W. T. and Lambert J. M. (1959) - Multivariate methods in plant ecology. (I) Association analysis in plant communities. *J. Ecology*, **47**, p 83-101.
- Williams W. T., Lance G. N. (1965) - Logic of computer based intrinsic classifications. *Nature*, **207**, p 159-161.
- Wishart D. (1969) - Mode analysis : a generalization of nearest neighbour which reduces chaining effects. *Numerical Taxonomy* (A.J. Cole ed.) p 282-311, Academic Press, London, .
- Wishart J. (1928) - The generalized product-moment distribution in samples from a normal multivariate population. *Biometrika*, **20A**, p 32-43.

- Wold S. (1976) - Pattern recognition by means of disjoint principal component models. *Pattern Recognition*, 8, p 127-139.
- Wold S. (1978) - Cross-validatory estimation of the number of components in factor and principal component models. *Technometrics*, 20, p 397-405.
- Wong M.A. (1982) - A hybrid clustering method for identifying high density clusters. *J. of Amer. Statist. Assoc.*, 77, p 841-847.
- Wong M. A. (1985) - A bootstrap testing procedure for investigating the number of subpopulations. *J. Statist. Comput. and Simul.*, 22, p 99-112.
- Worsley K. J. (1987) - Un exemple d'identification d'un modèle log-linéaire grâce à une analyse des correspondances. *Revue de Statist. Appl.* 35, p 13-20.
- Yenyukov I. S. (1988) - cf. Enyukov I. S. (1988).
- Young G. A. (1994) - Bootstrap : More than a stab in the dark. *Statistical Science*. 9, p 382-418.

## Index des auteurs

### A

*Agrawala A. K.* 251  
*Agresti A.* 3, 284, 287  
*Aitichison J.* 268  
*Aitken C. G. G.* 268  
*Aitkin M. A.* 290  
*Akaike H.* 245, 289  
*Aluja Banet T.* 329  
*Amari S.* 283  
*Anastassakos I.* 374  
*Anderberg M.R.* 145, 151  
*Anderson J. A.* 284, 292, 293,  
*Anderson T. W.* 32, 250, 251, 266, 360,  
375, 376  
*Andrews D.F.* 7  
*Arabie P.* 9, 29  
*Art D.* 333  
*Atkinson A. C.* 223, 245

### B

*Babeau A.* 199  
*Bailey R. A.* 237  
*Balbi S.* 326  
*Baldi P.* 281  
*Ball G. H.* 148, 152  
*Ballif J.-F.* 347, 353  
*Banfield C. F.* 403  
*Bardos M.* 277  
*Barnett B.* 392  
*Bartlett M. S.* 365, 375  
*Beltrami E.* 17  
*Benali H.* 336  
*Benasseni J.* 380  
*Brent R. P.* 390  
*Benzécri J.-P.* 1, 67, 94, 108, 138, 145,  
151, 167, 171, 172, 190, 191, 193,  
194, 271, 281, 298, 330, 336, 368, 369  
*Beran R.* 395  
*Berge C.* 164  
*Besley D. A.* 223  
*Besse P.* 395  
*Birch M. W.* 284  
*Bishop Y.* 284  
*Blasius J.* 369  
*Bock H. H.* 400, 401, 402, 404  
*Boeswillwald E.* 57

*Bolasco S.* 337  
*Bourlard H.* 282  
*Bouroche J.-M.* 302  
*Bourret P.* 283  
*Box G. E. P.* 360  
*Breiman L.* 302, 313  
*Brillouin L.* 373  
*Bruynooghe M.* 173  
*Burt C.* 108, 111, 135  
*Burtschy B.* 333, 334

### C

*Cacoullos T.* 251  
*Callant C.* 270  
*Caraux G.* 9  
*Carlier A.* 334, 337  
*Carrol J. D.* 108, 338, 339, 347, 398  
*Casin P.* 347  
*Cattel R. B.* 341,374  
*Caussinus H.* 332, 333  
*Cazes P.* 100, 122, 135, 190, 193, 194,  
334, 336  
*Celeux G.* 251, 269, 276, 284, 299, 302,  
313, 401  
*Chabanon C.* 283  
*Chandon J.-L.* 145  
*Chang J. J.* 338  
*Chang T. C.* 361  
*Chateau F.* 268  
*Chatterjee S.* 223  
*Cheng B.* 278, 282  
*Chernoff H.* 7  
*Chessel D.* 324, 336  
*Choudary Hanumara R.* 361  
*Christensen R.* 284, 289  
*Clemm D. S.* 361  
*Cliff A. D.* 334  
*Cliff N.* 339  
*Cochran W. G.* 237  
*Cook R. D.* 245  
*Cooper M. C.* 403, 404  
*Coppi R.* 337  
*Cormack R. M.* 145, 212  
*Cornfield J.* 292  
*Corsten L. C. A.* 362

Cottrell M. 282  
 Cox D. R. 237, 245, 290, 292, 300  
 Cramer H. 378  
 Critchley F. 380

## D

D'Ambra L. 326  
 Dagnélie P. 237  
 Darmois G. 319  
 Darroch J. N. 290  
 Daudin J.-J. 295, 395  
 Davis A. W. 376  
 Davis C. 381  
 Day N. E. 401  
 De Leeuw J. 108, 295, 298, 299, 337  
 Decarli A. 295  
 Delecroix M. 267  
 Deming W. E. 384  
 Devaud J.-M. 292  
 Devijver P. 251  
 Deville J.-C. 384  
 Diaconis P. 395  
 Diday E. 148, 151, 152, 403  
 Dillon W. R. 251  
 Dobson A. 284  
 Dodge Y. 29  
 Dolédec S. 336  
 Domenges D. 29  
 Draper N. R. 223  
 Droesbeke J.-J. 290  
 Drouet d'Aubigny G. 374, 398  
 Dubes R. C. 400, 402  
 Dubuisson B. 268, 283  
 Duda R. O. 251  
 Dugué D. 360

## E

Eastment H. T. 395  
 Eckart C. 16, 23, 53, 298, 400  
 Efron B. 387, 388, 392, 395, 398  
 Escofier B. 67, 81, 108, 121, 270, 299,  
 331, 334, 336, 344, 352, 353, 379, 381  
 Escoufier 89, 190, 295, 299, 332, 342,  
 343  
 Everitt 401

## F

Falguerolles (de) A. 290  
 Falissard B. 29  
 Faraj A. 334  
 Farebrother R. W. 225

Fénelon J.-P. 5, 359  
 Ferré L. 395  
 Fichet B. 29, 225  
 Fienberg S. E. 284  
 Finch M. D. 374  
 Finden C.R. 335  
 Fine J. 245, 247, 290  
 Fisher R. A. 1, 52, 67, 237, 251, 360  
 Fisher W.D 121, 152  
 Fix E. 267, 268  
 Florek K. 165  
 Flury B. 376  
 Forgy E. W. 148  
 Fort J.-C. 282  
 Fourgeaud C. 223  
 Francisco C. A. 374  
 Friedman J. 266, 270, 302, 332  
 Fukunaga K. 251  
 Furnival G. M. 244

## G

Gabriel K. R. 17, 290  
 Gallego F. J. 122  
 Gallinari P. 283  
 Garnett J. C. 246  
 Geary R. C. 331, 334  
 Geisser S. 269  
 Gifi A. 16, 380, 395, 398  
 Gilula Z. 190, 295  
 Girshick M. A. 360, 375, 376  
 Gnanadesikan R. 251, 333  
 Goldstein M. 251  
 Good P. 3  
 Goodman L. A. 284, 295, 300, 310  
 Gordon A. D. 145, 155, 335, 400  
 Govaert G. 190  
 Gower J. C. 17, 166, 339, 403  
 Green P. J. 392  
 Greenacre M. 190, 369, 382, 395  
 Grelet Y. 122  
 Grizzle J. E. 289  
 Guegen A. 302  
 Guttman L. 67, 108, 374, 398

## H

Haberman S. J. 284, 288, 289  
 Hall D. J. 148, 152  
 Hand D. J. 251, 267, 269, 388, 401  
 Hardy A. 404  
 Harman H. H. 32, 246  
 Harshman R. A. 337, 338  
 Hart P. E. 251  
 Harter H. L. 223

Hartigan J. A. 145, 400  
 Hatabian G. 398  
 Hayashi C. 67, 108, 338, 380  
 Hayashi F. 338  
 Heiser W. J. 94, 189  
 Hertz J. 283  
 Highleyman W. H. 269  
 Hill M. O. 67  
 Hirschfeld H. D. 67  
 Hodges J. L. 267, 268  
 Holland P. 284  
 Holmes-Junca S. 382, 392, 395  
 Hornik K. 281, 283  
 Horst P. 32, 108, 347  
 Hosmer D. W. 292  
 Hotelling H. 32, 213  
 Houzel Y. 199  
 Hsu P. L. 360  
 Huber P. J. 225  
 Hudon G. 295  
 Hunter W. G. 237  
 Hurley J. R. 341

## I

Iemba A. F. 234

## J

Jain A. K. 400, 404  
 Jambu M. 145, 162, 190  
 Jeffreys H. 371  
 Jmel S. 290  
 Johnson S. C. 156  
 Jolliffe I. 374  
 Jones M. C. 332  
 Jordan C. 17  
 Joreskog K. 250  
 Joussein B. 383

## K

Kahan W. M. 381  
 Kaiser H. F. 374  
 Kamp Y. 282  
 Kanal L. 251  
 Karhunen J. 281  
 Kato T. 270, 379  
 Kaufman L. 145  
 Kazmierczak J.-B. 53  
 Kendall M. G. 51, 362, 404  
 Kettnering R. J. 195  
 Kharchaf I. 195  
 Kiers H. A. L. 337  
 Kittler J. 251  
 Kohonen T. 281, 282  
 Krishnaiah P. R. 251, 361

Kroonenberg P. M. 337  
 Kruskal J. B. 9, 164, 398  
 Kruskal W. 310  
 Krzanowski W. J. 380, 395  
 Kshirsagar A. M. 362  
 Kubushishi T. 404  
 Kuder D. F. 67  
 Kullback S. 289, 371, 372

## L

L'Hermier des Plantes H. 342  
 Lachenbruch P. A. 251, 269  
 Lafosse R. 339  
 Lancaster H. O. 362, 378  
 Lance G. N. 155, 404  
 Laplace P. S. 225  
 Lauro N. C. 295, 326  
 Lavit C. 342  
 Lawley D. N. 250, 375  
 Le Calvé G. 225  
 Le Foll Y. 334  
 Lebeaux M.-O. 145, 190  
 Lebreton J.-D. 324  
 Lechevallier Y. 313  
 Leclerc A. 132, 135, 276, 295  
 Lelu A. 283  
 Leibler R. A. 289  
 Lemeshow S. 292  
 Lerman I. C. 145  
 Leroux B. 270, 379, 381

## M

Mahalanobis P. C. 251, 351  
 Malinvaud E. 223, 378  
 Mallows C. L. 2, 245, 382  
 Marcotorchino F. 9  
 Masson M. 108, 347  
 Matheron G. 331  
 Matusita K. 404  
 Maxwell A. E. 250  
 McCullagh P. 245, 289  
 McLachlan G. J. 251, 268  
 McQueen J. 148, 152  
 McQuitty L. L. 172  
 Mehta C. R. 3  
 Mehta M. L. 361  
 Méot A. 334  
 Messenger R. C. 302  
 Meulman J. 395  
 Meyer R. 29  
 Michelat G. 205  
 Mickey M. R. 269  
 Milgram M. 283  
 Miller R. G. 385  
 Milligan G. W. 403, 404

Mirkin B. 195  
 Mollière J.-L. 403  
 Mom A. 331, 332, 334  
 Mood A. M. 360  
 Moran P. A. P. 334  
 Moreau J. 334  
 Moreau J. V. 404  
 Morgan J. M. 302  
 Mosteller F. 223  
 Muirhead R. J. 360, 376  
 Mulaik S. A. 246  
 Murtagh F. 145

## N

Nakache J.-P. 108, 251, 276, 284, 299,  
 302  
 Nakhilé F. 131  
 Neave H. R. 390  
 Nelder J. A. 245, 289, 292  
 Newman T. G. 387  
 Nishisato S. 108

## O

O'Neill M. E. 362  
 Odell P. L. 387  
 Ohsumi N. 181, 404  
 Oja E. 281, 282  
 Olshen R. A. 302  
 Ord J. K. 334

## P

Pagès J.-P. 342  
 Pagès J. 344, 352  
 Palm R. 234  
 Parzen E. 266, 267, 289  
 Pearson K. 1, 32, 73  
 Perruchet C. 400  
 Pillai K. C. S. 361  
 Pinson S. 145  
 Pousse A. 376  
 Price B. 223  
 Prim R.C. 165  
 Proriot J. 279

## Q

Quenouille M. 385

## R

Ramsay J. O. 398  
 Rao C. R. 32, 53, 223, 319, 378, 398

Rasson J.-P. 404  
 Richardson M. 67  
 Ripley B. D. 278, 334, 387, 402  
 Ritov Y. 295  
 Ritter H. 282  
 Robert Ch. 264  
 Robert P. 343  
 Romeder J.-M. 269  
 Rosenblatt M. 266, 267  
 Ross G. 166  
 Rousseau R. 195  
 Rousseeuw P. J. 145  
 Roux M. 145  
 Roy S. N. 360  
 Rubin H. 250  
 Ruiz A. 333  
 Rumelhart D. E. 279

## S

Sabatier R. 323, 334  
 Saporta G. 223, 275, 276, 347, 360, 398  
 Sarle W. S. 403  
 Stärndal C.-E. 384  
 Schiffman S. S. 9, 398  
 Schönemann P. H. 339  
 Searle S. E. 223  
 Seber G. A. F. 223  
 Shepard R. N. 9, 398  
 Sibson R. 332  
 Silverman B. W. 267  
 Simon M. 205  
 Strat J. A. 282  
 Smith H. 223  
 Sneath P. H. A. 145, 155, 156  
 Snee R. D. 68  
 Sokal R. R. 145, 155  
 Sonquist J. A. 302  
 Spearman C. 51, 246  
 Srivastava M. S. 395  
 Stauffer D. F. 395  
 Steinberg D. M. 237  
 Stephan F. F. 384  
 Stone M. 269, 290, 302  
 Stuart A. 362  
 Sylvester J. J. 16

## T

Tabard N. 5, 108, 329, 332  
 Taguri M. 7  
 Tanaka Y. 380  
 Tenenhaus M. 94, 108, 295, 302

Ter Braak C. J. F. 324  
Theil H. 223  
Thionet P. 384  
Thom R. 373  
Thompson W. A. 361  
Thorndike L. M. 148, 152  
Thurstone L. L. 246  
Tibshirani R. J. 282, 387  
Titterington D. M. 278, 282  
Tomassone R. 223, 237, 251, 267, 385  
Toussaint G. T. 269  
Trécourt P. 295  
Tucker L. R. 337, 339  
Tukey J. W. 2, 223, 282, 332, 382, 385  
Turlot J.-C. 347

## V

Van Buuren S. 189  
Van Cutsem B. 29  
Van der Heijden P. G. M. 295, 298, 337  
Van Rijkevorsel J. 94, 122  
Volle M. 29

## W

Waikar V. B. 361  
Wakimoto K. 7

Ward J. H. 191, 194, 195  
Waternaux C. M. 376  
Wedderburn R. W. M. 245, 289, 292  
Weinberg S. L. 398  
Weisberg S. 245  
Werbos P. J. 279  
Wermuth N. 245, 290  
Whittaker J. 245, 290  
Wilkinson J. H. 270, 379, 380  
Williams W. T. 155, 404  
Wilson R. W. 244  
Wish M. 9, 398  
Wishart J. 359, 360, 361  
Wishart D. 168  
Wold S. 271, 395  
Wong M. A. 177, 404  
Worsley K. J. 295, 299

## Y

Yates F. 52  
Yenyukov I. S. 333  
Young F. W. 108,  
Young G. 16, 23, 53, 108, 298, 400  
Young G. A. 388

## Z

Zeng G. 402



# Index des matières

## A

- Agrégation
  - autour de centres mobiles 146, 148
  - hiérarchique 155
  - mixte 177
  - selon la variance 167
- Aides à l'interprétation (cf. règles)
- Algorithme
  - de classification 156
  - de classification mixte 203
  - de Florek 165
  - de Kruskal 164
  - de Prim 165
  - de segmentation 304
  - mixte 147, 177
- Analyse canonique 67, 210, 213,  
des correspondances 324  
généralisée 212, 338, 347, 348
- Analyse de contiguïté 331, 332, 333
- Analyse de covariance 210, 223, 241,  
242
- Analyse de variance 210, 223, 237,  
254, 269, 284, 347  
multidimensionnelle, 268
- Analyse des corrélations partielles 53
- Analyse des correspondances 14, 67-  
135
  - internes 336
  - multiples 14, 89, 108, 113, 185,  
219, 220, 269, 276, 347, 374, 392
  - multiples conditionnelles, 336
  - non-symétrique 325
- Analyse des différences locales 334,  
336
- Analyse des rangs 51, 359
- Analyse (factorielle) discriminante
  - 67, 210, 213, 218, 251-283,  
barycentrique 133, 276
  - quadratique 263, 265, 273
  - qualitative 275
- Analyse du nuage résiduel 322
- Analyse en axes principaux 262
- Analyse en composantes principales  
13, 26, 32- 57, 198, 234, 272, 281,  
333, 345, 357, 380, 395
- Analyse factorielle classique des  
psychologues (analyse en facteurs  
communs et spécifiques) 210, 246

- Analyse factorielle multiple 212, 338,  
344, 345, 352
- Analyse générale 14, 15 -31
- Analyse
  - inter-classes 335
  - interne (intra, intra-classes) 212,  
272, 335, 336
  - lissée 334
  - locale 212, 327, 331
  - logarithmique 53
  - mono-(multi) factorielle 246
  - multivariée descriptive généralisée  
(MDSG) 347
  - partielle/projetée 212, 319, 323,  
333, 336
- Analyse procrustéenne 212, 339-341,  
346, 391
  - orthogonale 338, 339, 340
  - sans contrainte 341
- Apprentissage 264
  - par coeur 275
- Arbre
  - de décision binaire. 302
  - de longueur minimale 159, 163-166
- Articulation exploration-inférence 300
- Automatic Interaction Detection* 302
- Axiome de réductibilité 173

## B

- Back-propagation.* 279
- Bande du tableau de Burt 276
- Base orthogonale hiérarchisée 234
- Bootstrap* 4, 7, 265, 379, 385, 386, 389,  
395, 398, 403
- Budget-temps 59, 200

## C

- Calculs de stabilité 358, 379, 389
- CART (méthode) 302
- CESP 57, 103
- Classement 252, 317
- Classification

- à partir des facteurs 187
- ascendante hiérarchique 155-175
- autour de centres mobiles 148
- mixte 177
- supervisée 317
- Codage
  - condensé 109
  - des variables 382, 383
  - des variables nominales 238
  - disjonctif (complet) 368, 117, 259
- Coefficient
  - de contiguïté 329, 330
  - de corrélation 37, 41, 42, 197, 214, 387
  - de corrélation canonique 216, 262
  - de corrélation multiple 228, 229, 230, 245
  - de corrélation partielle 319, 320
  - de régression 231, 258
  - Rv d'Escoufier et Robert 343
- Colinéarités 234, 238
- Comparaisons multiples 125, 290, 301
- Complémentarité analyses / classification 189, 199
- Compression de signal 281
- Conditions de vie et aspirations des Français 135, 199
- Contacts-média 103, 104, 389, 390
- Contribution 94, 95, 105, 121
  - relative (cf. cosinus carrés)
- Cosinus carrés 95, 96, 105
- Couche cachée (cf. perceptron)
- Covariance
  - locale 328, 331
  - partielle 320
- Credit-scoring 251
- Critère
  - d'affinité 404
  - d'agrégation de Ward 190
  - de Kullback-Leibler (déviance) 299
  - de la médiane 173
  - de pureté maximale 305, 310
  - de variance résiduelle minimale 307
  - de Ward généralisé 170, 173, 403
  - de Wilks 402
  - d'agrégation 156
  - externe (cf. procédures)

## D

- Décomposition aux valeurs singulières
  - 16, 271, 281, 300, 337, 340, 348, 400
- Delete-d Jackknife 387
- Dendrogramme 155, 175, 191, 199, 203

- Description statistique des classes 177, 181, 202, 188
- Diagnostic automatique 251
- DISQUAL (méthode) 275
- Dissection 404
- Distance
  - de Hellinger 29
  - de Mahalanobis (ou généralisée) 258, 265, 266, 352, 372
  - de Mahalanobis globale 264, 272, 274
  - de Mahalanobis locale 264, 269, 273, 274
  - du  $\chi^2$  73, 81, 100, 114, 119, 148, 299
  - du plus petit saut maximal 162
  - euclidienne 34, 73, 148, 271, 274
  - $L_1$  28, 225
  - ultramétrique 159
- Distribution du  $\chi^2$  (cf. loi)
- Divergence de Jeffreys 371

## E

- Échantillon d'apprentissage 252, 268, 269, 273, 274, 280, 313
- Échantillon-test 273, 303, 313, 373, 375
- Écologie (relevés écologiques) 324, 370
- Effet
  - Guttman 93, 94, 330
  - de chaîne 167
- Éléments supplémentaires 27, 43, 99, 107, 122, 296, 391
- Éléments terminaux 158
- Enquêtes 57, 199, 384
- Entropie de Shannon 310
- Épreuves empiriques de stabilité 382
- Équi-divisantes (divisions) 316, 318
- Équi-réductrices (divisions) 318
- Équivalence distributionnelle 74, 81, 82
- Erreurs de mesure 382, 383
- Estimation directe (non-paramétrique)
  - de la densité 266, 267
- Exemples d'application 57, 103, 135, 273

## F

- Facteur de taille 56, 199
- Facteurs communs/ spécifiques 247
- Fluctuations d'échantillonnage 382, 384, 385

Formes fortes (groupements stables)  
152, 153

Formule

de reconstitution des données 22,  
23, 89, 192, 340, 365  
de transition (cf. relations)

## G

Graphe

associé à une partition 336  
complet, valué, connexe, partiel  
162-165  
régulier 330, 332

## H

Hiérarchie de partitions, 155  
Hiérarchie indicée 158, 159  
Histogramme des indices 180, 403  
*Homogeneity analysis* 108, 395  
*Hybrid clustering* 177  
Hypothèse  
d'homogénéité spatiale 402  
d'indépendance/ nulle 70, 91, 183,  
294, 299, 359, 378, 399  
de normalité 231, 266

## I

Indépendance des taux d'inertie et de la  
trace 364  
Indice  
de diversité de Gini 310  
de niveau 159, 170  
Individus supplémentaires (cf.  
éléments)  
*INDSCAL* (modèle) 338, 398  
Information de Shannon-Wiener 371  
Interstructure 342, 346  
Intervalles de confiance d'Anderson  
376  
Intrastructure 343, 345  
*Iterative proportional fitting* 384

## J

*Jackknife* 4, 379, 385, 398  
*Joint correspondence analysis* 369

Juxtapositions de tableaux de  
contingence 391

## K

*K-means* (méthode) 152

## L

Lecture directe (algorithmes) 148, 281  
Ligne supplémentaire (ou illustrative)  
(cf. éléments)

Loi

binomiales 246  
de Fisher 233  
de Laplace-Gauss (cf. loi normale)  
de Poisson 246, 288  
de Student 232  
de Wishart 359, 360, 361  
des valeurs propres 359, 362  
du  $\chi^2$  360  
gamma 246  
hypergéométrique 184  
multinomiale 288, 359  
normale (ou de de Laplace-Gauss)  
52, 182, 184, 225, 245, 359  
normale sphérique 267

## M

Matrice associée à un graphe  
symétrique 369  
Matrice  
de contiguïté 328  
de Wishart 196, 362, 364, 375  
des corrélations 61  
des corrélations locales 331  
Matrice des covariances  
49, 213, 248, 254, 265, 269, 270,  
272  
locales 270, 331  
partielles 323  
intra classes et interclasses 327  
Matrices idempotentes (cf. opérateur  
projection) 217  
Méthode  
de Newton-Raphson 246, 289  
de validation 379  
*delta* 398  
des  $m$  (ou  $k$ ) plus proches voisins  
(cf. règle)

des moindres carrés 226, 234  
 des moindres carrés itératifs 289  
 des moindres carrés pondérés  
 itératifs 246  
 des noyaux 266  
 du maximum de vraisemblance 246,  
 288, 293  
 Méthodes neuronales 266, 277  
 Méthode STATIS 212, 338, 342  
 Métrique (cf. distance)  
 Modalités supplémentaires (cf. aussi  
 éléments) 123, 140  
 Modèle  
 bayésien d'affectation 264, 265  
 linéaire 210, 223, 265  
 linéaire généralisé 245, 246, 289,  
 292  
 log-linéaire 211, 246, 284, 285, 295  
 log-linéaire hiérarchique 287  
 log-linéaire non-saturé 299  
 logistique 291-295, 299  
 Modèles  
 à erreurs sur les variables 246  
 auto-organisés 280  
 de mélanges (en classification) 401  
 de partitions fixes 401  
 de variables latentes 246  
 décomposables 290  
 fonctionnels (effet fixes) 246  
 graphiques 290  
 supervisés / non-supervisés 280  
 RC de Goodman, 300  
 structurels (à effets aléatoires) 246  
*Multidimensional scaling* 338, 398  
 MULTISCALE (modèle) 398

## N

Neurones, neuro-mimétique 277-282  
 Non-réponses 109, 204  
 Normes (cf. distances)  
 Nuage moyen ou compromis 343  
 Nuées dynamiques 146, 148

## O

Opérateur de projection 217, 228, 272,  
 348  
 Ordonnance 158  
 Orthogonalisation de Gram-Schmidt  
 218

## P

PARAFAC (modèle) 338  
 Perceptron multi-couches 278  
 Polythétique (classe) 204  
 Pourcentage  
 de bien classés 268, 273, 375  
 de variance (d'inertie) 24, 54, 61  
 Pouvoir discriminant 256  
 Préordonnance 158  
 Problèmes  
 mal posés 234, 270  
 pauvrement posés 234, 270  
 Procédure  
 d'élagage 303  
 de consolidation 180  
 Procédures externes (de validation)  
 374, 389, 400, 404  
 Processus  
 de Poisson généralisé 402  
 de Poisson stationnaire 404  
*Projection poursuit* (projections  
 révélatrices) 331, 332

## R

Racine canonique (cf. coeff. de  
 corrélation)  
 Reconnaissance des formes 251, 268,  
 277  
 Reconstitution des données  
 (cf. formule de reconstitution)  
 Rééchantillonnage 4, 265, 358, 379, 385  
 Règle  
 d'affectation 309  
 de Bayes 266  
 des  $m$  ( ou  $k$ ) plus proches voisins  
 267, 268, 333  
 d'affectation 263, 265  
 Règles d'interprétation 89, 120, 181  
 Régression  
 logistique 266, 279, 284, 290, 317  
 (simple, multiple) 210, 213, 223,  
 227, 257, 269, 280, 303, 306, 347  
 multiple (variables nominales) 237,  
 284  
 pas-à-pas 317  
 régularisée (cf. régularisation)  
 sur composantes principales 233  
 sur variables mixtes 241  
 Régularisation  
 de la régression 233, 234, 238  
 en analyse discriminante 269

en analyse canonique généralisée 353  
 Relation  
 de contiguïté 329  
 de contiguïté *a posteriori* 332  
 Relations de transition 26, 75, 78, 85, 115, 261  
 Réplications 390, 391, 395  
 Représentation simultanée 45, 46, 67, 78, 87, 99, 118  
 Réseaux de neurones (neuro-mimétiques) 277-282  
 Robuste (robustesse) 225, 82, 268

## S

Saut  
 maximal 156, 173  
 minimal (*single linkage*) 156, 159-166, 173  
 Score (discriminant) (*scoring*) 265, 276, 277  
 Segmentation 211, 302, 317  
*Self organizing maps* 282  
 Séries chronologiques de tableaux 211  
 Simulation (d'échantillon) 4, 196, 385, 403  
 Sous-tableau du tableau de Burt 131  
 Stabilité  
 des axes (des facteurs) 379  
 des fonctions discriminantes 265  
 des formes 373, 379  
 externe/interne (cf. aussi procédures de validation) 382  
 STATIS (méthode) 344, 352, 395  
 Statistique  
 de Fisher 240, 241  
 de Student 293  
 de Wald 293  
 Stratégie de classification mixte 177  
 Structure  
*a priori* de l'ensemble des individus 327  
 de chaîne 327  
 de partition 335  
 de graphe 212, 327

## T

Table de Fisher 243  
 Tableau (table)  
 de contingence 67, 104, 390

de contingence de Burt 111, 112, 126, 137, 138, 276, 335, 369, 395  
 disjonctif complet 108, 110, 111, 113, 126, 272, 276  
 de contingence multidimensionnel 211, 298  
 multiple 212, 337  
 Tableaux pseudo-aléatoires. 391  
 Taux d'erreur apparent 268, 311  
 Taux d'erreur par resubstitution 268  
 Taux d'inertie 24, 129, 137, 368  
 Test  
 d'interaction 241  
 de sphéricité 367  
 du  $\chi^2$  91, 285, 289, 359  
 du rapport de vraisemblance 289  
 fishérien 2  
 Thémascope 202  
 Théorème  
 d'Eckart et Young (cf. analyse générale, décomposition aux valeurs singulières) 337  
 de Bayes 264, 291  
 de Gauss-Markov 230  
 de la limite centrale 182  
 de la médiane 171  
 de Wielandt-Hoffman 380, 381  
 Théorie  
 de l'information 289  
 de la perturbation 270, 395  
 des variables régionalisées 331  
 Tirage pseudo-aléatoire 387  
 Trajectoires 344, 346

## U

Ultramétrique 159, 162  
 inférieure maximale (sous-dominante) 159, 161

## V

Valeur pratique de l'information 373, 383  
 Valeur-test 64, 66, 123, 124, 125, 140, 181, 184  
 Validation croisée 4, 265, 268, 303, 375, 379, 385, 388, 389, 395  
 Validation des classifications 399

Validité et portée de résultats 89, 296,  
383

Variables

endogènes (dépendantes) 223  
actives 28, 296  
canoniques 215  
exogènes (explicatives) 223, 319  
instrumentales 53, 319  
supplémentaires (cf. aussi:  
éléments) 42, 57, 236

Variance

externe 256  
inter-classes 403

intra-classes 255

locale 328, 329

Variogramme 331

Voisins réciproques 171-173

## **Z**

Zones de confiance 358, 379, 389, 395

Zones de garde 396, 397

Achévé d'imprimer  
sur les presses de l'imprimerie  
Arts Graphiques du Perche  
28240 Meaucé

Dépôt légal Octobre 1995  
Imprimé en France

**Ludovic Lebart**

**Alain Morineau**

**Marie Piron**

# **Statistique exploratoire multidimensionnelle**

**C**e livre s'adresse aux chercheurs, ingénieurs, professeurs, étudiants qui sont confrontés dans leurs travaux aux recueils de données multidimensionnelles (ou multivariées). Les enquêtes socio-économiques, épidémiologiques et de marketing en sont des exemples courants. Mais les relevés écologiques, les bases documentaires, les données de télédétection, les mesures de contrôle de qualité constituent à l'heure actuelle des domaines en développement rapide.

Destiné à un public assez large, appuyé sur de nombreux exemples, l'ouvrage présente les concepts de base et les fondements des méthodes exploratoires de base (analyses factorielles, classification), et rend compte aussi des travaux les plus récents. Les auteurs insistent sur la complémentarité de ces méthodes et sur leur insertion dans l'arsenal des méthodes statistiques plus classiques. L'articulation avec les techniques explicatives et prévisionnelles, la validité et la stabilité des résultats sont deux points sensibles largement développés.

L'ouvrage peut être lu à plusieurs niveaux : celui du praticien, celui de l'utilisateur exigeant, enfin celui du chercheur en méthodologie statistique.



Code 042886  
ISBN 2 10 002886 3

