

REVUE DE STATISTIQUE APPLIQUÉE

BRIGITTE ESCOFIER

Une représentation des variables dans l'analyse des correspondances multiples

Revue de statistique appliquée, tome 27, n° 4 (1979), p. 37-47

http://www.numdam.org/item?id=RSA_1979__27_4_37_0

© Société française de statistique, 1979, tous droits réservés.

L'accès aux archives de la revue « *Revue de statistique appliquée* » (<http://www.sfds.asso.fr/publicat/rsa.htm>) implique l'accord avec les conditions générales d'utilisation (<http://www.numdam.org/conditions>). Toute utilisation commerciale ou impression systématique est constitutive d'une infraction pénale. Toute copie ou impression de ce fichier doit contenir la présente mention de copyright.

NUMDAM

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

UNE REPRÉSENTATION DES VARIABLES DANS L'ANALYSE DES CORRESPONDANCES MULTIPLES

Brigitte ESCOFIER *

RESUME

Nous proposons ici une représentation des variables à l'issue de l'analyse des correspondances d'un tableau disjonctif complet. En effet, dans les programmes classiques, seules les modalités des variables apparaissent, mais la variable elle-même, prise globalement, n'apparaît pas.

Pour cela, nous associons à chaque variable V un opérateur Ω_V , l'opérateur de projection orthogonale sur le sous espace des indicatrices de la variable V , et nous associons à chaque facteur F l'opérateur de projection orthogonale sur F , Ω_F . L'espace des opérateurs étant muni d'une métrique euclidienne, nous pouvons projeter Ω_V sur les opérateurs associés aux facteurs.

Cette représentation met en évidence, à la fois les liaisons entre une variable et un facteur et les liaisons entre les variables. En effet, la coordonnée de l'opérateur Ω_V sur l'opérateur Ω_F est proportionnelle à la contribution de l'ensemble des modalités de la variable V à l'inertie de F . Et, le cosinus de l'angle entre les opérateurs Ω_V et $\Omega_{V'}$, associés à deux variables V et V' est une mesure de la liaison entre V et V' .

Un exemple commenté illustre la technique que nous proposons.

I – LES DONNEES, LE TABLEAU DISJONCTIF COMPLET ET SON ANALYSE (cf. [1] et [4])

Pour analyser un ensemble de variables définissant chacune une partition d'un ensemble d'individus, on lui associe souvent un tableau disjonctif complet que l'on traite par l'analyse des correspondances. Ces variables peuvent être des variables qualitatives à plusieurs modalités, des questions avec plusieurs réponses possibles ou bien encore des variables quantitatives codées par classe.

Le tableau disjonctif complet associé à l'ensemble de variables V_1, \dots, V_q est le tableau des indicatrices des partitions définies par ces q variables. Il croise donc l'ensemble des individus I , et l'ensemble J des modalités des q variables. Il se compose uniquement de 0 et de 1. La valeur 1 au croisement de la j -ième colonne indique que l'individu i a la modalité j . Dans le groupe de colonnes associé aux modalités d'une même variable, il y a sur une ligne un 1 et un seul. Le tableau a donc la forme suivante :

* Département de Mathématiques – I.N.S.A., B.P. 14A – 35031 RENNES CEDEX.

	$\overbrace{\hspace{10em}}^J$		
	$\overbrace{\hspace{2em}}^{V_1}$	$\overbrace{\hspace{2em}}^{V_2}$	$\overbrace{\hspace{2em}}^{V_q}$
$\left\{ \begin{array}{l} \\ \\ \end{array} \right.$	001	0001	01
	100	0100	01

Dans l'analyse des correspondances de ce tableau disjonctif complet, on obtient des couples de facteurs définis respectivement sur l'ensemble I des individus et l'ensemble J des modalités de toutes les variables. Mais, les variables elles-mêmes n'apparaissent pas dans les résultats.

Ceux qui appliquent l'analyse des correspondances aux tableaux disjonctifs complets savent que pour interpréter les résultats il est nécessaire de connaître les relations entre un facteur et une variable. Comme ils ne disposent pas de techniques systématiques, ils font quelques calculs simples. Par exemple, ils mesurent l'importance d'une variable V pour un facteur F en calculant la contribution de l'ensemble des modalités de V à l'inertie de F. Ou bien, puisqu'une variable est une partition de l'ensemble I en classes, ils calculent, pour ces classes l'inertie interclasse du facteur F.

Jusqu'à présent les résultats de l'analyse des correspondances n'étaient pas utilisés pour estimer les liaisons entre les variables, seules les liaisons entre les modalités de ces variables étant mises en évidence.

II – OPERATEUR Ω_V ASSOCIE A LA VARIABLE V (cf. [2] et [5])

Les variables que nous considérons définissent une partition de l'ensemble I. La donnée d'une partition de I en s éléments est équivalente à la donnée des s fonctions indicatrices de cette partition. Ces s fonctions sont des éléments de l'espace vectoriel des fonctions numériques définies sur I, espace noté \mathbf{R}^I . Elles engendrent un sous espace vectoriel \mathcal{E}_V de \mathbf{R}^I de dimension s. Et ce sous espace \mathcal{E}_V détermine les s fonctions indicatrices.

La somme de toutes les fonctions indicatrices est la fonction constante sur I, égale à 1. Le sous espace \mathcal{E}_V contient donc cette fonction notée 1_Δ . Il suffit donc de connaître le sous espace E_V de \mathcal{E}_V orthogonal à 1_Δ pour connaître la partition de I. La notion d'orthogonalité dans \mathbf{R}^I est celle de la métrique usuelle.

Un sous espace est caractérisé par l'opérateur de projection orthogonale sur ce sous espace. A une variable V on peut donc associer par une application injective, l'opérateur de projection orthogonale sur le sous espace E_V que nous venons de définir. Nous notons cet opérateur Ω_V .

III – OPERATEUR Ω_F ASSOCIE AU FACTEUR F

Un facteur F sur I est un élément de l'espace \mathbf{R}^I , il engendre donc un sous espace {F} de \mathbf{R}^I de dimension 1. A un facteur F, nous associons l'opérateur de projection orthogonale sur ce sous espace, et nous noterons cet opérateur Ω_F .

IV – PRODUIT SCALAIRE DE DEUX OPERATEURS SYMETRIQUES

La trace du composé de deux opérateurs symétriques définit un produit scalaire sur l'espace des opérateurs symétriques de \mathbf{R}^I (cf. [2]).

Rappelons que si E et F sont deux sous espaces de \mathbf{R}^I , le produit scalaire des opérateurs de projection orthogonale sur ces sous espaces est la somme des cosinus carrés des angles canoniques entre E et F. Le carré de la norme d'un opérateur de projection est donc égal à la dimension de son image.

V – PROJECTION DE $\Omega_V/\|\Omega_V\|$ SUR Ω_F . LIENS ENTRE F ET V

Nous proposons de projeter les opérateurs normés associés aux variables sur les opérateurs Ω_F . Nous montrons dans ce paragraphe que la coordonnée de $\Omega_V/\|\Omega_V\|$ sur Ω_F met en évidence les liens entre F et V.

Puisque Ω_F et Ω_V sont des projecteurs, leur produit scalaire est égal à la somme des carrés des corrélations canoniques entre l'espace {F} et l'espace E_V . La dimension de {F} étant égale à 1, il y a un seul coefficient de corrélation canonique qui est le cosinus de l'angle θ entre F et E_V .

La dimension de E_V étant égale à s-1, où s est le nombre de modalités de V, on a $\|\Omega_V\| = \sqrt{s-1}$. La projection de l'opérateur normé $\Omega_V/\|\Omega_V\|$ sur Ω_F est donc égale à :

$$\cos(\Omega_F, \Omega_V) = \frac{\langle \Omega_F, \Omega_V \rangle}{\|\Omega_F\| \|\Omega_V\|} = \frac{\cos^2 \theta}{\sqrt{s-1}}$$

puisque F étant de dimension 1, $\|\Omega_F\| = 1$.

Relation avec l'inertie des modalités de la variable

Dans l'analyse factorielle des correspondances du tableau disjonctif complet, le nuage de points représentant les modalités des variables est situé dans l'espace \mathbf{R}_I , dual de \mathbf{R}^I . Ces deux espaces sont munis de la métrique usuelle aux coefficients $1/\text{card}I$ et $\text{card}I$ près. Ces métriques induisent un isomorphisme entre \mathbf{R}^I et \mathbf{R}_I .

L'image du facteur F par cet isomorphisme est un axe d'inertie f du nuage des modalités. Et l'image d'une fonction indicatrice de la variable V est un vecteur proportionnel au vecteur joignant l'origine à la modalité en question de la variable V. Le sous espace \mathcal{E}_V des fonctions indicatrices de V a donc pour image le sous espace engendré par l'origine et les modalités de V.

L'angle entre F et \mathcal{E}_V est donc égal à l'angle entre l'axe d'inertie f associé à F et le sous espace engendré par les modalités de V. Or l'angle entre F et \mathcal{E}_V est aussi égal à l'angle θ entre F et E_V puisque F est orthogonal au sous espace 1_{Δ} des fonctions constantes sur I.

Les vecteurs joignant l'origine aux modalités de V sont orthogonaux. Si q est le nombre total de variables étudiées, l'inertie, par rapport à l'origine, de chaque modalité est $1/q$. Donc, l'inertie des projections des modalités de la variable V dans une direction quelconque du sous espace qu'elles engendrent est constante et égale à $1/q$. Dans une direction quelconque de \mathbf{R}_I , elle sera égale au produit par $1/q$ du cosinus carré de l'angle entre cette direction et leur sous espace.

Pour la direction f, ce sera $\cos^2\theta/q$. La contribution de la variable V à l'inertie du facteur F est donc, en pourcentage de l'inertie totale λ de F, $\cos^2\theta/(\lambda q)$. Ce pourcentage se calcule facilement avec les résultats de l'analyse des correspondances, car les pourcentages d'inertie fournis par chaque modalité figurent sur les listings. Ceci permet de calculer $\cos^2\theta$ et la coordonnée de $\Omega_V/\|\Omega_V\|$ sur Ω_F .

Cette coordonnée est liée aussi au pourcentage, extrait par F, de l'inertie de l'ensemble des modalités de V. L'inertie de ces s modalités, par rapport au centre de gravité du nuage – qui est aussi leur centre de gravité – est égale à $(s - 1)$, puisqu'elles engendrent un sous espace de dimension $s - 1$. Le pourcentage d'inertie extrait par la projection sur f est donc $\cos^2\theta/(s - 1)$.

Autres interprétations

La variable V définit des classes sur I. Pour ces classes, l'inertie interclasse du facteur, quotienté par son inertie totale est exactement $\cos^2\theta$. Lorsque $\cos^2\theta$ vaut 1, l'inertie interclasse est égale à l'inertie totale, tous les éléments d'une même classe sont confondus sur le facteur F. Donc, plus la coordonnée de $\frac{\Omega_V}{\|\Omega_V\|}$ s'approche de $1/\sqrt{s - 1}$ plus les éléments de ces classes se retrouvent groupés. Notons que $\cos^2\theta$ est aussi le rapport de corrélation de F par rapport à V.

En résumé

$$\begin{aligned}
 \text{Coordonnée de } \frac{\Omega_V}{\|\Omega_V\|} \text{ sur } \Omega_F &= \cos(\Omega_V, \Omega_F) \\
 &= \cos^2\theta/\sqrt{s - 1} \\
 &= \text{Contribution de V à F} \times \frac{q}{\sqrt{s - 1}} \\
 &= \text{Pourcentage d'inertie de V extrait par F} \\
 &\quad \times \sqrt{s - 1} \\
 &= \frac{1}{\sqrt{s - 1}} \times \frac{\text{Inertie interclasse de F}}{\text{Inertie de F}}
 \end{aligned}$$

VI — LIAISON ENTRE 2 VARIABLES. SA REPRESENTATION SUR LES GRAPHIQUES

Rappelons d'abord quelques mesures classiques de la liaison entre deux variables définissant chacune une partition d'un même ensemble I. Cette liaison est décrite par le tableau de contingence croisant les deux positions.

Elle est généralement mesurée par le Ψ^2 de ce tableau de contingence, ou le ϕ^2 qui est égale au χ^2 divisé par le cardinal de l'ensemble I. Ou bien encore par le T^2 de Tchuprov qui, si les variables V et V' ont respectivement s et t modalités, vaut $\phi^2/\sqrt{(s-1)(t-1)}$.

Ce dernier, le T^2 est compris entre 0 et 1. Il vaut 1 si, et seulement si, les deux partitions sont identiques. Le tableau de contingence est alors carré et diagonal. Il vaut 0 lorsque ces deux variables sont indépendantes. Plus T^2 est grand, plus les variables sont dépendantes. Le T^2 mesure la dépendance de deux variables qualitatives, comme le coefficient de corrélation linéaire mesure la liaison entre deux variables quantitatives.

Dans l'espace des opérateurs, on retrouve le ϕ^2 et le T^2 des deux variables V et V'. le produit scalaire des opérateurs Ω_V et $\Omega_{V'}$, associés à V et V' est égal à leur ϕ^2 . Et le produit scalaire des opérateurs normés $\frac{\Omega_V}{\|\Omega_V\|}$ et $\frac{\Omega_{V'}}{\|\Omega_{V'}\|}$ est égal à leur T^2 .

Les opérateurs associés à deux facteurs distincts sont orthogonaux. En effet, la métrique induite sur R^I par un tableau disjonctif complet est, à un coefficient près, la métrique usuelle. Deux facteurs distincts sont donc orthogonaux pour la métrique usuelle et les composés des opérateurs qui leur sont associés sont donc nuls.

La projection sur deux opérateurs Ω_F et $\Omega_{F'}$, associés à deux facteurs F et F' donne donc la projection sur le plan qu'ils engendrent. Nous obtenons ainsi une représentation graphique plane des opérateurs associés aux variables.

Comme il y a une bijection entre les facteurs et les opérateurs qui leur sont associés, on peut représenter directement les variables sur les graphiques de l'analyse des correspondances. Mais, nous proposons plutôt de tracer des graphiques annexes car la projection des modalités d'une variable et la projection de la variable elle-même sont de natures différentes puisqu'elles se font dans des espaces différents.

En résumé, le cosinus de l'angle entre les opérateurs normés représentant les variables V et V' dans l'espace Σ des opérateurs symétriques est égal au T^2 entre ces deux variables. Et nous projetons ces opérateurs sur les opérateurs associés aux facteurs qui forment une base incomplète orthonormée de Σ .

Si deux variables V et V' sont bien représentées par leur projection, l'angle entre les opérateurs Ω_V et $\Omega_{V'}$, et donc le T^2 entre les variables V et V' sera visible directement sur les graphiques, de la même façon qu'en analyse en composantes principales, on voit la corrélation entre deux variables quantitatives bien représentées sur les composantes principales.

La qualité de représentation d'une variable sur un sous espace est mesurée par la longueur de sa projection, puisque l'opérateur associé que l'on projette est normé.

Mais, nous avons montré au paragraphe précédent que la coordonnée de $\frac{\Omega_V}{\|\Omega_V\|}$ sur Ω_F était égale à $\cos^2\theta/\sqrt{s-1}$. Il faudra donc considérer un nombre de facteurs au moins égal à $s-1$ pour que l'opérateur $\frac{\Omega_V}{\|\Omega_V\|}$ puisse être bien représenté. Ceci s'explique bien, puisque l'opérateur associé à un facteur est de rang 1 alors que l'opérateur associé à la variable est de rang $s-1$. Mais cela implique que pour des variables ayant plus de 3 modalités, nous devons consulter plusieurs graphiques plans pour juger de leur proximité.

VII – LIAISON ENTRE UNE VARIABLE ET UN PLAN FACTORIEL

On peut se demander si la représentation de la variable V dans le plan $\{\Omega_F, \Omega_{F'}\}$ permet aussi de visualiser une liaison entre V et le plan factoriel $\{F, F'\}$. C'est ce que nous étudions ici.

La liaison entre V et $\{F, F'\}$ peut se mesurer par des pourcentages d'inertie : le pourcentage, extrait par le plan $\{F, F'\}$ de l'inertie des s modalités de V qui vaut $\frac{\cos^2\theta + \cos^2\theta'}{s-1}$; ou bien, le pourcentage d'inertie fourni par ces modalités à l'inertie du plan $\{F, F'\}$, qui vaut $\frac{\cos^2\theta + \cos^2\theta'}{q(\lambda + \lambda')}$; ou bien encore, par le pourcentage

d'inertie interclasse – pour les classes de V – dans l'inertie du plan $\{F, F'\}$: $\frac{\lambda \cos^2\theta + \lambda' \cos^2\theta'}{\lambda + \lambda'}$.

Toutes ces notions font intervenir plus ou moins directement la somme des coordonnées, $\frac{\cos^2\theta}{\sqrt{s-1}}$ et $\frac{\cos^2\theta'}{\sqrt{s-1}}$, de V . Elles n'apparaissent donc pas sur le graphique plan.

La liaison entre V et $\{F, F'\}$ peut aussi être mesurée par la liaison entre de sous espace E_V et le sous espace $\{F, F'\}$. Donc, par le cosinus de l'angle entre les opérateurs Ω_V et l'opérateur de projection orthogonale sur $\{F, F'\}$. L'opérateur de projection orthogonale sur $\{F, F'\}$ est la somme des opérateurs Ω_F et $\Omega_{F'}$. Il est donc situé dans le plan $\{\Omega_F, \Omega_{F'}\}$, sur la première bissectrice. Le cosinus de l'angle entre Ω_V et cet opérateur est donc égal à la longueur de la projection de Ω_V sur la première bissectrice du plan.

En conclusion, certains liens entre une variable V et le plan $\{F, F'\}$ sont visibles sur les graphiques, mais ce ne sont sans doute pas les plus faciles à interpréter.

VIII – REMARQUES

La description des variables qualitatives est tout à fait analogue à celle qui est obtenue dans [2] et [5] par une analyse en composantes principales des opérateurs associés aux variables.

Mais les opérateurs sur lesquels on projette les opérateurs représentant les variables ne sont pas les mêmes. Dans [2] et [5], ce sont les composantes principales obtenues en diagonalisant la matrice de T^2 . Alors qu'ici, ce sont les opérateurs associés aux facteurs de l'analyse des correspondances. Le but de cette technique n'est pas de décrire au mieux les T^2 , mais d'apporter une aide à l'interprétation des résultats de l'analyse des tableaux disjonctifs complets.

Cette comparaison amène cependant une remarque. Dans l'analyse en composantes principales des opérateurs, la première composante C maximise la quantité :

$$\sum_{i=1}^q \cos^2(\Omega_{V_i}, C)$$

L'opérateur attaché au premier facteur de l'analyse des correspondances possède aussi une propriété extrême. En effet, F est le vecteur de \mathbf{R}^1 maximisant :

$$\sum_{i=1}^q \cos^2 \theta_i$$

où θ_i est l'angle entre F et le sous-espace associé à la variable V_i . Donc, l'opérateur Ω_F associé à F est l'opérateur symétrique de rang 1 maximisant la quantité :

$$\sum_{i=1}^q \cos(\Omega_{V_i}, \Omega_F) \times \sqrt{s_i - 1}$$

IX – EXEMPLE

L'exemple que nous avons choisi est un tableau disjonctif complet de 15 variables quantitatives transformées en variables qualitatives par division en 4 classes. Ces variables sont diverses mesures prises sur 120 plantes.

Nous donnons dans le tableau 1 les valeurs des 5 premiers facteurs pour les 60 modalités des 15 variables. Les modalités sont désignées par une lettre caractérisant la variable suivi d'un chiffre indiquant la classe de la variable.

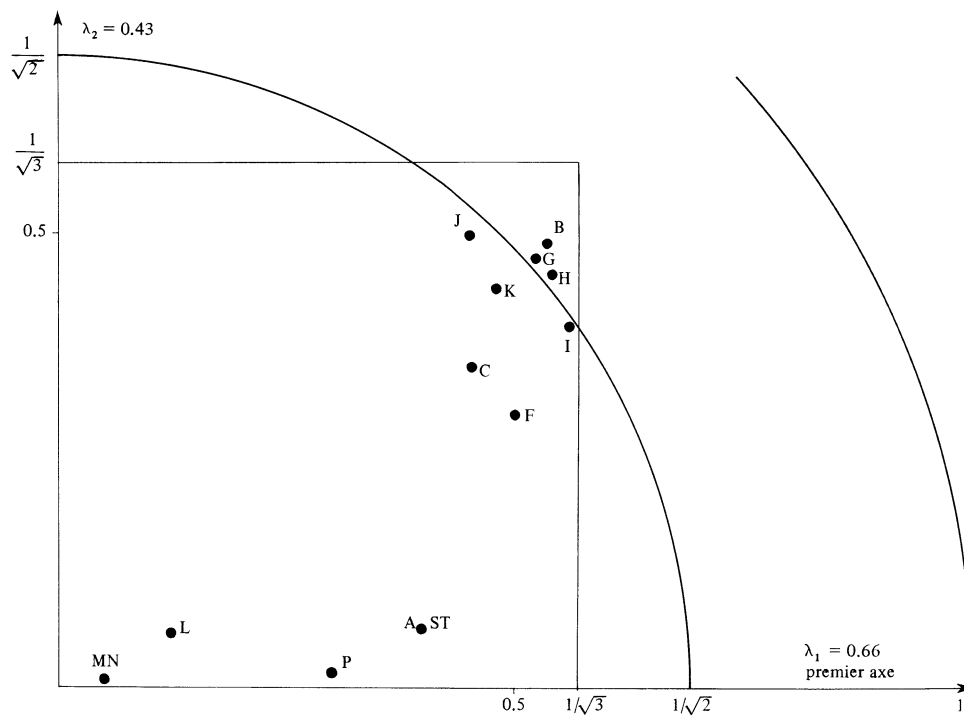
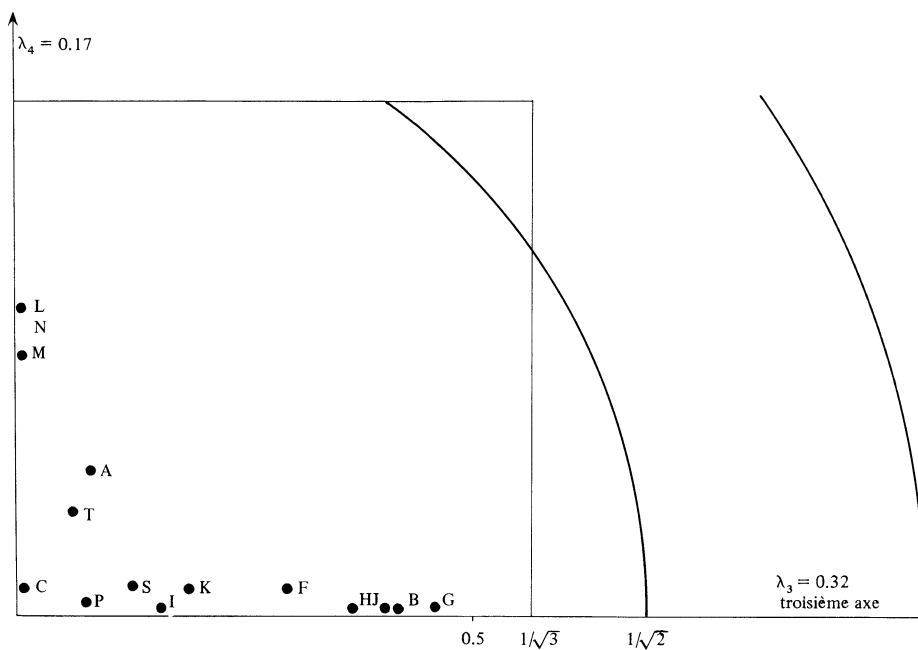
Nous avons calculé les "projections des variables" sur les 4 premiers facteurs, à partir de ce tableau et des valeurs propres. Commentons maintenant les graphiques obtenus.

Plan $[F_1, F_2]$

Les variables ayant chacune 4 modalités, la coordonnée d'une variable V sur un facteur est égale à $\cos^2 \theta / \sqrt{3}$. L'angle θ est celui défini au § V et $\cos^2 \theta$ est égal au pourcentage d'inertie inter-classe.

Chaque coordonnée est donc comprise entre 0 et $1/\sqrt{3}$. Nous avons fait figurer, sur le plan $[F_1, F_2]$ le carré de côté $1/\sqrt{3}$.

Quand une variable est proche d'un côté extérieur de ce carré, $\cos^2 \theta$ est proche de 1. Cette variable est liée au facteur : les éléments de chacune de ses classes sont bien regroupés, le facteur est proche du sous-espace engendré par ses fonctions indicatrices, et le pourcentage d'inertie que leur extrait ce facteur est presque 1/3.



C'est le cas pour le premier facteur pour les variables (I, B, H, G, F, K, C) pour lesquels $\cos^2 \theta$ est supérieur à 0,78. C'est aussi le cas, à un degré moindre pour les variables S, T, A, ($\cos^2 \theta \cong 0.71$).

Pour le deuxième facteur ceci est encore vérifié par le groupe de variables (J, K, H, B, G).

Par contre, les variables L, M et N sont très proches de l'origine. Ce qui met en évidence leur indépendance avec les 2 premiers facteurs. Il sera donc inutile d'étudier la position de leurs modalités sur les graphiques d'Analyse des correspondances. D'autre part, S, T, A, P sont un peu liées au premier facteur, mais absolument pas au second.

Sur ce graphique, on lit aussi très bien l'importance des variables dans la détermination des facteurs. En effet, les coordonnées des variables sur un facteur sont aussi proportionnelles au pourcentage d'inertie fourni par les variables à ce facteur. Les 2 premiers facteurs sont donc essentiellement déterminés par le groupe de variables situées en haut à droite.

Certaines variables sont très proches sur le plan (1 - 2). Mais cette proximité n'est pas suffisante pour conclure sur leur distance réelle. En effet, le plan des 2 premiers facteurs ne peut suffir pour bien représenter des variables qui ont 4 modalités, et qui sont donc associées à des opérateurs de rang 3. Pour cela, il faudra au moins l'espace des 3 premiers facteurs.

Pour conclure à leur sujet, nous considérerons donc les 2 graphiques des plans $[F_1, F_2]$ et des plans $[F_3, F_4]$. Pour faciliter la lecture, nous avons tracé dans ces deux plans, le cercle de rayon $1/\sqrt{2}$: une variable sera bien représentée dans le sous-espace des 4 premiers facteurs si ses projections dans chacun des deux plans est proche de ce cercle.

Plan $[F_3, F_4]$

Le graphique montre immédiatement que les variables importantes pour le 3^{ème} facteur sont les mêmes que celles qui déterminent les 2 premiers.

Et que les variables L, M, N apparaissent au niveau du 4^{ème} facteur, et le déterminent quasiment à elles seules. Les variables A et T ne sont pas complètement indépendantes de ce 4^{ème} facteur.

D'autre part, on peut maintenant, à l'aide de ces 2 graphiques, affirmer que les variables G, B, J, H qui sont très bien représentées dans $[F_1, F_2, F_3]$ sont réellement très proches. Et que le sous-espace de R^I engendré par leurs 4 modalités est pratiquement le même et très proche du sous-espace $[F_1, F_2, F_3]$.

En conclusion, un simple coup d'oeil sur ces graphiques permet d'orienter le dépouillement de l'analyse. L'interprétation des 3 premiers facteurs devra s'appuyer sur les groupes de variables (G, B, J, H) et aussi (K, I, C, F). Le premier groupe est si homogène qu'il est vraisemblablement inutile d'en différencier les éléments. Il n'est pas étonnant de retrouver sur les graphiques de l'analyse un effet Guttman dans les plans 1 - 2 et 1 - 3 pour les variables qui sont concernées.

Par contre, le 4^{ème} facteur est indépendant des 3 premiers, on doit donc l'interpréter séparément en se basant sur les variables L, M et N.

Vraisemblablement, une étude poussée des graphiques aurait permis d'aboutir aux mêmes conclusions, mais après un long travail.

BIBLIOGRAPHIE

- [1] BENZECRI et collaborateurs. – “L’analyse des données” – Dunod, 1973.
- [2] CAILLIEZ F. PAGES J.P. – “Introduction à l’analyse des données”. SMASH, 1976,
- [3] LEBART L. et FENELON J.P. – “Statistique et informatique appliquées”. Dunod, 1973.
- [4] LEBART L. MORINEAU A., TABARD N. – “Technique de la description statistique”.
- [5] PAGES J.P., ESCOUFIER Y. CAZES P. – “Opérateur et analyse des tableaux de plus de deux dimensions”. *Cahier du B.U.R.O.*, 1976.