

# Indicateur informationnel agrégé pour évaluer la qualité de l'analyse des données sans avoir manque d'information

*Sabina Popescu – Spineni*

Institut de Santé Publique & Université de Médecine  
et Pharmacie "Carol Davila" - Bukarest, Roumanie.

*analyse factorielle, statistique informationnelle, santé publique*

[spopescu@ispb.ro](mailto:spopescu@ispb.ro)

## 1. Introduction

Dans les sciences médico-sociales, aussi que dans le management sanitaire, l'analyse multivariée des données a été imposée par des besoins opérationnels. Mais l'analyse des données utilise beaucoup de méthodes d'optimisation qui proposent des algorithmes très rigoureux pour établir la partition d'un set  $n$  d'objets caractérisés par  $k$  variables, qui caractérise une population ou un group de référence. L'analyse multidimensionnelle des données (MDA) peut inclure deux principales méthodes d'analyse: l'analyse cluster (AC) et l'analyse linéaire des données (composantes principales, canonique, factorielle, analyse des correspondances), simples ou multiples. On se propose une représentation synthétique avec le minimum de manque d'information.

Le but de l'analyse factorielle comme technique d'analyse multivariée est de déterminer si la corrélation d'un grand nombre de variables observées peut être expliquée par un petit nombre de facteurs fondamentales et combien de tels facteurs on a besoin. Des méthodes classiques d'optimisation sont adaptées pour suivre le type des données. [ 1]

Les solutions des méthodes mentionnées sont appropriées pour déterminer un nombre minimal de facteurs capables d'expliquer toute la variabilité, par une symétrie parfaite entre les profils-colonnes et les profils-lignes analysées, en considérant les contraintes et en envisageant les priorités. Le concept de *distance* joue un rôle essentiel dans la plupart des modèles d'analyse multivariée, en considérant les différentes directions de dispersion de la variabilité, mais la théorie des *valeurs propres* prends la plus grande importance dans l'analyse globale.

Le principal but de l'analyse factorielle est de réduire les dimensions de la représentation (duale) spatiale des observations, en réduisant le nombre d'axes factorielles. Mais il est très difficile d'établir un critère adéquate pour mesurer la qualité de la représentation dans l'espace factoriel réduit, après avoir réduit le nombre des axes factorielles. [ 5]

On a beaucoup de techniques empiriques pour réduire le nombre des axes factorielles (Kaiser, Cattell), mais dans cet étude on présente un bon critère pour trouver et tester un indicateur informationnel agrégé (Popescu-Spineni, 1998), celui-ci calculé par l'aide d'un concept statistique informationnel, sans avoir manque d'information. Sont bien proposés beaucoup d'exemples, pour appliquer et tester cet indicateur informationnel agrégé, avec des références théoriques appropriées.

## 2. L'analyse factorielle: la réduction du nombre des facteurs

La réduction de la dimension de l'espace (cartésien) des observations dans l'analyse factorielle ou des composantes principales revient à la réduction du nombre des axes factorielles, jusqu'à présent, sans une solution rigoureuse.

Le but de l'analyse étant d'obtenir une représentation des observations dans l'espace de dimension réduite (soit  $< p$ ), on doit apprécier la perte d'information après avoir calculé le nombre des facteurs retenus. On observe dans ce cas que le nuage des points n'est pas encore centré dans le centre de gravité (G). Après avoir calculé les valeurs propres ( $\lambda_i$ , avec  $1 \leq i \leq p$ ) du système des  $k$  équations, on considère le critère du pourcentage de l'inertie totale ( $\tau$ ) comme mesure pour apprécier la qualité de la représentation factorielle (Benzécri, 1979), par :

$$\tau = \frac{\lambda_1 + \lambda_2 + \dots + \lambda_p}{\lambda_1 + \lambda_2 + \dots + \lambda_K} = \frac{\lambda_1 + \lambda_2 + \dots + \lambda_p}{\phi^2}.$$

Par exemple, avec  $(\lambda_1 + \lambda_2)/(\phi^2) = 0,9$ , donc 90%, on résulte que le nuage des points est applatisé sur un sous-espace avec deux dimensions (Saporta, 1996). On peut observer que, dans l'appréciation du pourcentage de l'inertie, on doit tenir compte du nombre des variables initiales: un pourcentage de 10% n'a pas la même signification pour un tableau avec 20 variables comme pour un autre de 100 variables. [ 1]

Ont été élaborés des critères théoriques (Saporta, 1996) qui proposent déterminer si les valeurs propres sont significativement différentes, en commençant d'un certain rang, mais si non, on retient les premières valeurs propres ( $\lambda_i$ ). La réduction de dimension n'est pas possible seulement s'il existe une redondance entre les variables  $x_1, \dots, x_p$ , ou non, s'il sont indépendantes (Saporta). Par la méthode des approximations successives, on estime les valeurs propres et les vecteurs propres associés à l'équation caractéristique  $|R - \lambda I| = 0$ , qui explique la pluspartie de la variation totale. En supposant que le processus c'est interrompu après avoir estimé  $p$  ( $p < K$ ) valeurs propres, on doit vérifier si les plus petites des valeurs propres sont égales, pour maximiser la variance des composantes principales de l'analyse. [ 1]

En réalité, on s'applique des critères empiriques (Kaiser, 1965), étant le cas par exemple, de l'analyse des composantes principales pour les valeurs propres supérieures à 1.

En général, pour résoudre l'équation statistique matricielle  $Y = AX + \varepsilon$ , dans l'analyse factorielle, on doit proposer un critère d'optimisation et la méthode d'estimation des paramètres. Pour l'estimation optimale des paramètres, il existent deux conditions principales: la reproduction des corrélations observées et l'explication de la variabilité (Torrens-Ibern, 1972). Les méthodes statistiques pour l'estimation des paramètres dépendent de la technique de résoudre le système d'équation linéaires: (1) la méthode des moindres carrés; (2) le principe de la vraisemblance maximale (maximum likelihood); - les deux étant équivalentes dans le cas normal (Gauss-Laplace).

En général, dans l'analyse factorielle, le problème de tester le nombre des facteurs (ou des axes) qu'on doit retenir se pose de deux directions: soit on suppose l'existence d'un nombre  $p$  de facteurs communs en vérifiant l'hypothèse nulle  $H_0$ , soit on détermine le nombre intégral des facteurs, en calculant leurs nombres en partant du nombre des observations. [ 2]

Les plus importantes critères théoriques proviennent de Bartlett (1951) et Lawley (1940), mais Jöreskog (1963) propose une procédure spécifique pour réduire le nombre des facteurs dans le cas des observations non-normales ("sans hypothèses préalables") [1]

### 3. L'approximation de la distribution des valeurs propres

La réduction du nombre des facteurs dans l'analyse factorielle ou dans l'analyse des correspondances (Benzécri, 1980) dépend de l'analyse de la distribution théorique des valeurs propres ( $\lambda_i$ ). Dans un tableau de contingence, sous l'hypothèse de l'indépendance des lignes et des colonnes, cette distribution s'approche d'une loi de distribution connue d'une matrice Wishart (Lebart, Morineau, Piron, 1995). On se souligne que la loi de distribution des valeurs propres résultées de l'analyse a été erronée appréciée, aussi que l'inertie totale du nuage des points, comme de type  $\chi^2$ , pendant que différentes simulations ont démontré au contraire, en restant un problème théorique encore. L'utilisation du taux de l'inertie (pourcentage de la variance) pour la "qualité de la représentation", proposé de Benzécri (1979) est très difficile d'être globalisée. Enfin, pour mesurer la variation de l'information, Kullback (1959) a proposé la théorie de Shannon-Wiener pour utiliser un indice de divergence dans les problèmes d'inférence pour la statistique multivariée (Lebart et autres, 1995).

Basée sur le théorème de Bayes, on mesure la distance  $J(H_1:H_2)$  entre les hypothèses  $H_1$  versus  $H_2$  (alternatives), en tenant compte de la matrice de covariance ( $\Gamma$ ) et des valeurs propres ( $\lambda_i$ ) :

$$J(I, \Gamma) = \frac{1}{2} \left( \sum_{i=1}^p \lambda_i + \sum_{i=1}^p \frac{1}{\lambda_i} \right) - p > 0,$$

où  $I(1: 2)$  est l'information moyenne après la discrimination dans les deux espaces des deux échantillons, sous les deux hypothèses contraires (Lebart et autres, 1995).

Si les deux inerties théoriques totales sont égales sous les deux hypothèses  $H_1$  et  $H_2$ , on a :

$$\sum_{i=1}^p \lambda_i = p$$

Donc, le seul terme qui peut avoir une variation est:  $\sum_{i=1}^p \frac{1}{\lambda_i}$ ,

fait qui pose en question les petites valeurs propres, en temps que l'analyse factorielle ne retient en général que les plus grandes (Lebart et autres, 1995). D'après Lebart et autres, on peut aussi construire, basé sur la théorie de Kullback (1959), un domaine critique:

$$\left\{ \sum_{\alpha=1}^p \frac{1}{\lambda_{\alpha}} - \sum_{\alpha=1}^p \lambda_{\alpha} \right\} \geq C.$$

On peut observer que la divergence entre les deux hypothèses augmente dans le cas où les valeurs propres de la matrice  $\Gamma$  se rapprochent de zéro. Conformément au théorie de l'information, les valeurs propres infiniment petites auront un impact plus grand par rapport aux celles qui peuvent expliquer près de 80% de l'inertie totale du sous-espace des deux facteurs mis en correspondance, avec une perte considérable d'information (Lebart et autres, 1995). [ 2 ]

#### 4. Indicateur informationnel agrégé :

En général, dans l'analyse factorielle:  $\tau_j = \lambda_j / \sum_{j=1}^p \lambda_j$ , avec:  $\sum_{j=1}^p \tau_j = 1$ ,

représente l'indicateur de la "qualité de la représentation" (Benzécri, 1979). En tenant compte de ces considérations, aussi que d'autres théories informationnelles (de Kullback, 1959 et de Onicescu, 1963), pour une évaluation globale de l'analyse factorielle, je propose dans le présent travail un indicateur informationnel agrégé (IS-Popescu-Spineni, 1998), d'omogénéité-hétérogénéité, que moi même j'ai construit, en partant des deux rapports (informationnels), construites dans ce but: [ 3 ] [ 4 ]

$$(1) \quad E(\tau) = \sum_{j=1}^K \frac{\lambda_j^2}{\left(\sum_{j=1}^K \lambda_j\right)^2}, \quad \text{et} \quad (2) \quad E(\tau') = \sum_{j=1}^K \frac{\frac{1}{\lambda_j^2}}{\left(\sum_{j=1}^K \frac{1}{\lambda_j}\right)^2},$$

Ces rapports m'ont conduit à élaborer l'ISPS, un "indicateur synergique" (ou domaine critique), étant en mesure d'évaluer la signification globale de la représentation, par comparaison avec le nombre des variables retenues ( $1/k$ , avec  $k < p$ ):

$$(ISPS): \quad \{|E(\tau) - E(\tau')|\} \geq \frac{1}{K}, \quad (3)$$

L'indicateur  $E(\tau)$ , est mis en balance (dans ISPS) avec l'indicateur,  $E(\tau')$ , en tenant compte qu'ils ont une variation informationnelle, stricte entre  $1/K$  et 1, car la constante  $C = 1/K$  est précisément établie, sans avoir besoin d'une simulation. *Interprétation*, sans avoir manque d'information: dans le cas où la

valeur de  $ISPS > 1/K$ , on a une signification statistique de l'analyse ( $H_0$  acceptée), au contraire, pour  $ISPS \leq 1/K$ , l'analyse globale n'est pas significative ( $H_0$  rejetée). [ 3 ] [ 4 ]

## 5. Applications

**5.1.** Pour le commencement, je donnerai un exemple d'utilisation comparative des deux indicateurs informationnels  $E(\tau_j)$  et  $E(\tau'_j)$ , aussi que de *ISPS* (l'indicateur synergique) sur les valeurs propres sorties d'une application effectuée de A. Rizzi ("Analisi dei Dati", Roma, 1989, p.188), sur un tableau "4x10":

Exemple No.1:

Ex.1	$\lambda_j$	$\tau_j$	cum.	cum.	$1/\lambda_j$	$\tau'_j$
1	0,0443	0,532		1,000	22,573	0,004
2	0,0202	0,242	0,532	0,996	49,506	0,008
3	0,0088	0,106	0,772	0,998	113,636	0,018
4	0,0054	0,065	0,887	0,970	185,785	0,029
5	0,0023	0,028	0,942	0,941	434,783	0,069
6	0,0021	0,025	0,970	0,872	476,191	0,076
7	0,0002	0,003	0,997	0,796	5000,000	0,796
	0,0833	1,000	1,000	—	6282,474	1,000

En calculant:

$$(1) E(\tau_j) = \sum_{j=1}^K \frac{\lambda_j^2}{(\sum_{j=1}^K \lambda_j)^2}; \text{ et } (2) E(\tau'_j) = \sum_{j=1}^K \frac{1/\lambda_j^2}{(\sum_{j=1}^K 1/\lambda_j)^2};$$

$$\text{En calculant aussi l'indicateur } ISPS: \quad \{|E(\tau) - E(\tau')|\} \geq \frac{1}{K}, \quad (3)$$

on obtient un résultat global significatif:

$$E(\tau_j) = 0,357; E(\tau'_j) = 0,645; \quad \text{où: } k = 7;$$

$$\{|E(\tau) - E(\tau')|\} = 0,645 - 0,357 = 0,288 > 0,143 = \frac{1}{7}.$$

En conclusion, il faut rejeter  $H_0$ , car le tableau n'est pas équilibré. [ 3 ] [ 4 ]

**5.2.** Dans un autre recherche, j'ai fait un étude sur la situation de la répartition par districts (41 districts), pour 15 spécialités médicales, des médecins de la Roumanie de l'an 1985 ((a)-données en chiffres absolues, (b)-donnees rapportées pour 10000 habitants), en utilisant l'analyse des correspondances, aussi que l'analyse cluster, en parallèle, sur un tableau de contingence de "41x15". A l'aide de l'analyse des correspondances, faite sur la répartition des médecins par districts, pour les 15 spécialités, dans le cas (b)- données rapportées par 10000 habitants, j'ai obtenu une représentation duale sur les deux axes factorielles F1 (avec  $\lambda_1 = 0,01449356$ ) et F2 (avec  $\lambda_2 = 0,00876388$ ), en caractérisant la variabilité totale avec 31,8% et respectivement 19,2%, ayant une qualité de la représentation de  $\tau_j = 51,0\%$ , où  $K = 15$  ( $j = 1, \dots, K$ ). [ 3 ] [ 4 ]

La représentation duale a montré une agglomération des «objets» (districts) et des «variables» (spécialités) autour l'origine des axes factorielles, plus accentuée dans l'analyse (b) que dans l'analyse (a), due au rapport des spécialistes par habitants ( $^0/_{0000}$ ). Autour du centre des axes, l'inertie du nuage est plus faible sans avoir des directions privilégiées. Dans l'analyse (a) a été observée une forte corrélation des districts avec les centres universitaires sur la première axe F1, qui ne se mentient que

partiellement dans l'analyse (b), où la distribution des districts et des spécialités est plus équilibrée autour du centre des axes. Cette situation plus équilibrée montrée par l'analyse duale dans le cas (b) se vérifie avec les indicateurs informationnels décrits plus haut, surtout avec *l'indicateur synergique ISPS*, appliqués sur les valeurs propres (cas (b)- pour  $0/0000$ ), sorties de l'analyse d'un tableau "41x15" :

Exemple No. 2 :

	A1 (a) $\lambda_j$	inertie(%)	A2 (b) $\lambda_j$	inertie (%)
$\lambda_1$	0,0154	42,4	0,0145	31,9
$\lambda_2$	0,0046	12,8	0,0088	19,3
$\lambda_3$	0,0042	11,6	0,0067	14,7
$\lambda_4$	0,0027	7,5	0,0037	8,2
$\lambda_5$	0,0023	6,4	0,0029	6,4
$\lambda_6$	0,0017	4,6	0,0019	4,2
$\lambda_7$	0,0012	3,4	0,0017	3,7
$\lambda_8$	0,0011	3,1	0,0014	3,2
$\lambda_9$	0,0009	2,4	0,0011	2,4
$\lambda_{10}$	0,0007	1,9	0,0009	2,0
$\lambda_{11}$	0,0005	1,4	0,0007	1,5
$\lambda_{12}$	0,0004	1,1	0,0005	1,1
$\lambda_{13}$	0,0003	0,9	0,0004	0,8
$\lambda_{14}$	0,0002	0,5	0,0003	0,6
$\lambda_{15}$	0,0000	0,5	0,0000	0,0
	$\varphi = 0,03619$	100.	$\varphi = 0,04551$	100.

(b) En calculant les indicateurs (1)  $E(\tau_j)$  et (2)  $E(\tau'_j)$ ,  
on obtient:  $E(\tau_j) = 0,179$  et  $E(\tau'_j) = 0,136$ ; où:  $k = 14$ ;  
En calculant aussi l'indicateur *ISPS*, on obtient un résultat non-significatif:  
 $\{|E(\tau) - E(\tau')|\} = 0,179 - 0,136 = 0,043 < 0,071 = 1/14$ .

Il faut accepter  $H_0$ , car l'analyse n'a pas montré des directions privilégiées, pour le cas des médecins spécialistes rapportés par habitants ( $0/0000$ ). [ 3] [ 4]

#### Références:

1. Benzécri, J.P. – *L'Analyse des données*, Dunod, Paris, 1980;
2. Lebart, L., Morineau, A., Piron, M. – *Statistique exploratoire multidimensionnelle*, Dunod, Paris, 1995;
3. Popescu-Spineni Sabina – *Statistical Methods for Multidimensional Data Analysis and Classification*; Editura Universitara "Carol Davila", 2000, (145 pages);
4. Popescu-Spineni Sabina - *Hierarchy Techniques of Multidimensional Data Analysis (MDA) in Social Medicine Research*; Studies in Classification, Data Analysis and Knowledge Organization, Springer Verlag (A. Rizzi, M. Vichi, H.H. Bock Editors), 641-646, 1998;
5. Rizzi, A., ed. – *Some Relation between Matrices and Structures of Multidimensional Data Analysis*; Giardini Editori e Stampatori, Pisa, 1995.

