

STATISTIQUE ET ANALYSE DES DONNÉES

PABLO GONZALEZ VICENTE

Recherche de la dimension de l'espace latent en ACP

Statistique et analyse des données, tome 11, n° 3 (1986), p. 19-29

http://www.numdam.org/item?id=SAD_1986__11_3_19_0

© Association pour la statistique et ses utilisations, 1986, tous droits réservés.

L'accès aux archives de la revue « Statistique et analyse des données » 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/>

RECHERCHE DE LA DIMENSION DE
L'ESPACE LATENT EN ACP

GONZALEZ VICENTE Pablo

Electricité de France
Direction des Etudes et Recherches

Résumé : On étudie une condition nécessaire et une condition suffisante de positivité pour des matrices symétriques et on donne une interprétation géométrique de ces résultats. On propose ensuite une interprétation de l'ACP et un critère simple pour la recherche de la dimension de l'"espace latent".

Abstract : In this paper we give one sufficient condition and one necessary for symmetric matrices to be positive. We also give a geometrical interpretation to these results. Finally we suggest a new way of considering Principal Component Analysis and a simple criterion for determining the latent space dimension.

Indices de classification STMA : 06-070, 06-080.

1 - GENERALITES

Soit $S(n)$ l'espace vectoriel réel des matrices symétriques $n \times n$ à coefficients réels. Nous fournissons cet espace du produit interne :

$$\langle A, B \rangle = \text{Trace}(AB')$$

$$A, B \in S(n)$$

et de la norme associée :

$$|A| = (\text{Trace}(AA'))^{\frac{1}{2}}$$

Manuscrit reçu le 22 juin 1986

Révisé le 1 juin 1987

Si τ dénote le sous-ensemble de $S(n)$ des matrices symétriques et définies (semi)-positives, alors :

- * τ est un cône convexe fermé avec sommet à l'origine,
- * τ° est l'ensemble des matrices strictement définies positives,
- * Frontière (τ) est l'ensemble des matrices semi-définies positives de rang non complet.

Rappelons finalement que pour toute matrice $A \in S(n)$ existe P orthogonale tel que : $P'AP = D$ D matrice diagonale $PP' = P'P = I$

et, en plus : $|A| = |D|$, $\text{trace}(A) = \text{trace}(D)$.

2 - CONDITIONS DE SEMI-POSITIVITE

Théorème 1 :

Soit $A \in S(n)$
 Si $\frac{\text{Trace}(A)}{|A|} \geq \sqrt{n-1}$ alors A est semi-définie positive
 L'inégalité stricte implique la stricte positivité de A

Démonstration :

Soit $A \in S(n)$ tel que : $\frac{\text{Trace}(A)}{|A|} \geq \sqrt{n-1}$

Si A n'est pas semi-définie positive alors A a une valeur propre négative ; soit celle-ci λ_1 .

Soit $D^* = \begin{bmatrix} 0_1 & & \\ & \dots & \\ & & 1 \end{bmatrix}$ et $A = PDP'$, $D = \begin{bmatrix} \lambda_1 & & \\ & \dots & \\ & & \lambda_n \end{bmatrix}$ avec $\lambda_1 < 0$.

Alors : $\langle I - D^*, D \rangle = \lambda_1 < 0$
 $\Rightarrow \langle I, D \rangle < \langle D^*, D \rangle \leq |D^*| |D|$

$$\text{Donc : } \frac{\text{Trace}(A)}{|A|} = \frac{\text{Trace}(D)}{|D|} < |D^*| = \sqrt{n-1}$$

ce qui contredit l'hypothèse du départ.

Théorème 2 :

Soit $A \in S(n)$

Si A est semi-définie positive alors $\frac{\text{Trace}(A)}{|A|} \geq 1$

Démonstration :

Soit A semi-définie positive ; alors :

$$A = PDP^t, \quad D = \begin{bmatrix} \lambda_1 & & \\ & \dots & \\ & & \lambda_n \end{bmatrix} \quad \text{avec } \lambda_i \geq 0 \quad \forall_i = 1, \dots, n$$

Soit $\{D_i\}_{1 \leq i \leq n}$ la base canonique des matrices diagonales de dim.n

$$\text{Alors } |A| = |D| = \left| \sum_{i=1}^n \lambda_i D_i \right| \leq \sum_{i=1}^n \lambda_i |D_i| = \sum_{i=1}^n \lambda_i = \text{Trace}(D) = \text{Trace}(A).$$

$$\text{Donc : } \frac{\text{Trace}(A)}{|A|} \geq 1.$$

Corolaire :

Soit $A \in S(n)$; alors :

$$* \text{ Si } n=2, A \text{ est semi-définie positive } \Leftrightarrow \frac{\text{Trace}(A)}{|A|} \geq 1$$

$$* \text{ Si } \text{rang}(A) = K \text{ alors } A \text{ semi-définie positive } \Rightarrow \frac{\text{Trace}(A)}{|A|} \geq \sqrt{K}$$

3 - INTERPRETATION GEOMETRIQUE

Remarquons que :

$$\frac{\text{Trace}(A)}{|A|} \geq \sqrt{n-1} \Leftrightarrow \frac{\langle A, I \rangle}{\sqrt{n} |A|} \geq \frac{\sqrt{n-1}}{\sqrt{n}} \Leftrightarrow \frac{\langle A, I \rangle}{|A| |I|} \geq \frac{\sqrt{n-1}}{\sqrt{n-1} \sqrt{n}}$$

Notons $E_1 = I - D_1 = \begin{bmatrix} 0 & & \\ & \dots & \\ & & 1 \end{bmatrix}$ alors :

$$\frac{\text{Trace}(A)}{|A|} \geq \sqrt{n-1} \Leftrightarrow \frac{\langle A, I \rangle}{|A| |I|} \geq \frac{\langle E_1, I \rangle}{|E_1| |I|} \Leftrightarrow \text{angle}(A, I) \leq \text{angle}(E_1, I).$$

Le théorème 1 se réécrit donc :

$A \in S(n)$ et $\text{angle}(A, I) \leq \text{angle}(E_1, I) \Rightarrow A$ semi-définie positive.

D'autre part :

$$\frac{\text{Trace}(A)}{|A|} \geq 1 \Leftrightarrow \frac{\langle A, I \rangle}{|A| |I|} \geq \frac{1}{|I|} = \frac{|D_1|}{|D_1| |I|} = \frac{\langle D_1, I \rangle}{|D_1| |I|}$$

Le théorème 2 se réécrit donc :

$A \in S(n)$ et A semi-définie positive $\Rightarrow \text{angle}(A, I) \leq \text{angle}(D_1, I)$.

Nous pouvons conclure que le cône \mathcal{C} des matrices semi-définies positives est contenu dans un cône de révolution maximal, dont la génératrice est la matrice D_1 et l'axe de rotation la matrice identité I .

Evidemment toutes les matrices dans ce cône-là ne sont pas semi-définies positives, mais il existe un cône de révolution minimal, dont la génératrice est la matrice $E_1 = 1 - D_1$ et l'axe de rotation est l'identité I , tel que toutes les matrices contenues dans lui sont semi-définies positives.

Dans la figure 1 nous donnons cette interprétation pour les matrices diagonales de 3×3 :

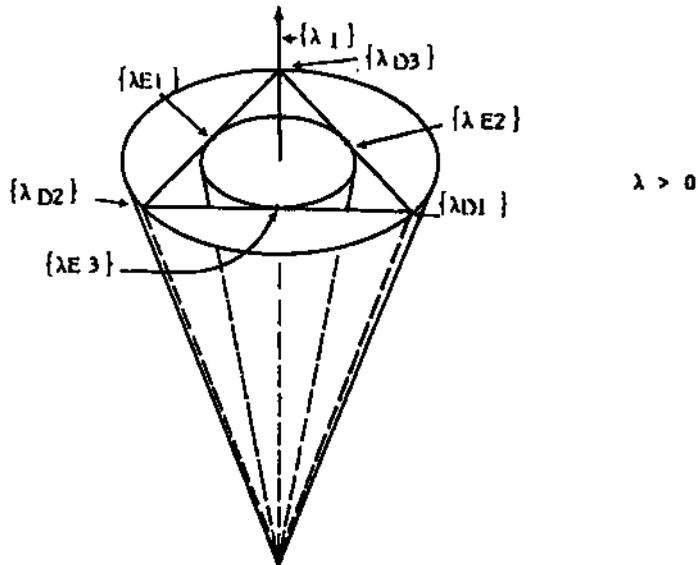


Figure 1

Une interprétation géométrique très intéressante peut être donnée dans $S(2)$ grâce à l'identification :

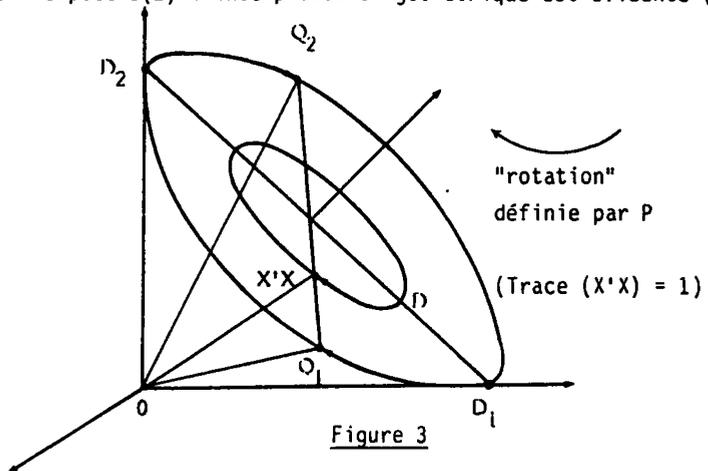
$$\begin{bmatrix} X & Z \\ Z & Y \end{bmatrix} \longleftrightarrow \begin{bmatrix} X \\ Y \\ Z \end{bmatrix} \in \mathbb{R}^3$$

Dans ce cas particulier \mathcal{B} est exactement le cône de révolution :

$$\{ X \in S(2) \text{ tel que } \text{angle}(X, I) \leq \text{angle}(D_1, I) = 45^\circ \}$$

La base $\{Q_i\} \ 1 \leq i \leq q$ est une base orthonormée et chaque Q_i peut être écrit en fonction des vecteurs propres P_i de $X'X$: $Q_i = P_i P'_i$.

Dans l'espace $S(2)$ l'interprétation géométrique est évidente (figure 3) :



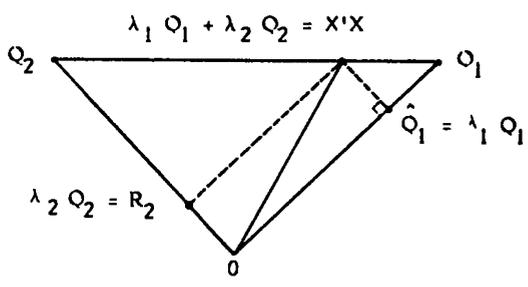
Retenir les n premiers axes dans l'ACP revient à projeter la matrice $X'X$ sur le sous espace vectoriel engendré par les matrices Q_i associés aux n plus grandes valeurs propres λ_i . Etant donné que $\{Q_i\} \ 1 \leq i \leq q$ est une base orthonormée, cette projection est :

$$\hat{Q}_n = \sum_{i=1}^n \lambda_i Q_i \quad \text{rang}(\hat{Q}_n) = n$$

et le "résidu" R_n :

$$R_n = Q - \hat{Q}_n = \sum_{i=n+1}^q \lambda_i Q_i \quad \text{avec} \ ||Q - \hat{Q}_n||^2 = \sum_{i=n+1}^q \lambda_i^2 \text{ et } \hat{Q}_n \perp Q - \hat{Q}_n$$

Reprenant la fig. 3 nous avons :



Le problème qui se pose tout naturellement est le choix optimal de n dans la décomposition $X'X = \hat{Q}_n + R_n$.

Si l'on interprète \hat{Q}_n comme la matrice de variance expliquée par les premiers n facteurs, il est naturel de supposer que la matrice de variance de résidus soit de la forme $R_n = \sigma^2 I$. Evidemment cette condition sera difficilement réalisée dans la pratique. Mais lorsque l'on utilise l'ACP comme approximation du modèle factoriel, on peut chercher n , de façon que cette condition soit le mieux satisfaite.

La proximité entre R_n et $\sigma^2 I$ est mesurée par le cosinus de l'angle θ entre ces deux matrices :

$$\cos \theta (R_n, \sigma^2 I) = \frac{\langle R_n, I \rangle}{|R_n| |I|} = \frac{\text{Trace}(R_n)}{q |R_n|}$$

Le critère sera donc :

$$\text{Max}_n \frac{\sum_{i=n+1}^q \lambda_i}{\left(\sum_{i=n+1}^q \lambda_i^2\right)^{\frac{1}{2}}} \Leftrightarrow \text{Min}_n \theta (R_n, I)$$

5 - APPLICATION

Nous reprenons ici les exemples donnés par ANASTASSAKOS et D'AUBIGNY (1) ; nous arrivons au mêmes résultats avec un critère plus simple que le test de sphéricité empirique (VELICER) pour eux exemplifié.

Cas 1 :

On considère le modèle :

$$Y_j = \sum_{i=1}^3 a_{ij} X_i = u_j \quad j = 1, \dots, 6$$

Les variables X_i ($1 \leq i \leq 3$) ont été générées de 3 échantillons gaussiens indépendants ($N = 100$, $m = 0$, $\sigma = 1$) et les u_j ont été générés de 6 échantillons indépendants (entre eux et des échantillons précédents) gaussiens ($m = 0$, $\sigma = 0.01$). Les coefficients sont :

$$\begin{aligned} y_1 &= .5X_1 + .3X_2 + .2X_3 \\ y_2 &= .7X_1 + .2X_2 + .1X_3 \\ y_3 &= .9X_1 + .8X_2 + .02X_3 \\ y_4 &= .3X_1 + .1X_2 + .6X_3 \\ y_5 &= .01X_1 + .8X_2 + .19X_3 \\ y_6 &= .2X_1 + .35X_2 + .45X_3 \end{aligned}$$

Matrice de
correlations =

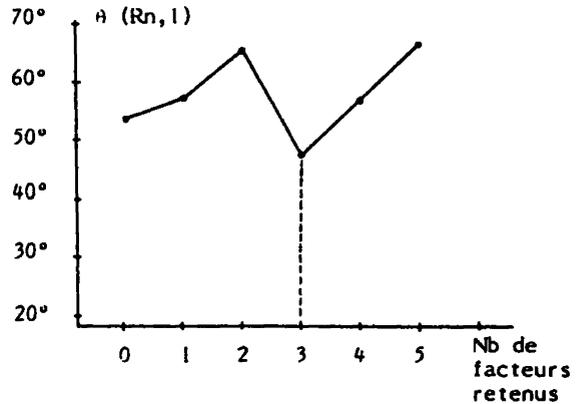
$$\begin{bmatrix} 1 & & & & & \\ .928 & 1 & & & & \\ .821 & .964 & 1 & & & \\ .677 & .532 & .416 & 1 & & \\ .448 & .150 & -.066 & .215 & 1 & \\ .752 & .502 & .309 & .872 & .645 & 1 \end{bmatrix}$$

L'ACP fournit les valeurs propres :

$$\begin{aligned} \lambda_1 &= 3.874 & \lambda_2 &= 1.425 & \lambda_3 &= 0.675 \\ \lambda_4 &= 0.012 & \lambda_5 &= 0.009 & \lambda_6 &= 0.006 \end{aligned}$$

L'angle entre $R_n = X'X\hat{Q}_n$ et I :

Nb de facteurs retenus	0	1	2	3	4	5
$\theta(R_n, I)$	54°	57°	65°	47°	56°	66°



Cas 2 :

Nous étudions maintenant une situation d'équicorrélation :

$$\text{Matrice de correlations} = \begin{bmatrix} 1 & P & \dots & P \\ P & 1 & \dots & P \\ \vdots & \vdots & \ddots & \vdots \\ P & P & \dots & 1 \end{bmatrix} \quad 0 < P \leq 1$$

L'ACP de cette matrice fournit les valeurs propres :

$$\lambda_1 = 1 + (q-1)P \quad \lambda_2 = \dots = \lambda_q = 1 - P$$

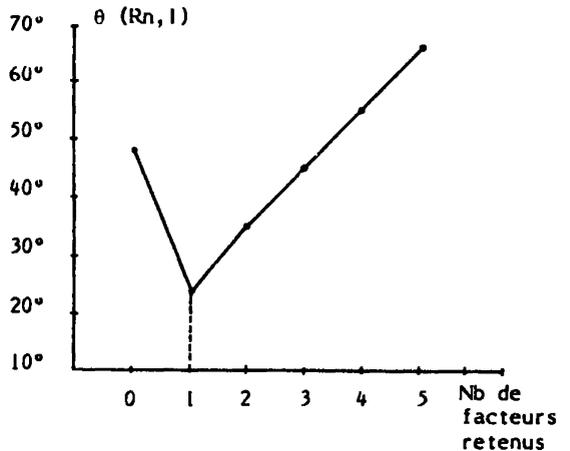
Calcul de $\theta (R_n, I)$:

$$n=0 \Rightarrow R_n = X'X \cos \theta = \frac{q}{\sqrt{q} \sqrt{q(q-1) p^2 + q}}$$

$$n>0 \Rightarrow R_n = X'X - \hat{Q}_n \cos \theta = \sqrt{\frac{q-n}{q}}$$

Dans le graphique nous considérons
 $q = 6$

Nb de facteurs retenus :					
0	1	2	3	4	5
$\theta (R_n, I) : 18^\circ \ 24^\circ \ 35^\circ \ 45^\circ \ 55^\circ \ 66^\circ$					
pour $P=0,5$					



Cas 3 :

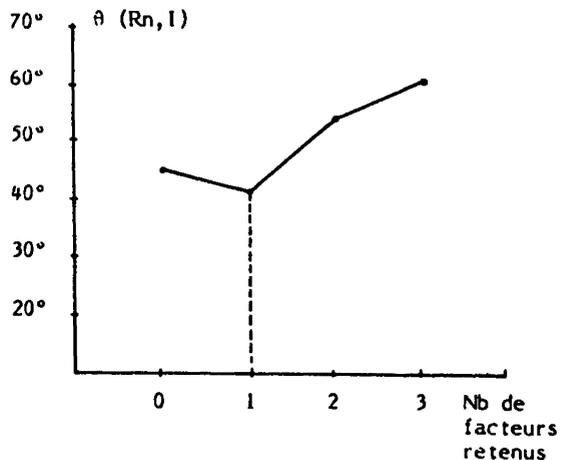
La matrice de corrélations $\begin{bmatrix} 1 & 0.7 & 0.6 & 0.4 \\ & 1 & 0.4 & 0.6 \\ & & 1 & 0.7 \\ & & & 1 \end{bmatrix}$ correspond à une

situation où les coefficients de corrélation multiple de chacune des q variables, sachant les $q-1$ autres, sont tous égaux ; les vecteurs propres fournissent les q contrastes orthogonaux du plan d'expérience à 2^q facteurs.

Les valeurs propres fournies par l'ACP sont :

$$\lambda_1=2.7 \quad \lambda_2=0.7 \quad \lambda_3=0.5 \quad \lambda_4=0.1$$

Nb de facteurs retenus :			
0	1	2	3
$\theta (R_n, I) : 45^\circ \ 41^\circ \ 54^\circ \ 60^\circ$			



REFERENCES

(1) ANASTASSAKOS I., D'AUBIGNY G.

"L'utilisation des tests de sphéricité pour la recherche de la dimension de l'espace latent en analyse factorielle classique et en analyse en composantes principales",

Revue de Statistique Appliquée, 1984, Vol XXXII, n° 2, pp 45-57.