

JOURNAL DE LA SOCIÉTÉ STATISTIQUE DE PARIS

JEAN-JACQUES BOULANGER

Le problème des unités en analyse statistique multidimensionnelle appliquée à la macro-économie

Journal de la société statistique de Paris, tome 121, n° 2 (1980), p. 90-95

http://www.numdam.org/item?id=JSFS_1980__121_2_90_0

© Société de statistique de Paris, 1980, tous droits réservés.

L'accès aux archives de la revue « Journal de la société statistique de Paris » (<http://publications-sfds.math.cnrs.fr/index.php/J-SFdS>) 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/>

LE PROBLÈME DES UNITÉS EN ANALYSE STATISTIQUE MULTIDIMENSIONNELLE APPLIQUÉE A LA MACRO-ÉCONOMIE

Jean-Jacques BOULANGER

Dans l'étude qui suit, l'auteur analyse le problème posé par le choix des unités en analyse statistique multidimensionnelle appliquée à la macro-économie et suggère une solution permettant d'en amoindrir les effets sur la détermination des vecteurs propres et des composantes principales.

In the following paper, the author analyzes the problem raised by the choice of the units in multivariate statistical analysis applied to macroeconomics, and suggests a solution enabling to lessen the effects on the determination of the latent vectors and of the principal components.

In der Studie, die folgt, analysiert der Verfasser das Problem, das die multidimensionale statistische Analyse für ein gegebenes Problem stellt und zwar für die Mikroökonomie. Er schlägt eine Lösung vor, die die Wirkung auf die charakteristischen Vektoren vermindert und ebenso auf die hauptsächlichsten Komponenten.

Le choix des unités lors de l'analyse en composantes principales intrigue les statisticiens. Faut-il, par exemple, utiliser les mesures propres à un échantillon, ou bien remplacer les variables x_i par leurs valeurs réduites $\frac{x_i}{\sigma_i}$, c'est-à-dire effectuer un changement d'unité. Or on obtient des valeurs propres et des vecteurs propres dans les deux cas très différents, et il en sera de même pour tout autre changement d'unité.

Nous avons essayé d'analyser le problème en cherchant la première composante principale de la matrice des (variances et) covariances relatives aux 3 variables α_1 , α_2 et α_3 définies de la façon suivante :

α_1 = taux d'accroissement annuel en volume du produit national brut (PNB)

α_2 = taux d'accroissement annuel en volume de la formation brute de capital (FBC)

α_3 = taux d'accroissement annuel en volume des importations et revenus versés au reste du monde (Imp + etc.).

Le choix s'est porté sur ces 3 variables parmi les agrégats normalisés de l'ancienne nomenclature SEEF, parce que ce sont celles qui présentent entre elles les corrélations les plus importantes. On est ainsi parti, pour la période 1960-71, de la matrice des covariances suivantes :

$$A = \begin{bmatrix} 0,911 & 3,706 & 3,245 \\ 3,706 & 21,083 & 16,880 \\ 3,245 & 16,880 & 22,469 \end{bmatrix}$$

et l'on obtient comme première valeur propre λ_1 et premier vecteur propre u_1 les résultats suivants :

$$\lambda_1 = 39,30 \quad u_1 = \begin{bmatrix} 0,1268 \\ 0,6880 \\ 0,7146 \end{bmatrix}$$

On en déduit la première composante principale :

$$y_1 = 0,1268 (\alpha_1 - \bar{\alpha}_1) + 0,6880 (\alpha_2 - \bar{\alpha}_2) + 0,7146 (\alpha_3 - \bar{\alpha}_3)$$

qui peut être considérée comme un facteur général de croissance de l'économie.

Mais au lieu de considérer les variables centrées $x_i = \alpha_i - \bar{\alpha}_i$ définies plus haut, on peut s'intéresser aux variables centrées $x_i = k_i (\alpha_i - \bar{\alpha}_i)$, k_i étant le rapport au PNB du même PNB, de la FBC et des (Imp+ etc.), ceux-ci étant mesurés en milliards de francs par exemple, ce qui conduit à estimer les k_i par $k_1 = 1$, $k_2 = 0,27$, $k_3 = 0,15$.

Dans le premier cas on cherche à déterminer un facteur général de croissance à partir des taux de croissance $\alpha_1, \alpha_2, \alpha_3$ des 3 agrégats définis plus haut, donc à partir d'accroissements relatifs. Dans le second cas on essaie de déterminer un facteur général de croissance à partir des accroissements absolus (à un facteur près) du PNB, de la FBC et des (Imp+ etc). Ces deux choix sont tous deux extrêmement défendables. Cependant le deuxième choix consiste à majorer l'importance du PNB par rapport aux 2 autres agrégats, le premier choix consiste à la diminuer.

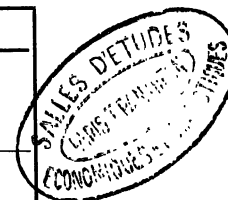
Le deuxième choix conduit à la matrice des covariances suivantes :

$$A = \begin{bmatrix} 0,9111 & 1,0006 & 0,4868 \\ 1,0006 & 1,5369 & 0,6836 \\ 0,4868 & 0,6836 & 0,5056 \end{bmatrix}$$

et l'on obtient comme première valeur propre λ_1 et premier vecteur propre u_1 les résultats suivants :

$$\lambda_1 = 2,607 \quad u_1 = \begin{bmatrix} 0,5485 \\ 0,7494 \\ 0,3708 \end{bmatrix}$$

Tableau 1	Cas 1	Cas 2	Cas 3	Cas 4
$u_1 = \begin{bmatrix} u_{11} \\ u_{12} \\ u_{13} \end{bmatrix}$	0,127 0,688 0,715	0,280 0,783 0,556	0,548 0,749 0,371	0,579 0,592 0,560
$\frac{\lambda_1}{\text{Trace}}$	0,884	0,881	0,883	0,853



On voit que les composantes du premier vecteur propre sont dans les deux cas nettement différentes et conduisent à deux axes factoriels sans rapport entre eux. Plus généralement on a analysé les 4 cas suivants caractérisés par les choix :

Cas 1	$x_i = \alpha_i - \bar{\alpha}_i$
Cas 2	$x_i = \sqrt{k_i} (\alpha_i - \bar{\alpha}_i)$
Cas 3	$x_i = k_i (\alpha_i - \bar{\alpha}_i)$
Cas 4	$x_i = \frac{\alpha_i - \bar{\alpha}_i}{\sigma_i}$ (tableau 1)

La première composante principale s'exprime donc dans les 4 cas par :

$$1^{\text{er}} \text{ cas : } y_1 = u_{11} (\alpha_1 - \bar{\alpha}_1) + u_{12} (\alpha_2 - \bar{\alpha}_2) + u_{13} (\alpha_3 - \bar{\alpha}_3)$$

$$2^{\text{e}} \text{ cas : } y_1 = u_{11} (\alpha_1 - \bar{\alpha}_1) + u_{12} \sqrt{k_2} (\alpha_2 - \bar{\alpha}_2) + u_{13} \sqrt{k_3} (\alpha_3 - \bar{\alpha}_3)$$

$$3^{\text{e}} \text{ cas : } y_1 = u_{11} (\alpha_1 - \bar{\alpha}_1) + u_{12} k_2 (\alpha_2 - \bar{\alpha}_2) + u_{13} k_3 (\alpha_3 - \bar{\alpha}_3)$$

$$4^{\text{e}} \text{ cas : } y_1 = u_{11} \frac{\alpha_1 - \bar{\alpha}_1}{\sigma_1} + u_{12} \frac{\alpha_2 - \bar{\alpha}_2}{\sigma_2} + u_{13} \frac{\alpha_3 - \bar{\alpha}_3}{\sigma_3}$$

Pour rendre les résultats plus comparables on a entrepris de présenter dans les 4 cas y_1 sous la forme :

$y_1 = l_1 (\alpha_1 - \bar{\alpha}_1) + l_2 (\alpha_2 - \bar{\alpha}_2) + l_3 (\alpha_3 - \bar{\alpha}_3)$, l_1, l_2, l_3 étant les composantes d'un vecteur unitaire, ce qui conduit à multiplier dans chacun des cas les 3 u_{11}, u_{12} et u_{13} par un même nombre. On obtient ainsi pour les 4 cas les résultats suivants :

$$(1) \begin{cases} \text{Cas 1 : } y_1 = 0,127 (\alpha_1 - \bar{\alpha}_1) + 0,688 (\alpha_2 - \bar{\alpha}_2) + 0,715 (\alpha_3 - \bar{\alpha}_3) \\ \text{Cas 2 : } y_1 = 0,520 (\alpha_1 - \bar{\alpha}_1) + 0,755 (\alpha_2 - \bar{\alpha}_2) + 0,400 (\alpha_3 - \bar{\alpha}_3) \\ \text{Cas 3 : } y_1 = 0,934 (\alpha_1 - \bar{\alpha}_1) + 0,344 (\alpha_2 - \bar{\alpha}_2) + 0,095 (\alpha_3 - \bar{\alpha}_3) \\ \text{Cas 4 : } y_1 = 0,961 (\alpha_1 - \bar{\alpha}_1) + 0,204 (\alpha_2 - \bar{\alpha}_2) + 0,187 (\alpha_3 - \bar{\alpha}_3) \end{cases}$$

On voit que l'on obtient comme expressions de la première composante principale des résultats encore plus différents que pour les u_{11}, u_{12}, u_{13} , ce qui conduit à douter de l'intérêt de la méthode. Par contre, si l'on rapporte λ_1 à la trace, on obtient des valeurs extrêmement voisines, qui confirment dans les 4 cas l'importance de la première composante comme facteur général de croissance. Mais ceci ne suffit pas pour exprimer les variables x_1, x_2, x_3 en fonction de facteurs généraux et spécifiques, tels que le suggère l'analyse factorielle.

Si nous revenons à l'expression de y_1 en fonction des variables centrées $\alpha_1 - \bar{\alpha}_1, \alpha_2 - \bar{\alpha}_2, \alpha_3 - \bar{\alpha}_3$ nous avons vu que le premier cas défavorise la variable PNB au profit des 2 autres; au contraire les 3^e et 4^e cas favorisent la variable PNB et le bon sens conduit à penser que la meilleure estimation de y_1 doit être aussi éloignée des cas 3 et 4 que du cas 1, donc peu éloignée de celle du cas 2. Mais quelle expression préférer? En réalité il est impossible d'effectuer un choix basé sur la mathématique.

Considérons une loi de Gauss à p variables x_1, x_2, \dots, x_p . Sa densité de probabilité est égale $Ke^{-\frac{Q}{2}}$ avec :

$$Q = \text{forme quadratique d'ordre } p = x'Ax, \quad K = (2\pi)^{-\frac{p}{2}} |A|^{-\frac{1}{2}}$$

$$A = \begin{bmatrix} a_{11} & \dots & a_{1p} \\ \dots & \dots & \dots \\ a_{p1} & \dots & a_{pp} \end{bmatrix}, \quad x = \begin{bmatrix} x_1 \\ \vdots \\ x_p \end{bmatrix}, \quad x' = \text{vecteur ligne transposé de } x \text{ et } |A| = \text{déter-}$$

minant formé des éléments de la matrice A .

On a d'autre part $A = V^{-1}$ (V étant la matrice des covariances).

L'équation $Q = 1$ représente un hyperellipsoïde dont les axes principaux ont comme directions les p vecteurs propres de la matrice A , alors que les p demi-axes principaux sont égaux aux inverses des racines carrées des p valeurs propres $\lambda_1, \lambda_2, \dots, \lambda_p$ de la matrice A . Les valeurs propres et vecteurs propres de la matrice A permettent donc d'estimer les $1/2$ axes principaux de l'hyperellipsoïde $Q = 1$. Mais ceci suppose que les variables x_1, x_2, \dots, x_p suivent une loi de Gauss à p variables. En effet pour obtenir à partir d'une matrice symétrique d'ordre p , p valeurs propres et p vecteurs propres, il suffit que la matrice soit définie positive.

D'autre part il existe toujours une loi de Gauss à p variables ayant pour matrice des covariances la matrice des covariances expérimentalement observées et ce sont les $1/2$ axes principaux de l'hyperellipsoïde correspondant ($Q = x'V^{-1}x = 1$)* que les valeurs propres et vecteurs propres permettent d'estimer. Or on peut imaginer un très grand nombre de fonctions ayant pour matrice des covariances la matrice des covariances observées. Dans tous les cas les vecteurs propres sont les mêmes. Et si la loi des x_1, x_2, \dots, x_p est trop éloignée d'une loi de Gauss à p variables, les estimations n'ont plus qu'une valeur réduite.

Considérons maintenant une loi de Gauss à 3 variables x_1, x_2, x_3 . A cette loi de Gauss correspond un ellipsoïde dont les directions principales sont les vecteurs propres de la matrice $A = V^{-1}$. On peut imaginer cet ellipsoïde tracé dans l'espace ainsi que ses directions principales. Supposons que sans toucher à x_1 et x_2 on choisisse comme unité de comptage sur l'axe Ox_3 une unité 10 fois plus petite. Tout se passera comme si l'on avait effectué une affinité parallèlement à Ox_3 de rapport $1/10$. On obtiendra un ellipsoïde très aplati selon le plan $x_1 Ox_2$, et les axes de symétrie de l'ellipsoïde seront très différents de ce qu'ils étaient avant l'affinité. Or garder la variable x_3 et diviser l'unité de mesure sur Ox_3 par 10, ou garder l'unité de mesure sur Ox_3 et diviser x_3 par 10 conduit au même résultat : Les axes de symétrie de l'ellipsoïde sont grandement transformés et il en est de même des vecteurs propres de la matrice A .

Reprenons maintenant notre recherche de la meilleure première composante principale. Nous avons vu que, du point de vue de la mathématique, il est impossible de fixer un choix. On peut cependant s'appuyer sur certaines remarques. Par exemple on pourra penser que plus une variable a une corrélation multiple importante avec les 2 autres, plus sa contribution sera importante dans la détermination de la première composante principale. Le coefficient de corrélation multiple entre α_1 d'une part, α_2 et α_3 d'autre part est égal à $R_{1,23} = 0,851$. De même on a : $R_{2,13} = 0,880$ et $R_{3,12} = 0,784$. On est donc conduit à penser que c'est d'abord α_2 qui contribue le plus à la formation de la première composante principale y_1 , puis α_1 et enfin seulement α_3 , ce qui est conforme à ce qui est communément admis.

Posons

$$x_1 = \alpha_1 - \bar{\alpha}_1, \quad x_2 = (k_2)^z (\alpha_2 - \bar{\alpha}_2), \quad x_3 = (k_3)^z (\alpha_3 - \bar{\alpha}_3).$$

Nous allons maintenant analyser l'évolution du vecteur unitaire l (l_1, l_2, l_3) en fonction de z . Laissons de côté le cas 4 et considérons le cas 5 défini par $z = -1/2$, qui conduit à :

Cas 5

$$y_1 = 0,022 (\alpha_1 - \bar{\alpha}_1) + 0,442 (\alpha_2 - \bar{\alpha}_2) + 0,897 (\alpha_3 - \bar{\alpha}_3)$$

Nous avons porté sur le graphique 1 les valeurs de l_1, l_2, l_3 en fonction des valeurs de z : $-0,5, 0, 0,5, 1$ et nous avons ajusté par les courbes l_1, l_2, l_3 en fonction de z . Nous voyons que si nous imposons comme conditions $l_2 > l_1 > l_3$ il reste pour le choix de z un domaine assez étroit comprenant $z = 0,5$, ce qui conduit à choisir z voisin de $0,5$ et plus simplement $z = 0,5$.

Mais on peut apporter une autre limitation aux variations du vecteur l . En effet si k_2 et k_3 étaient égaux à $k_1 (=1)$, $(k_2)^z$ et $(k_3)^z$ seraient constants quel que soit z et le vecteur l demeurerait fixe. On est donc amené à choisir des agrégats, des branches, des sous-branches, etc., qui ont à peu près la même valeur monétaire, ce qui n'est pas toujours possible, car il n'est pas toujours possible de scinder une branche (ou de regrouper deux sous-branches) en des parties économiquement valables. Mais cette réserve faite, on doit chercher à déterminer des unités économiques ayant la valeur monétaire la plus voisine possible. La combi-

* On peut aussi considérer l'hyperellipsoïde inverse correspondant à $Q' = x' V x = 1$, qui a mêmes axes principaux et des valeurs de demi-axes inverses des précédents.

