

# REVUE DE STATISTIQUE APPLIQUÉE

L. BONIFAS

Y. ESCOUFIER

P. L. GONZALEZ

R. SABATIER

## **Choix de variables en analyse en composantes principales**

*Revue de statistique appliquée*, tome 32, n° 2 (1984), p. 5-15

[http://www.numdam.org/item?id=RSA\\_1984\\_\\_32\\_2\\_5\\_0](http://www.numdam.org/item?id=RSA_1984__32_2_5_0)

© Société française de statistique, 1984, 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/>

# CHOIX DE VARIABLES EN ANALYSE EN COMPOSANTES PRINCIPALES

L. BONIFAS, Y. ESCOUFIER, P.L. GONZALEZ, R. SABATIER

*Unité de Biométrie 9, Place Pierre Viala 34000 Montpellier Cedex*

---

## RESUME

Un bref rappel des propriétés d'approximation de l'Analyse en Composantes Principales permet de poser clairement le problème de l'Analyse en Composantes Principales par rapport à des variables instrumentales et d'en expliciter les solutions. On en déduit une démarche de choix de variables en Analyse en Composantes Principales. Les résultats d'une application sont présentés.

## INTRODUCTION

Cet article voudrait atteindre deux objectifs : en premier lieu, il souhaite contribuer à une promotion de l'emploi des méthodes "d'analyse en Composantes Principales par rapport à des variables instrumentales", et de "Choix de variables en Analyse en Composantes Principales". Pour cela il apporte les justifications mathématiques de ces méthodes et fait référence aux programmes disponibles auprès des auteurs.

Le second objectif est de rappeler clairement que la mise en œuvre d'une Analyse en Composantes Principales (ACP) repose non seulement sur des données, mais aussi sur le choix d'une métrique. Il en découle que toutes les procédures qui élaborent des cheminements compliqués pour choisir des variables en ACP puis mettent en œuvre l'ACP sur les variables choisies sans se poser la question de la métrique à utiliser avec ces variables, sont améliorables. Les résultats des paragraphes II et III montrent clairement la marche à suivre dans la résolution de ce problème.

## I. RAPPEL SUR L'ANALYSE EN COMPOSANTES PRINCIPALES (ACP)

Disposant d'un tableau  $T$ ,  $n \times p$ , des mesures de  $p$  variables quantitatives faites sur  $n$  individus, le chercheur qui veut explorer les liaisons linéaires entre les variables et les ressemblances entre les individus peut faire appel à l'ACP. Pour mettre en œuvre cette méthode, il devra faire des choix que la présentation des techniques d'Analyse des Données popularisée en France par le livre de F. CAILLIEZ et J.P. PAGES [2], rend explicites.

L'étude des liaisons linéaires entre les variables suppose le calcul de coefficients de covariance et de corrélation, ce qui nécessite l'affectation de poids aux individus. Une matrice diagonale positive  $D$  telle que  ${}^t\delta_n D\delta_n = 1$  où  $\delta_n$  est un vecteur de  $\mathbf{R}^n$  à  $n$  composantes égales à l'unité, définit ces poids.

De même, l'étude des ressemblances entre les individus suppose une quantification de ces ressemblances par le calcul de distances euclidiennes entre les individus rendu possible par le choix d'une métrique  $Q$  sur  $\mathbf{R}^p$ . Ces choix faits et  $X$  étant le tableau des données centrées pour les poids  $D(X = (I_n - \delta_n {}^t\delta_n D)T)$ , il est possible de calculer les axes principaux de  $\mathbf{R}^p$  et  $\mathbf{R}^n$  définis par :

$$U_\alpha \in \mathbf{R}^p; \quad \alpha = 1, \dots, p$$

tels que  ${}^tXDXQU_\alpha = \lambda_\alpha U_\alpha$  avec  ${}^tUQU = I_p$

$$\text{où} \quad U = (U_1 \begin{array}{c} \vdots \\ U_2 \end{array} \begin{array}{c} \vdots \\ \vdots \end{array}, \dots, U_p)$$

$$V_\alpha \in \mathbf{R}^n; \quad \alpha = 1, \dots, p$$

tels que  $XQ{}^tXDV_\alpha = \lambda_\alpha V_\alpha$  avec  ${}^tVDV = I_p$

$$\text{où} \quad V = (V_1 \begin{array}{c} \vdots \\ V_2 \end{array} \begin{array}{c} \vdots \\ \vdots \end{array}, \dots, V_p)$$

De ces axes, on déduit les coordonnées principales

$$\varphi_\alpha = \sqrt{\lambda_\alpha} U_\alpha \in \mathbf{R}^p \quad \text{et} \quad \psi_\alpha = \sqrt{\lambda_\alpha} V_\alpha \in \mathbf{R}^n$$

(aussi appelées composantes principales pour  $\psi_\alpha$ ) qui vérifient :

$$\|{}^tXDXQ - \sum_{\alpha=1}^q \varphi_\alpha {}^t\varphi_\alpha Q\|^2 = \sum_{\alpha=q+1}^p \lambda_\alpha^2 \quad \text{où} \quad \|AQ\|^2 = \text{Tr}[(AQ)^2]$$

$$\|XQ{}^tXD - \sum_{\alpha=1}^q \psi_\alpha {}^t\psi_\alpha D\|^2 = \sum_{\alpha=q+1}^p \lambda_\alpha^2 \quad \text{où} \quad \|AD\|^2 = \text{Tr}[(AD)^2]$$

Ainsi, dans le cas  $q = p$ , les coordonnées principales dans  $\mathbf{R}^p$  permettent une reconstitution exacte de  ${}^tXDX$ , matrice de variance des variables observées : la covariance des variables  $i$  et  $i'$  est égale au produit scalaire  $\sum_{\alpha=1}^p \varphi_{\alpha i} \varphi_{\alpha i'}$  des vecteurs  $(\varphi_{\alpha i}; \alpha = 1, \dots, p)$  et  $(\varphi_{\alpha i'}; \alpha = 1, \dots, p)$ .

De même, dans ce cas, les composantes principales permettent une reconstitution de  $XQ{}^tX$ , matrice des produits scalaires entre individus : la ressemblance entre les individus  $j$  et  $j'$  est quantifiée par la distance  $\sum_{\alpha=1}^p (\psi_{\alpha j} - \psi_{\alpha j'})^2$  entre les points de coordonnées  $(\psi_{\alpha j}; \alpha = 1, \dots, p)$ ,  $(\psi_{\alpha j'}; \alpha = 1, \dots, p)$ .

Si  $Q$  est l'identité de  $\mathbf{R}^p$  et  $D = (1/n)I_n$ , les représentations obtenues en se limitant à  $(\varphi_1, \varphi_2)$  et  $(\psi_1, \psi_2)$  fournissent des approximations visuelles des covariances et des ressemblances entre individus.

On sait par ailleurs [10], [13] que pour  $q$  donné les approximations de  ${}^tXDXQ$  et  $XQ{}^tXD$  fournies respectivement par  $\sum_{\alpha=1}^q \varphi_{\alpha} {}^t\varphi_{\alpha} Q$  et  $\sum_{\alpha=1}^q \psi_{\alpha} {}^t\psi_{\alpha} D$  sont les meilleures qui puissent fournir  $q$  vecteurs  $Q$ -orthogonaux de  $R^p$  d'une part et  $q$  vecteurs  $D$ -orthogonaux de  $R^n$  d'autre part quelles que soient les normes unitairement invariantes utilisées dans  $R^p$  ou  $R^n$  pour apprécier la qualité de ces approximations [14].

Un résultat analogue, connu sous le nom de "formule de reconstitution des données" existe pour  $X$ . Pour  $Q = L {}^tL$ ,  $A$  de dimensions  $n \times p$  et  $\|A\|^2 = \text{Tr}({}^tAA)$ , on a :

$$\left\| D^{1/2} XL - \sum_{\alpha=1}^q D^{1/2} \frac{\psi_{\alpha} {}^t\varphi_{\alpha}}{\sqrt{\lambda_{\alpha}}} L \right\|^2 = \sum_{\alpha=q+1}^p \lambda_{\alpha}$$

Il sera utile pour la suite de remarquer que :

$$\|XQ{}^tXD - \sum_{\alpha=1}^q \psi_{\alpha} {}^t\psi_{\alpha} D\|^2 = \|XQ{}^tXD\|^2 (1 - RV^2(XQ{}^tXD, \sum_{\alpha=1}^q \psi_{\alpha} {}^t\psi_{\alpha} D))$$

où  $RV$  est le coefficient défini entre deux opérateurs  $D$ -symétriques de  $R^n$ ,  $AD$  et  $BD$  par

$$RV(AD, BD) = \frac{\text{Tr}(ADBD)}{\|AD\| \|BD\|}$$

On remarque aisément que ce coefficient peut s'interpréter comme le cosinus des vecteurs de  $R^{n \times n}$  dont les éléments sont respectivement fournis par ceux de  $AD$  et  $BD$ . De plus, la symétrie de  $A$  et  $B$  permet d'assurer la positivité de  $\text{Tr}(ADBD)$ . Il en résulte que  $0 \leq RV(AD, BD) \leq 1$ . D'autres propriétés de ce coefficient sont exposées en (3).

## II. L'ANALYSE EN COMPOSANTES PRINCIPALES PAR RAPPORT A DES VARIABLES INSTRUMENTALES

**II.1.** — D'après ce qui précède, à une étude  $(X, Q, D)$  nous savons associer un opérateur de  $R^n$ ,  $WD = XQ{}^tXD$  qui en est caractéristique au sens où ses vecteurs propres  $(\psi_{\alpha}; \alpha = 1, \dots, p)$  convenablement normés permettent une représentation des individus dans laquelle les produits scalaires entre individus sont égaux aux éléments de  $W$ , tandis que pour  $q < p$ , les vecteurs  $(\psi_{\alpha}, \alpha = 1, \dots, q)$  permettent la meilleure représentation possible de  $W$ .

Ce point étant acquis, on peut alors s'interroger sur le problème suivant : une étude  $(X, Q, D)$  étant donnée et des données  $Y$  étant observées sur les mêmes individus munis des mêmes poids, peut-on trouver une matrice symétrique définie ou semi-définie positive, c'est-à-dire une métrique "au sens large"  $M$ , à utiliser avec  $Y$  pour que les études  $(X, Q, D)$  et  $(Y, M, D)$  conduisent aux mêmes opérateurs, donc aux mêmes composantes principales et par conséquent aux mêmes représentations des individus. La suite montrera que la réponse à cette question est en

général négative. La formulation adoptée permet cependant de bien mettre en évidence ce qu'il est possible de faire.

Supposons en effet que X soit de dimension  $n \times p$  et Y de dimension  $n \times q$ .  
Posons :

$$S_{YY} = {}^tYDY \text{ supposée inversible}$$

$$S_{XY} = {}^tXDY \text{ et } S_{YX} = {}^tS_{XY}$$

$$B = S_{YY}^{-1} S_{YX} Q S_{XY} S_{YY}^{-1}$$

Le calcul permet d'établir le résultat suivant

**Lemme** : Pour toute matrice symétrique R,  $q \times q$ , on a :

$$\text{Tr}(YB^t YD YR^t YD) = \text{Tr}(XQ^t X D YR^t YD)$$

Il en découle :

**Propriété 1** : Pour toute métrique M,  $q \times q$

$$\|XQ^t X D - YM^t YD\|^2 = \|XQ^t X D - YB^t YD\|^2 + \|YB^t YD - YM^t YD\|^2$$

**Démonstration** : On peut écrire

$$\begin{aligned} \|XQ^t X D - YM^t YD\|^2 &= \|XQ^t X D - YB^t YD\|^2 + \|YB^t YD - YM^t YD\|^2 \\ &\quad + 2 \text{Tr} [(XQ^t X D - YB^t YD)(YB^t YD - YM^t YD)] \end{aligned}$$

et le résultat s'obtient en appliquant le lemme à  $R = B - M$ .

Cette propriété montre que la norme de la différence entre  $XQ^t X D$  et  $YM^t YD$  dépend de deux termes :

- *d'une part*  $\|XQ^t X D - YB^t YD\|^2$  qui ne contient pas la métrique M. C'est dire que ce terme mesure ce qui est présent dans (X, Q, D) et ne pourra jamais être reconstruit à partir d'une étude bâtie sur Y.
- *d'autre part*  $\|YB^t YD - YM^t YD\|^2$  qui dépend explicitement du choix de M, le meilleur choix étant B.

## II.2.

**Propriété 2**  $\|XQ^t X D - YB^t YD\|^2 = \|XQ^t X D\|^2 (1 - RV^2(XQ^t X D, YB^t YD))$

**Démonstration** : Le lemme appliqué à la matrice B elle-même donne

$$\|YB^t YD\|^2 = \text{Tr} [(YB^t YD)^2] = \text{Tr}(XQ^t X D YB^t YD)$$

Il en découle :

$$\|XQ^t X D - YB^t YD\|^2 = \|XQ^t X D\|^2 - \|YB^t YD\|^2$$

d'où le résultat.

Posons  $P_Y = Y S_{YY}^{-1} {}^tYD$  le projecteur D-orthogonal sur le sous-espace  $E_Y$  de  $R^n$  engendré par les variables qui définissent Y.

$$YB^t YD = P_Y XQ^t X {}^tP_Y D$$

On trouve donc que la meilleure reconstitution de la représentation des individus fournie par  $(X, Q, D)$  sur la base des variables  $Y$  est fournie par l'étude  $(P_Y X, Q, D)$  c'est-à-dire l'étude faite sur les projections orthogonales des variables  $X$  sur le sous-espace  $E_Y$  pour la métrique  $Q$ . Cette remarque donne un point de vue sur la régression peu souvent mis en avant. De plus si  $X$  se réduit à une seule variable, le calcul montre que  $RV(X^tXD, YB^tYD)$  est égal au carré du coefficient de corrélation linéaire entre  $X$  et  $P_Y X$ , c'est-à-dire au carré du coefficient de corrélation multiple entre  $X$  et  $Y$ .

**II.3.** Intéressons-nous maintenant au second terme du second membre de l'égalité de la propriété 1.

Soit  $U$  la matrice  $q \times r$  des  $r$  axes principaux de  $R^q$  associés aux plus grandes valeurs propres de l'étude  $(Y, B, D)$ .

**Propriété 3 :** Parmi toutes les métriques  $M$ ,  $q \times q$ , de rang  $r$  le meilleur choix est fourni pour  $M = BU^tUB$ .

*Démonstration :* On sait que  ${}^tYDYBU = U\Lambda$  avec  ${}^tUBU = I_r$ . Un calcul classique utilisé dans l'établissement des "formules de transition" permet d'établir que  $YBU = \psi$  où  $\psi$  est la matrice  $n \times r$  des  $r$  premières composantes principales de  $(Y, B, D)$ . Choisissons alors  $M = BU^tUB$ .

On a  $YM^tYD = \psi^t \psi D = \sum_{\alpha=1}^r \psi_\alpha {}^t\psi_\alpha D$ . Un retour aux résultats rappelés dans le paragraphe I établit la propriété.

De plus, la remarque finale du paragraphe I permet d'écrire pour ce choix particulier de  $M$  :

$$\|YB^tYD - YM^tYD\|^2 = \|YB^tYD\|^2 (1 - RV^2(YB^tYD, YM^tYD))$$

**Corollaire :** Pour  $M = BU^tUB$

$$\|XQ^tXD - YM^tYD\|^2 = \|XQ^tXD\|^2 (1 - RV^2(XQ^tXD, YB^tYD) RV^2(YB^tYD, YM^tYD))$$

*Démonstration :* on a

$$\begin{aligned} \|YB^tYD - YM^tYD\|^2 &= \|XQ^tXD\|^2 (1 - RV^2(YB^tYB, YM^tYD)) \\ &= \|XQ^tXD\|^2 \frac{\|YB^tYD\|^2}{\|XQ^tXD\|^2} (1 - RV^2(YB^tYB, YM^tYD)) \\ &= \|XQ^tXD\|^2 RV^2(XQ^tXD, YB^tYB) (1 - RV^2(YB^tYD, YM^tYD)) \end{aligned}$$

et le résultat en découle en utilisant les propriétés 1 et 2.

En conclusion, l'Analyse en Composantes Principales des données  $Y$  par rapport à l'étude de référence  $(X, Q, D)$  revient à réaliser l'ACP du triplet  $(Y, B, D)$  qui est équivalente du point de vue de la représentation des individus à l'ACP de  $(P_Y X, Q, D)$ . Un programme permettant la mise en oeuvre de cette méthode a été

intégré dans la bibliothèque ANADO disponible au CNUSC. Des copies peuvent être fournies par les auteurs.

L'annexe 2 de la thèse de R. SABATIER [13] fournit une description précise de ce programme.

*Remarque 1*

Posons  $Z = S_{YY}^{-1}U$

Alors  ${}^tYDYBU = U\Lambda$  avec  ${}^tUBU = I_r$

implique  $S_{YX} QS_{XY}Z = S_{YY}Z\Lambda$  avec  ${}^tZS_{YX} QS_{XY}Z = I_r$   
 et  ${}^tZS_{YY}Z = \Lambda^{-1}$

Les solutions de l'ACPVI telles que proposées pour  $Q = I_p$  par RAO [9] correspondent donc à  $Z \Lambda^{1/2}$ . Elles conduisent à des composantes principales

$$\xi = YZ \Lambda^{1/2}.$$

Or

$$\begin{aligned} \xi &= YZ \Lambda^{1/2} = Y S_{YY}^{-1} S_{YX} Q S_{XY} S_{YY}^{-1} U \Lambda^{-1/2} \\ &= Y B U \Lambda^{-1/2} = \psi \Lambda^{-1/2} \end{aligned}$$

*Remarque 2*

On peut écrire :

$$\begin{aligned} \text{Tr}(YM{}^tYD) &= \text{Tr}(YBU{}^tUB{}^tYD) \\ &= \text{Tr}({}^tUBU\Lambda) \\ &= \text{Tr}(\Lambda^{1/2} {}^tZ{}^tYDXQ{}^tXDYZ\Lambda^{1/2}) \\ &= \text{Tr}({}^t\xi DXQ{}^tXD\xi) \end{aligned}$$

Si on note  $\xi = (\xi_1, \dots, \xi_q)$ ,  $X = (X_1, \dots, X_p)$  et  $Q_{ij}$  l'élément (i, j) de la métrique Q, on peut écrire :

$$\text{Tr}(YM{}^tYD) = \sum_{\alpha=1}^q \sum_{i=1}^p \sum_{j=1}^p Q_{ij} \text{cov}(\xi_\alpha, X_i) \text{cov}(\xi_\alpha, X_j)$$

Si  $Q = I_p$ , on a donc la somme des carrés des covariances entre les variables  $\xi_\alpha$  et  $X_i$ . On retrouve la présentation faite par OBADIA en termes de "composantes Explicatives" [8].

### III. CHOIX DE VARIABLES

L'étude (X, Q, D) étant donnée, nous extrayons r colonnes de X qui permettent la construction d'une matrice notée  $Y_{[r]}$ . Soit  $B_{[r]}$  la matrice à associer à  $Y_{[r]}$  pour que les études (X, Q, D) et  $(Y_{[r]}, B_{[r]}, D)$  soient les plus proches possibles.  $B_{[r]}$  a été définie au paragraphe précédent.

Soit  $\mathcal{R}$  l'ensemble des matrices à  $r$  colonnes extraites de  $X$ , nous souhaitons résoudre le problème suivant :

$$\left\| \begin{array}{l} \text{Trouver } Y_{[r]}^* \in \mathcal{R} \text{ telle que} \\ \text{RV}(XQ^tXD, Y_{[r]}^* B_{[r]}^* {}^tY_{[r]}^* D) = \sup_{Y_{[r]} \in \mathcal{R}} \text{RV}(XQ^tXD, Y_{[r]} B_{[r]} {}^tY_{[r]} D) \\ Y_{[r]} \in \mathcal{R} \end{array} \right.$$

D'après la propriété 2, ce problème est équivalent au suivant :

$$\left\| \begin{array}{l} \text{Trouver } Y_{[r]}^* \in \mathcal{R} \text{ telle que} \\ \|XQ^tXD - Y_{[r]}^* B_{[r]}^* {}^tY_{[r]}^* D\|^2 = \inf_{Y_{[r]} \in \mathcal{R}} \|XQ^tXD - Y_{[r]} B_{[r]} {}^tY_{[r]} D\|^2 \\ Y_{[r]} \in \mathcal{R} \end{array} \right.$$

Pour  $r$  et  $p$  suffisamment petits, on peut toujours envisager le calcul exhaustif sur l'ensemble des possibilités offertes. Lorsque  $r$  et  $p$  deviennent raisonnablement grands, cette solution n'est plus possible et l'optimum ne pouvant être défini de façon analytique, on propose un algorithme progressif pour sa construction.

Le lemme permet d'écrire :

$$\text{RV}^2(XQ^tXD, Y_{[r]} B_{[r]} {}^tY_{[r]} D) = \frac{\text{Tr}(XQ^tXD Y_{[r]} B_{[r]} {}^tY_{[r]} D)}{\|XQ^tXD\|^2}$$

Posons  ${}^tY_{[r]} D Y_{[r]} = L_{[r]} {}^tL_{[r]}$  ; on a :

$$\text{Tr}(XQ^tXD Y_{[r]} B_{[r]} {}^tY_{[r]} D) = \text{Tr}[(L_{[r]}^{-1} {}^tY_{[r]} DXQ^tXD Y_{[r]} ({}^tL_{[r]})^{-1}]^2$$

Choisissons pour  $L_{[r]} {}^tL_{[r]}$  la décomposition de CHOLESKI de  ${}^tY_{[r]} D Y_{[r]}$  ( $L_{[r]}$  est triangulaire inférieure) et posons

$$A_{[r-1]} = L_{[r-1]}^{-1} {}^tY_{[r-1]} DXQ^tXD Y_{[r-1]} ({}^tL_{[r-1]})^{-1}$$

Connaissant les  $[r-1]$  premières variables retenues qui constituent  $Y_{[r-1]}$ , on connaît  $L_{[r-1]}^{-1}$  et il est facile d'en déduire  $L_{[r]}^{-1}$  due à l'introduction d'une  $r$ ème ligne de la matrice symétrique  $A_{[r]}$ .

On a alors :

$$\text{Tr}[(A_{[r]})^2] = \text{Tr}[(A_{[r-1]})^2] + 2 \sum_{i=1}^r (A_{[r]})_{ii}^2 - (A_{[r]})_{rr}^2$$

ce qui sert de base à la construction de l'algorithme progressif qui maximisera  $\text{Tr}(XQ^tXD Y_{[r]} B_{[r]} {}^tY_{[r]} D)$  donc  $\text{RV}^2(XQ^tXD, Y_{[r]} B_{[r]} {}^tY_{[r]} D)$ .

Un programme est également disponible dans la bibliothèque ANADO (utilisable au CNUSC) ou auprès des auteurs. Il permet de s'occuper du choix des variables sans s'occuper du choix de  $B_{[r]}$  (option métrique identité) de limiter la recherche de la métrique à utiliser avec  $Y_{[r]}$  à celle d'une métrique diagonale (option métrique diagonale) ; ou de résoudre le problème général tel qu'exposé ici. Le programme permet d'imposer la prise en compte de certaines variables dans le choix.



#### IV. EXEMPLE

Dans une étude concernant la population atmosphérique, on dispose d'un réseau de 18 stations de mesures numérotées de 1 à 18. On dispose de 562 mesures journalières de fumées noires concomitantes sur ces 18 stations qui permettent de calculer la matrice des corrélations fournies à la fin de l'article.

Ayant décidé d'extraire de ce réseau, un réseau restreint à 10 stations, le programme de choix de variables décrit au paragraphe III fournit les résultats consignés dans la table I.

TABLE I

Station introduite	Valeur de RV après l'introduction de la station
14	0,866
17	0,923
18	0,943
1	0,968
5	0,976
12	0,982
15	0,986
4	0,989
16	0,991
11	0,992

La valeur du dernier coefficient RV obtenu est en elle-même une information sur la qualité de l'approximation faite qui s'interprète par référence à la propriété 2.

Dans la pratique, il apparaît que donnée sous cette forme, la solution ne satisfait pas pleinement le spécialiste de la pollution qui la trouve peu explicite. Pour l'éclairer, plusieurs possibilités sont envisageables :

1) Utilisant le support visuel de l'implantation géographique du réseau initial, on fait apparaître le réseau restreint. Bien que certaines informations puissent être tirées de cette représentation, elle n'apporte aucune précision sur les raisons du choix.

2) Une seconde solution consiste à calculer les coefficients de corrélation multiple entre les variables non retenues et les variables retenues. On passe ainsi d'un critère RV global à des critères plus individuels. Les valeurs obtenues pour l'exemple étudié sont présentées dans la table II.

3) Une dernière solution, qui a l'avantage de mettre en évidence les contraintes d'évolutions du programme, consiste à réaliser une classification automatique des stations sur la base de la matrice des corrélations puis à positionner sur cette classification les stations retenues. On voit alors que le programme cherche à énumérer les variables dans un ordre tel que, à tout moment, les différents partitions apparues soient représentées de façon équitable.

TABLE II

Stations non retenues	Coefficients de corrélation multiple avec les stations retenues
2	0,91
3	0,89
6	0,86
7	0,94
8	0,93
9	0,89
10	0,96
13	0,83

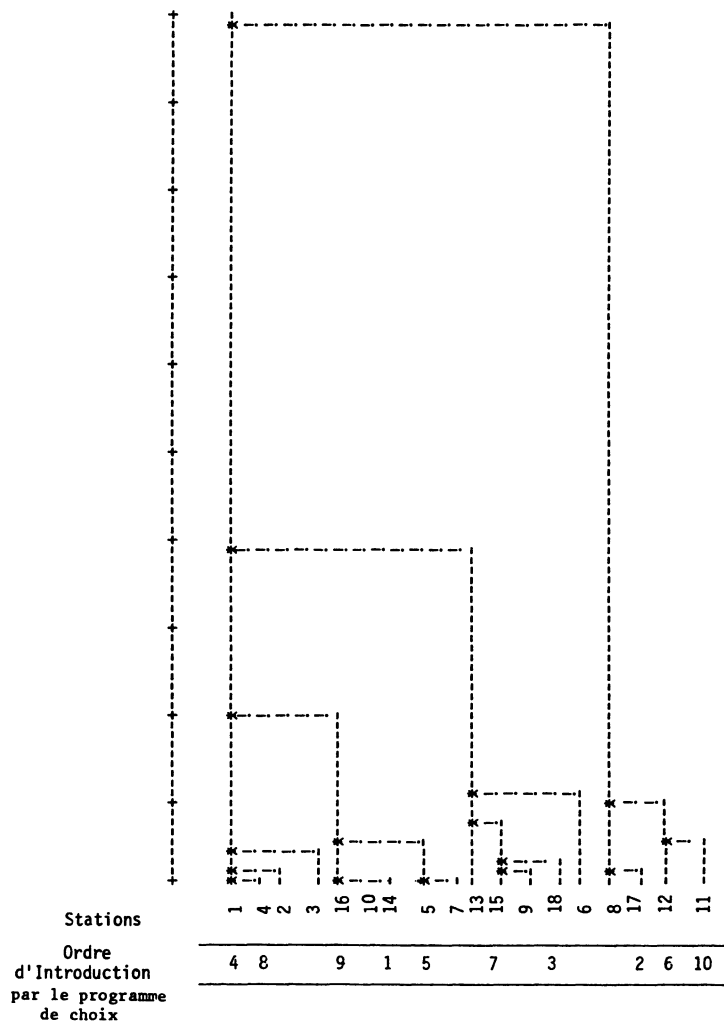


Figure 1. – Classification obtenue en appliquant le programme CAH2COORD du logiciel ADDAD à la matrice des corrélations



Pour interpréter les résultats obtenus sur l'exemple traité et présentés dans la figure 1, on notera que les stations 14, 17 et 18 avaient été imposées pour des raisons extérieures à la logique mathématique.

\*

L'application qui illustre cet article est tirée d'un travail réalisé dans le cadre d'une collaboration avec la direction environnement de la Société Nationale Elf Aquitaine. Les auteurs tiennent à remercier cette société pour le soutien qu'elle leur apporte.

## BIBLIOGRAPHIE

- [1] BERNARD C. DOCHI. – *Choix de variables en Analyse des Données*, 29 juin 1979, Thèse USTL Montpellier.
- [2] F. CAILLIEZ, J.P. PAGES (1976). – *Introduction à l'Analyse des Données*, SMASH, 9, rue Duban, 75016 Paris.
- [3] Y. ESCOUFIER (1973). – Le traitement des variables vectorielles, *Biometrics*, 29, p. 751-760.
- [4] Y. ESCOUFIER, P. ROBERT (1979). – *Choosing variables and metrics by optimizing the RV coefficient*, In: *Optimizing methods in Statistic*. Ed. J.S. Rustagi Academic Press INC, p. 205-219.
- [5] P.L. GONZALEZ (1981). – *Choix de variables, applications au choix de Stations dans un réseau*, Rapport Technique 8104, CRIG Montpellier.
- [6] P.L. GONZALEZ. – *Analyse statistique de données psycho-sensorielles*, 19 avril 1982. Thèse USTL Montpellier.
- [7] L. LEBART, A. MORINEAU, J.P. FENELON (1979). – *Traitement des données statistiques : Méthodes et Programmes*, Dunod.
- [8] J. OBADIA (1978). – L'analyse en composantes explicatives, *R.S.A.*, XXVI, n° 4, p. 5-28.
- [9] C.R. RAO (1964). – The use and the interpretation of principal component analysis in applied research, *Sankya, Sci. A.*, 26, p. 329-359.
- [10] C.R. RAO (1980). – Matrix approximations and reduction of dimensionality in multivariate statistical analysis, In: *Multivariate Analysis V*. Ed. P.R. Krishnaiah, NHPC.
- [11] P. ROBERT, Y. ESCOUFIER (1976). – A unifying tool for linear Multivariate statistical Methods: The RV coefficient, *Appl. Stat. C*, 25 (3), p. 257-265.
- [12] R. SABATIER (1981). – *Two examples of choosing variables and metrics with the RV coefficient*, Rapport Technique 8102, CRIG Montpellier.
- [13] R. SABATIER. – *Approximations d'un tableau de données. Application à la reconstitution des paléoclimats*, 21 mars 1983, Thèse USTL Montpellier.
- [14] R. SABATIER, Y. JAN, Y. ESCOUFIER (1983). – *Approximations d'applications linéaires et Analyse en Composantes Principales*, 3<sup>e</sup> Journées Internationales Analyse des Données et Informatique, Ed. I.N.R.I.A. Domaine de Voluceau-Rocquencourt, B.P. 105, 78153 Le Chesnay Cedex.