

REVUE DE STATISTIQUE APPLIQUÉE

MOHAMED HANAFI

ROGER LAFOSSE

Généralisations de la régression simple pour analyser la dépendance de K ensembles de variables avec un $K + 1$ ème

Revue de statistique appliquée, tome 49, n° 1 (2001), p. 5-30

http://www.numdam.org/item?id=RSA_2001__49_1_5_0

© Société française de statistique, 2001, 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/>

GÉNÉRALISATIONS DE LA RÉGRESSION SIMPLE POUR ANALYSER LA DÉPENDANCE DE K ENSEMBLES DE VARIABLES AVEC UN $K + 1$ ème

Hanafi Mohamed, Lafosse Roger

E.N.I.T.I.A.A., rue de la Géraudière B.P. 82225 44322 Nantes Cedex 3

hanafi@enitaa-nantes.fr

Lab. de Statistique et Probabilités, UMR C55830 Université Paul Sabatier,

118 rte de Narbonne, 31062 Toulouse Cedex 4 lafosse@cict.fr

RÉSUMÉ

Une analyse de dépendance simultanée considérée entre une matrice et K matrices (Lafosse & Hanafi, 1997) peut se présenter comme une analyse factorielle de deux tableaux, grâce à un projecteur particulier. Les régressions simples simultanées qui s'ensuivent permettent d'évaluer la contribution à cette dépendance de chaque ligne de chaque matrice. Les coefficients de régression simultanée sont globalement solution d'un problème d'optimisation nouvellement posé.

Mots-clés : corrélation linéaire, analyses factorielles, concordance entre matrices, moindres carrés, biplot.

ABSTRACT

By introducing a projector, the simultaneous dependency analysis considered between a matrix and K matrices by Lafosse & Hanafi (1997) boils down to an analysis between two matrices. Some simultaneous simple regressions then lead to a calculus of the contributions to this dependency for every rows of each matrix. The simultaneous calculus of the regression coefficients now comes from the solution of one optimization problem.

Keywords : linear correlation, factor analysis, agreement between matrices, least squares, biplot.

1. Introduction

Cazes & coll. (1977) ont proposé une étude de la dépendance simultanée entre une variable qualitative Y et K autres variables qualitatives. Cette étude revenait à définir un codage des modalités de la variable Y fortement dépendant de chacune des K variables qualitatives. Leur analyse n'est pas une analyse canonique classique entre deux ensembles d'indicatrices des modalités, car chacun des K sous-ensembles

d'indicatrices de modalités associé respectivement à chacune des K variables est muni de sa propre semi-métrique de Mahalanobis, de sorte que la particularité de chaque sous-ensemble est bien prise en compte. Leur analyse est cependant une analyse canonique de deux tableaux d'indicatrices, en convenant que la métrique relative à l'ensemble des modalités des K variables est la métrique bloc diagonale dont les blocs sont ces semi-métriques de Mahalanobis.

L'analyse de Tucker (1958) de deux tableaux est une analyse canonique si l'on convient que pour l'obtenir il suffit de remplacer les métriques de Mahalanobis par les métriques identité (par exemple, Lafosse (1997)), ou encore si l'on considère que les composantes de Tucker des tableaux formés des composantes principales réduites, obtenus par les deux ACP respectives, sont les composantes de l'analyse canonique classique des deux tableaux (Ten Berge, 1977).

Considérons l'analyse de Tucker généralisée visant une étude de la dépendance simultanée d'un tableau \mathbf{Y} avec K autres tableaux $\mathbf{X}_i, i = 1, \dots, K$, et qui reprend la démarche de Cazes & coll.. La métrique bloc diagonale dont il est question ci-dessus devient la métrique identité, de sorte que nous sommes ramenés à l'analyse de Tucker de \mathbf{Y} et du tableau \mathbf{X} concaténé des K tableaux. Il n'y a alors plus prise en compte de la partition en K sous-ensembles du tableau \mathbf{X} , comme cela se produisait chez Cazes à cause des métriques. La partition en K sous-ensembles n'est en fait pas considérée dans cette analyse de Tucker.

Lafosse et Hanafi (1997) ont proposé une généralisation de Tucker visant la même étude, où les composantes successives de \mathbf{Y} définies s'associent chacune à K composantes partenaires trouvées dans chacun des K tableaux. La première composante de cette analyse est cependant la première composante de la «généralisation» précédente. Cette première solution a en fait la particularité de pouvoir s'associer à toute partition du tableau concaténé, et donc en particulier à la partition de \mathbf{X} considérée. Ensuite, les autres solutions successives obtenues dans cette analyse dépendent de la partition, bien que les métriques soient les métriques identité.

Dans ce papier nous réécrivons cette dernière analyse avec des métriques quelconques, de sorte que la prise en compte de la partition peut aussi provenir des métriques. Cette analyse qui est une étude de dépendance entre triplets statistiques est nommée analyse CONCOR.

Comme nous l'avons souligné ci-dessus, l'analyse de Cazes & coll. est une analyse factorielle de deux tableaux avec un choix particulier de métriques. L'analyse CONCOR peut quant à elle se ramener à une analyse factorielle de deux tableaux de la manière suivante : une fois les termes de l'analyse obtenus, il devient possible de proposer un projecteur \mathbf{P} de sorte que l'analyse factorielle de \mathbf{Y} et de \mathbf{XP} redonnent ces termes. C'est cette propriété qui fonde la suite de l'article. En effet, pouvant se ramener à la situation de deux tableaux, il devient possible de faire le calcul des contributions partielles à la dépendance simultanée des lignes de chaque tableau. Cela revient à définir des régressions simples simultanées, généralisant à plus de deux tableaux un calcul proposé pour deux tableaux (Lafosse 1997). Une évaluation des participations des colonnes est alors envisageable.

La régression simple simultanée entre deux tableaux a été proposée pour répondre à un problème d'ajustement entre nuages d'individus (Lafosse, 1985)

ou pour proposer une matrice estimée caractérisant la participation d'un tableau à la dépendance (image concordante, Lafosse, 1997). Dans ce papier, l'obtention simultanée de l'ensemble des coefficients de régression simple est solution d'un problème d'optimisation. Le problème d'invariance posé vient en prolongement d'un problème d'invariance par rotation considéré notamment par Green (1952), Cliff (1966) et Ten Berge (1977). On suggère que ce qui caractérise la participation d'un tableau à la dépendance doit être invariant par transformation linéaire de l'autre tableau, et pas seulement par isométrie.

Un programme de l'analyse CONCOR avec les développements donnés dans cet article a été écrit en langage MATLAB (Lafosse, 1998). Pour quelques autres écritures on peut s'adresser aux deux auteurs.

2. (K+1)-uples de l'analyse CONCOR : Rappels et compléments

2.1. Définition des composantes

On rappelle la définition des $K + 1$ uples définis en analyse de la concordance d'un tableau avec K tableaux (Lafosse et Hanafi, 1997). Les écritures sont ici plus générales, $K + 1$ métriques euclidiennes respectives quelconques venant maintenant à la place des métriques identités implicitement considérées. Le calcul lui-même mérite alors d'être précisé.

Soient \mathbf{Y} un tableau centré $n \times q$ et K tableaux centrés $\mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_i, \dots, \mathbf{X}_K$ de dimensions respectives $n \times m_i$.

On note $\mathbf{X} = [\mathbf{X}_1 \mathbf{X}_2 \dots \mathbf{X}_K]$ le tableau $n \times m$, avec $m = \sum m_i$, formé par la concaténation des K tableaux \mathbf{X}_i , qui alors peuvent être considérés comme des sous-ensembles de variables de l'ensemble \mathbf{X} .

La matrice diagonale \mathbf{D} représente la métrique des poids affectés à chacun des n individus, la somme des poids valant 1.

La matrice \mathbf{N} désigne la métrique euclidienne considérée dans \mathcal{R}^q et les matrices \mathbf{M}_i celles respectivement considérées dans les espaces \mathcal{R}^{m_i} .

Un couple tel que le couple (\mathbf{Y}, \mathbf{N}) fait référence à l'ensemble des vecteurs lignes de \mathbf{Y} constitutif d'un nuage de points de \mathcal{R}^q représenté avec la métrique \mathbf{N} . Ainsi parlera-t-on du nuage (\mathbf{Y}, \mathbf{N}) et de relation entre nuages (plus de détails dans Hanafi, 1997, chapitre 2).

Un vecteur \mathbf{x} est considéré comme vecteur colonne dans sa notation matricielle et le vecteur ligne transposé est noté \mathbf{x}' . On note $(b'_j, \mathbf{a}'_{1j}, \dots, \mathbf{a}'_{ij}, \dots, \mathbf{a}'_{Kj})'$ le j -ème $K + 1$ uple d'axes de $\mathcal{R}^q \times \mathcal{R}^{m_1} \times \mathcal{R}^{m_2} \times \dots \times \mathcal{R}^{m_K}$ définis dans l'analyse CONCOR de \mathbf{Y} avec les K tableaux $\mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_i, \dots, \mathbf{X}_K, j = 1, 2, \dots, S$.

Le système $\{\mathbf{b}_j\}_{j=1, \dots, S}$ est formé de S vecteurs \mathbf{N} -orthonormés.

Les K systèmes $\{\mathbf{a}_{ij}\}_{j=1, \dots, S} i = 1, \dots, K$, sont chacun formés de S vecteurs respectivement \mathbf{M}_i -orthonormés. Il peut exister pour j fixé des vecteurs \mathbf{a}_{ij} nuls, mais au moins l'un d'entre eux est non nul. La valeur de S qui correspond au nombre de solutions trouvées est au moins égale à $\sup(\text{rang}(\mathbf{Y}'\mathbf{D}\mathbf{X}_j))$.

Pour j fixé, une composante $\mathbf{X}_i \mathbf{M}_i \mathbf{a}_{ij}$ est dite composante *partenaire* dans le tableau \mathbf{X}_i de la composante $\mathbf{Y} \mathbf{N} \mathbf{b}_j$. Les S composantes $\mathbf{Y} \mathbf{N} \mathbf{b}_j$ constituent les *composantes principales de l'analyse CONCOR*. Un j -ème $K + 1$ -uple $(\mathbf{Y} \mathbf{N} \mathbf{b}_j, \mathbf{X}_1 \mathbf{M}_1 \mathbf{a}_{1j}, \dots, \mathbf{X}_i \mathbf{M}_i \mathbf{a}_{ij}, \dots, \mathbf{X}_K \mathbf{M}_K \mathbf{a}_{Kj})$ de composantes est donc formé de la j -ème composante principale et de ses K composantes partenaires.

On précise maintenant le calcul de ces termes. Le critère optimisé en concordance (Lafosse et Hanafi, 1997), est ici considéré sous contraintes de norme des vecteurs \mathbf{a}_{ij} et \mathbf{b}_j avec les métriques euclidiennes respectives choisies, de sorte que l'optimum λ_j^2 obtenu lors du calcul d'un $K + 1$ -uple j vaut :

$$\lambda_j^2 = \sum_{i=1}^K \text{cov}^2(\mathbf{X}_i \mathbf{M}_i \mathbf{a}_{ij}, \mathbf{Y} \mathbf{N} \mathbf{b}_j). \quad (1)$$

On note \mathbf{M} la métrique euclidienne de \mathcal{R}^m formée des blocs diagonaux \mathbf{M}_i , les autres blocs étant nuls. Avec la convention que les matrices $\sum_{h=1}^{j-1} \mathbf{M}_i \mathbf{a}_{ih} \mathbf{a}'_{ih}$, $i = 1, \dots, K$, sont nulles quand $j = 1$, une composante $\mathbf{Y} \mathbf{N} \mathbf{b}_j$ est obtenue en calculant les deux vecteurs \mathbf{b}_j et \mathbf{a}_j , où \mathbf{b}_j est \mathbf{N} -normé dans \mathcal{R}^q et \mathbf{a}_j est \mathbf{M} -normé dans \mathcal{R}^m , associés à la plus grande valeur singulière λ_j de la matrice :

$$\begin{aligned} & [(\mathbf{X}_1 - \mathbf{X}_1 \sum_{h=1}^{j-1} \mathbf{M}_1 \mathbf{a}_{1h} \mathbf{a}'_{1h}) \dots (\mathbf{X}_i - \mathbf{X}_i \sum_{h=1}^{j-1} \mathbf{M}_i \mathbf{a}_{ih} \mathbf{a}'_{ih}) \\ & \dots (\mathbf{X}_K - \mathbf{X}_K \sum_{h=1}^{j-1} \mathbf{M}_K \mathbf{a}_{Kh} \mathbf{a}'_{Kh})]' \mathbf{D} \mathbf{Y} \end{aligned} \quad (2)$$

Le vecteur \mathbf{b}_j est ainsi égal au vecteur $\mathbf{N}^{-1/2} \mathbf{v}_j$, avec \mathbf{v}_j vecteur propre \mathbf{I}_q -normé associé à la plus grande valeur propre λ_j^2 de la matrice symétrique :

$$\mathbf{N}^{1/2} \mathbf{Y}' \mathbf{D} \left[\sum_{h=1}^K \mathbf{X}_i (\mathbf{I}_{mi} - \sum_{h=1}^{j-1} \mathbf{M}_i \mathbf{a}_{ih} \mathbf{a}'_{ih}) \mathbf{M}_i \mathbf{X}'_i \right] \mathbf{D} \mathbf{Y} \mathbf{N}^{1/2}. \quad (3)$$

Les K composantes partenaires d'une j -ème composante principale se déduisent de cette composante principale par les relations :

$$\begin{aligned} \mathbf{X}_i \mathbf{M}_i (\mathbf{I}_{mi} - \sum_{h=1}^{j-1} \mathbf{a}_{ih} \mathbf{a}'_{ih} \mathbf{M}_i) \mathbf{X}'_i \mathbf{D} \mathbf{Y} \mathbf{N} \mathbf{b}_j &= \text{cov}(\mathbf{X}_i \mathbf{M}_i \mathbf{a}_{ij}, \mathbf{Y} \mathbf{N} \mathbf{b}_j) \mathbf{X}_i \mathbf{M}_i \mathbf{a}_{ij}, \\ & i = 1, \dots, K. \end{aligned}$$

La moyenne pondérée $\mathbf{X} \mathbf{M} \mathbf{a}_j$ des composantes $\mathbf{X}_i \mathbf{M}_i \mathbf{a}_{ij}$ partenaires de $\mathbf{Y} \mathbf{N} \mathbf{b}_j$ est donnée par la relation

$$\mathbf{X} \mathbf{M} \mathbf{a}_j = \sum_{i=1}^K p_{ij} \mathbf{X}_i \mathbf{M}_i \mathbf{a}_{ij}, \quad (4)$$

Avec les poids p_{ij}

$$p_{ij} = \text{cov}(\mathbf{X}_i \mathbf{M}_i \mathbf{a}_{ij}, \mathbf{Y} \mathbf{N} \mathbf{b}_j) / \lambda_j. \quad (5)$$

On a ainsi :

$$\lambda_j^2 = \text{cov}^2(\mathbf{X} \mathbf{M} \mathbf{a}_j, \mathbf{Y} \mathbf{N} \mathbf{b}_j) \quad \text{avec} : \sum_{i=1}^K p_{ij}^2 = 1.$$

Un poids p_{ij} mesure la valeur relative de la participation du tableau \mathbf{X}_i lors du calcul du $K+1$ uple j . La composante moyenne $\mathbf{X} \mathbf{M} \mathbf{a}_j$ s'associe à la composante principale $\mathbf{Y} \mathbf{N} \mathbf{b}_j$ par la relation

$$\mathbf{Y} \mathbf{N} \mathbf{Y}' \mathbf{D} \mathbf{X} \mathbf{M} \mathbf{a}_j = \lambda_j \mathbf{Y} \mathbf{N} \mathbf{b}_j, \quad (6)$$

et d'après (5) et (4), le vecteur \mathbf{a}_j est le concaténé des vecteurs $p_{ij} \mathbf{a}_{ij}$

$$\mathbf{a}'_j = [p_{1j} \mathbf{a}'_{1j}, \dots, p_{ij} \mathbf{a}'_{ij}, \dots, p_{Kj} \mathbf{a}'_{Kj}]. \quad (7)$$

Remarques

La démarche en analyse CONCOR vise l'obtention d'un découpage de dépendance entre tableaux. La solution proposée est associée à deux optimisations indépendantes sous contraintes de norme, les deux critères respectifs étant :

$$f(\mathbf{u}, \mathbf{v}) = \text{cov}^2(\mathbf{X} \mathbf{M} \mathbf{u}, \mathbf{Y} \mathbf{N} \mathbf{v})$$

et

$$g(\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_K, \mathbf{f}) = \sum_{i=1}^K \text{cov}^2(\mathbf{X}_i \mathbf{M}_i \mathbf{u}_i, \mathbf{Y} \mathbf{N} \mathbf{f}).$$

Pour une solution j de l'analyse, les deux optimisations respectives ont mené à l'optimum :

$$\text{cov}^2(\mathbf{X} \mathbf{M} \mathbf{a}_j, \mathbf{Y} \mathbf{N} \mathbf{b}_j) = \sum_{i=1}^K \text{cov}^2(\mathbf{X}_i \mathbf{M}_i \mathbf{a}_{ij}, \mathbf{Y} \mathbf{N} \mathbf{b}_j).$$

L'égalité induit le découpage et fonde les développements donnés dans ce papier (section 2.2). Cependant on ne va pas se servir directement des valeurs de critère optimisé pour constituer les mesures de l'analyse. Notant ρ_j les corrélations linéaires $\rho(\mathbf{Y} \mathbf{N} \mathbf{b}_j, \mathbf{X} \mathbf{M} \mathbf{a}_j)$, et ρ_{ij} les coefficients de corrélation $\rho(\mathbf{Y} \mathbf{N} \mathbf{b}_j, \mathbf{X}_i \mathbf{M}_i \mathbf{a}_{ij})$, les S égalités suivantes, qui se déduisent des égalités précédentes entre covariances, seront

celles préférées par la suite (section 3) pour décomposer l'importance relative des composantes partenaires :

$$\rho_j^2 \text{var}(\mathbf{XMa}_j) = \sum_{i=1}^K \rho_{ij}^2 \text{var}(\mathbf{X}_i \mathbf{M}_i \mathbf{a}_{ij}), j = 1, \dots, S,$$

précisant ainsi le type de compromis défini par la composante moyenne \mathbf{XMa}_j .

De telles « variances expliquées » définissent des parts d'inertie impliquées dans la dépendance mesurée par des corrélations linéaires. Ces variances expliquées ne traduisent des parts de variance des variables des tableaux que lorsque les métriques sont métriques identité. L'égalité globale suivante, qui se déduit des précédentes, traduit ainsi une décomposition de l'importance relative des tableaux \mathbf{X}_i quant à la dépendance simultanée avec \mathbf{Y} :

$$\sum_{j=1}^S \rho_j^2 \text{var}(\mathbf{XMa}_j) = \sum_{i=1}^K \sum_{j=1}^S \rho_{ij}^2 \text{var}(\mathbf{X}_i \mathbf{M}_i \mathbf{a}_{ij}).$$

Les K cumulés :

$$\sum_{j=1}^S \rho_{ij}^2 \text{var}(\mathbf{X}_i \mathbf{M}_i \mathbf{a}_{ij}), i = 1, \dots, K,$$

peuvent donc être plus judicieux à considérer que les cumulés de covariances pour désigner le tableau jugé le plus important. Un résultat de simulation faite dans un contexte de reconnaissance de structure va dans ce sens (Lafosse 1999b).

Cependant, la distinction entre critères qui sous-tendent un découpage, et mesures par les variances expliquées, est inutile en analyse canonique ou en analyse de type ACPVI, puisqu'alors les covariances sont elles-mêmes des variances expliquées (ces contextes particuliers étant introduits en choisissant pour métriques des métriques de Mahalanobis).

Les écritures avec des métriques quelconques, permettent de définir l'analyse CONCOR dans toute sa généralité. Cependant métriques identité et (semi-)métriques de Mahalanobis constituent les choix les plus évidents dans ce contexte. Une autre façon de supprimer l'information sur la structure interne à un tableau, que par le choix de Mahalanobis, consiste à remplacer les variables initiales par un système de variables non corrélées et de variances 1, engendrant cependant le même sous-espace que les variables initiales. Ce qui avait été relevé par Ten Berge (1977, 1988) et Van de Geer (1984), pour l'analyse SUMCOR de Horst (1961), reste valide en analyse CONCOR. Un programme d'analyses CONCOR correspondant seulement aux deux types de métriques indiqués peut donc être simplifié, puisque finalement, avec cette substitution faite au départ et aisément produite depuis le calcul des composantes principales réduites d'ACP, seules les métriques identité peuvent alors être considérées.

2.2. Couples monogames

On définit maintenant des couples de composantes qui sont appelés *couples monogames* de l'analyse CONCOR. Ce sont eux qui réalisent le découpage de relation entre un tableau et K tableaux. La notion de couples monogames a été proposée en analyse de concordance de deux tableaux (Lafosse, 1997). Plus généralement en analyse CONCOR, les couples monogames définis sont associés à un couple de tableaux obtenu après le calcul d'un projecteur.

Les corrélations croisées suivantes sont nulles pour tous les indices h et j , $h \neq j$, de $[1, S]$:

$$\rho(\mathbf{XMa}_h, \mathbf{YNb}_j) = 0. \quad (8)$$

Les S couples $(\mathbf{XMa}_j, \mathbf{YNb}_j)$ sont donc tels qu'une composante principale \mathbf{YNb}_j n'est corrélée qu'à sa composante partenaire compromis \mathbf{XMa}_j , dans la mesure où elle est de corrélation nulle avec les $S - 1$ autres composantes partenaires compromis, tout comme sont nulles les corrélations de \mathbf{XMa}_j avec les $S - 1$ composantes \mathbf{YNb}_h , $h \neq j$.

De plus le système N -orthonormé $\{\mathbf{b}_j\}$ induit un découpage de l'inertie du nuage (\mathbf{Y}, \mathbf{N}) alors que le système M -orthonormé $\{\mathbf{a}_j\}$ induit un découpage de l'inertie du nuage (\mathbf{X}, \mathbf{M}) .

Par suite on considère que les S couples $(\mathbf{XMa}_j, \mathbf{YNb}_j)$ sont isolables les uns des autres et constituent les couples monogames de l'analyse. Ces couples vérifient les relations suivantes :

$$\mathbf{Y}'\mathbf{DXMa}_j = \lambda_j \mathbf{b}_j, \quad (9)$$

$$\mathbf{YNY}'\mathbf{DXMa}_j = \lambda_j \mathbf{YNb}_j.$$

Mais, excepté le cas où $K = 1$, dès que $j > 1$ et à cause des déflations à chaque pas des tableaux \mathbf{X}_j , on a :

$$\mathbf{X}'\mathbf{DYNb}_j \neq \lambda_j \mathbf{a}_j,$$

de sorte que les couples $(\mathbf{XMa}_j, \mathbf{YNb}_j)$ ne sont pas les couples monogames des deux nuages (\mathbf{Y}, \mathbf{N}) et (\mathbf{X}, \mathbf{M}) .

On note \mathbf{B} la matrice $q \times S$ constituée des vecteurs colonnes \mathbf{b}_j et \mathbf{A} la matrice $m \times S$ constituée des vecteurs colonnes respectifs \mathbf{a}_j . On note Δ la matrice diagonale $S \times S$ contenant sur la diagonale les valeurs respectives $\lambda_j, j = 1, \dots, S$. D'après (9), on a :

$$\mathbf{Y}'\mathbf{DXMA} = \mathbf{B}\Delta. \quad (10)$$

La matrice $\mathbf{AA}'\mathbf{M}$ est celle du projecteur sur le sous-espace de dimension S de \mathcal{R}^m impliqué dans la dépendance simultanée avec \mathbf{Y} . La relation établie au départ

entre Y et X est restreinte maintenant à une relation entre les individus de Y et ceux de X projetés sur ce sous-espace.

Propriété fondamentale 2.2. – *Les couples (XMa_j, YNb_j) sont les couples monogames de l'analyse de concordance des triplés statistiques (Y, N, D) et $(XMAA', M, D)$.*

Comme $AA'Ma_j = a_j$, la relation (9) peut encore s'écrire :

$$Y'DXMAA'Ma_j = \lambda_j b_j, \quad (11)$$

alors que, d'après (10) et sachant l'orthogonalité des vecteurs b_j , on a

$$(XMAA')'DYNb_j = AA'MX'DYNb_j = A\Delta B'Nb_j = \lambda_j a_j. \quad (12)$$

Cela prouve que les S couples de vecteurs (b_j, a_j) sont ceux associés aux S valeurs singulières de la décomposition en valeurs singulières de la matrice :

$$Y'DXMAA'.$$

L'analyse sera désormais conduite depuis le découpage de la relation entre les tableaux Y et $XMAA'$, découpage défini par les S couples (XMa_j, YNb_j) .

Par construction, le nombre de solutions S est donc inférieur ou égal au rang de $Y'DX$. Quand aucun partenaire n'est nul, $S = \sup(\text{rang}(Y'DX_j))$.

3. Régressions simultanées

3.1. Introduction

Utilisant par la suite des notions notamment définies dans Lafosse (1997), cette introduction en constitue un rappel original produit depuis un problème introduit par Green (1952), plus connu sous le nom de problème de Cliff (1966).

Soient X $n \times m$ et Y $n \times q$ deux matrices centrées. Pour obtenir un ajustement de Y à X , Cliff pose le problème (pour $m = q$) de la recherche d'une transformation R , application de \mathcal{R}^q dans \mathcal{R}^m vérifiant $R'R = I_q$ et maximisant le critère

$$f_1(R) = \text{tr}[X'DYR'].$$

Soit $U \Delta V'$ une décomposition en valeurs singulières usuelle de $X'DY$.

Une solution est $R' = VU'$ (Cliff, et plus de détails quand $m \neq q$ dans Ten Berge, 1977). L'intensité de la relation entre X et le tableau transformé pourrait alors s'évaluer par

$$\text{tr}(X'DYVU') = \text{tr}[U \Delta U'] = \text{tr}[\Delta],$$

mais le critère f_1 demande à être relativisé.

Le problème de Cliff peut ainsi être vu comme celui de la maximisation d'un critère g , équivalent au précédent quand $\mathbf{R}'\mathbf{R} = \mathbf{I}_q$:

$$g_1(\mathbf{R}) = \cos_D^2(\mathbf{X}, \mathbf{YR}') = \text{tr}^2[\mathbf{X}'\mathbf{D}\mathbf{YR}'] / \{\text{tr}[\mathbf{X}'\mathbf{D}\mathbf{X}]\text{tr}[(\mathbf{YR}')'\mathbf{D}\mathbf{YR}']\}.$$

La valeur du critère g_1 obtenue au maximum est indépendante de toute rotation orthogonale préliminaire faite sur \mathbf{Y} .

On suppose maintenant que le problème vise à définir une mesure du lien entre \mathbf{X} et \mathbf{Y} pouvant mieux caractériser la participation de \mathbf{X} à ce lien, parce qu'elle est aussi indépendante d'affinités orthogonales faites sur \mathbf{Y} . La solution de Cliff ne serait alors qu'une première étape à franchir (cette première étant en fait déjà une seconde, le centrage initial des tableaux revenant à effectuer une famille de translations optimale au sens des moindres carrés). En effet, réaliser des affinités orthogonales dans les directions des axes définis par les colonnes de \mathbf{U} ne modifie pas le problème de Cliff et sa solution. La question de ces affinités peut donc se traiter une fois la rotation de Cliff adoptée. On considère ainsi la famille des transformations $\mathbf{C}' = \mathbf{V}\mathbf{U}'\mathbf{U}\Delta_C\mathbf{U}'$ avec le critère à optimiser

$$h(\Delta_C) = \cos^2(\mathbf{X}, \mathbf{Y}\mathbf{C}').$$

Il s'agit alors de trouver une matrice diagonale Δ_C à éléments diagonaux strictement positifs maximisant

$$h(\Delta_C) = \text{tr}^2(\Delta\Delta_C) / \{\text{tr}(\mathbf{X}'\mathbf{D}\mathbf{X})\text{tr}(\Delta_Y\Delta_C^2)\}. \quad (*)$$

Dans l'expression précédente la matrice diagonale Δ_Y contient sur la diagonale les valeurs $\text{var}(\mathbf{Y}\mathbf{v}_j)$, alors que la matrice diagonale Δ contient sur la diagonale les valeurs $\text{cov}(\mathbf{X}\mathbf{u}_j, \mathbf{Y}\mathbf{v}_j)$.

Soit \mathbf{v} le vecteur dont les composantes sont les éléments diagonaux de Δ_C , \mathbf{u} le vecteur dont les composantes sont les éléments diagonaux de Δ .

La quantité (*) est équivalente à la quantité

$$\frac{\langle \mathbf{u}, \mathbf{v} \rangle^2}{\langle \mathbf{v}, \mathbf{v} \rangle_{\Delta_Y}}.$$

Maximiser le critère h revient ainsi à maximiser

$$\frac{\langle \mathbf{u}, \mathbf{v} \rangle^2}{\langle \mathbf{v}, \mathbf{v} \rangle_{\Delta_Y}}.$$

On pose $\tilde{\mathbf{v}} = \Delta_Y^{-\frac{1}{2}}\mathbf{v}$.

Sous contrainte que $\tilde{\mathbf{v}}$ soit normé, le problème revient à maximiser le critère g_2

$$g_2(\tilde{\mathbf{v}}) = \tilde{\mathbf{v}}'\Delta_Y^{-\frac{1}{2}}\mathbf{u}\mathbf{u}'\Delta_Y^{-\frac{1}{2}}\tilde{\mathbf{v}} = (\mathbf{u}'\Delta_Y^{-\frac{1}{2}}\tilde{\mathbf{v}})^2.$$

Une solution est alors obtenue pour

$$\tilde{\mathbf{v}} = \Delta_Y^{-\frac{1}{2}} \mathbf{u} / \|\Delta_Y^{-\frac{1}{2}} \mathbf{u}\|.$$

Ainsi $\mathbf{v} = \Delta_Y^{-1} \mathbf{u} / \|\Delta_Y^{-\frac{1}{2}} \mathbf{u}\|$ et comme le problème est invariant par homothétie, $\mathbf{v} = \Delta_Y^{-1} \mathbf{u}$ est aussi une solution, de sorte que $\Delta_C = \Delta_Y^{-1} \Delta$ est une solution de maximisation du critère h .

La matrice Δ_C a donc sur la diagonale les coefficients des régressions simples des composantes $\mathbf{X}\mathbf{u}_j$ sur $\mathbf{Y}\mathbf{v}_j$, $j = 1, \dots, r$.

La matrice transformée $\mathbf{Y}\mathbf{V}\Delta_Y^{-1}\Delta\mathbf{U}'$ caractérise maintenant d'une façon plus indépendante de \mathbf{Y} la participation de \mathbf{X} au lien de dépendance établi entre \mathbf{X} et \mathbf{Y} . C'est une matrice qui a été définie la première fois en terme d'ajustement entre nuages d'individus dans Lafosse (1985) et qui a été nommée image concordante dans Lafosse (1997).

Pour la solution $\Delta_C = \Delta_Y^{-1} \Delta$ le critère h devient

$$h(\mathbf{C}) = \cos^2(\mathbf{X}, \mathbf{Y}\mathbf{V}\Delta_Y^{-1}\Delta\mathbf{U}') = \text{tr}(\Delta^2 \Delta_Y^{-1}) / \text{tr}[\mathbf{X}'\mathbf{D}\mathbf{X}].$$

Ce rapport entre somme de variances expliquées et variance totale qui traduit la quantité relative d'inertie de \mathbf{X} impliquée dans la liaison entre \mathbf{X} et \mathbf{Y} quand les métriques sont les métriques identité avait été nommé LAI(\mathbf{X} / \mathbf{Y}), pour « Linear Agreement Index ». Chaque variance expliquée est décomposable en contributions partielles, contributions des individus de \mathbf{X} à la dépendance due à la corrélation linéaire au carré correspondante.

Depuis la régression simultanée de \mathbf{Y} sur \mathbf{X} , on peut aussi définir les contributions partielles des individus de \mathbf{Y} , contributions à la dépendance due à la même corrélation. Ces deux ensembles de contributions conduisent à associer deux biplots pour définir des quadriplots (Lafosse, 1999a).

L'analyse CONCOR est une généralisation des analyses factorielles de deux tableaux, que l'on retrouve pour $K = 1$. Quand la généralisation induit des problèmes particuliers, on situe par la suite cette particularité par rapport au cas « $K = 1$ ». Sinon, les démonstrations déjà produites pour $K = 1$ ne sont pas reprises.

3. 2. Régression de \mathbf{Y} sur les tableaux \mathbf{X}_j

3.2.1. Image concordante

Sans que de nouvelles justifications soient nécessaires puisqu'on s'est ramené dans la section 2.2 à une dépendance entre les deux triplés $(\mathbf{Y}, \mathbf{N}, \mathbf{D})$ et $(\mathbf{X}\mathbf{M}\mathbf{A}\mathbf{A}', \mathbf{M}, \mathbf{D})$, on définit une image concordante \mathbf{Y}_A depuis des régressions simples simultanées considérées entre composantes des couples monogames $(\mathbf{X}\mathbf{M}\mathbf{a}_j, \mathbf{Y}\mathbf{N}\mathbf{b}_j)$, $j = 1, \dots, S$.

On note Δ_b la matrice diagonale des S coefficients de régression, c'est-à-dire des valeurs :

$$\text{cov}(\mathbf{Y}\mathbf{N}\mathbf{b}_j, \mathbf{X}\mathbf{M}\mathbf{a}_j)/\text{var}(\mathbf{X}\mathbf{M}\mathbf{a}_j). \quad (13)$$

Ramenant à 1 les variances des composantes compromis $\mathbf{X}\mathbf{M}\mathbf{a}_j$ on range par colonnes ces composantes, une fois réduites, dans une matrice nommée \mathbf{K} . La régression simple simultanée de $\mathbf{Y}\mathbf{N}\mathbf{B}$ sur $\mathbf{X}\mathbf{M}\mathbf{A}$ conduit à définir l'image de \mathbf{Y} concordante avec les tableaux \mathbf{X}_i par la matrice $n \times q$:

$$\mathbf{Y}_A = \mathbf{X}\mathbf{M}\mathbf{A}\Delta_b\mathbf{B}' = \mathbf{K}\mathbf{K}'\mathbf{D}\mathbf{Y}. \quad (14)$$

L'indice global LAI mesure l'intensité de la participation de (\mathbf{Y}, \mathbf{N}) à la dépendance simultanée :

$$\begin{aligned} \text{LAI}[(\mathbf{Y}, \mathbf{N})/(\mathbf{X}_i, \mathbf{M}_i)] &= \text{tr}(\mathbf{Y}'_A\mathbf{D}\mathbf{Y}_A\mathbf{N})/\text{tr}(\mathbf{Y}'\mathbf{D}\mathbf{Y}\mathbf{N}) \\ &= \sum_{i=1}^S \rho_j^2 \text{var}(\mathbf{Y}\mathbf{N}\mathbf{b}_j)/\text{tr}(\mathbf{Y}'\mathbf{D}\mathbf{Y}\mathbf{N}). \end{aligned} \quad (15)$$

Aucune difficulté particulière n'apparaît quand $K > 1$ pour définir les contributions partielles des individus et les participations partielles des variables du tableau \mathbf{Y} (par $\text{diag}(\mathbf{Y}'\mathbf{D}\mathbf{Y}_A\mathbf{N}\mathbf{b}_j\mathbf{b}'_j)$).

3.2.2. Anti-image et discordance

Pour $K = 1$, le résidu de régression simple simultanée est décomposée en deux parties, l'anti-image qui a été définie par Guttman et la discordance qui traduit la non adéquation entre systèmes de covariations internes. Notant

$$\mathbf{P}_{\mathbf{X}\mathbf{M}\mathbf{A}} = \mathbf{X}\mathbf{M}\mathbf{A}[(\mathbf{X}\mathbf{M}\mathbf{A})'\mathbf{D}\mathbf{X}\mathbf{M}\mathbf{A}]^{-1}(\mathbf{X}\mathbf{M}\mathbf{A})'\mathbf{D}$$

le projecteur sur le sous-espace engendré par les colonnes de $\mathbf{X}\mathbf{M}\mathbf{A}$, cette anti-image est ici définie par

$$\mathbf{Y}_N = \mathbf{Y} - \mathbf{P}_{\mathbf{X}\mathbf{M}\mathbf{A}}\mathbf{Y}. \quad (16)$$

La matrice \mathbf{Y}_D étant l'image discordante définie par la suite, les trois matrices nulles suivantes

$$\mathbf{Y}'_N\mathbf{D}\mathbf{Y}_A = \mathbf{Y}'_N\mathbf{D}\mathbf{Y}_D = \mathbf{Y}'_N\mathbf{D}\mathbf{X}\mathbf{M}\mathbf{A}\mathbf{A}' = \mathbf{0}, \quad (17)$$

rapellent que l'anti-image définie est indépendante des matrices impliquées dans la dépendance simultanée.

L'image discordante est définie par :

$$\mathbf{Y}_D = \mathbf{P}_{\mathbf{X}\mathbf{M}\mathbf{A}}\mathbf{Y} - \mathbf{Y}_A. \quad (18)$$

L'existence de cette image provient de la prise en compte de la structure de \mathbf{X} . On nomme ici structure, l'ensemble des liens de corrélations entre les variables de \mathbf{X} . Considérant que c'est celle captée dans l'ACP de \mathbf{X} par les composantes principales, la structure de \mathbf{X} est donc invariante par isométries.

Quand toutes les métriques sont les métriques identité, la matrice \mathbf{Y}_D traduit une différence de structure interne, induite par l'existence de liens internes dans chacun des \mathbf{K} tableaux différents de ceux de \mathbf{Y} , alors que \mathbf{Y}_A traduit l'analogie de la structure du nuage (\mathbf{Y}, \mathbf{N}) avec les structures des nuages $(\mathbf{X}_j, \mathbf{M}_j)$.

Quand $K = 1$, alors $\mathbf{P}_{XMA}\mathbf{Y} = \mathbf{P}_X\mathbf{Y}$ (Lafosse, 1997), et pour un axe j de l'analyse (interprétable depuis les biplots de concordance), l'image discordante :

$$\mathbf{Y}_D\mathbf{N}b_j = \mathbf{P}_{XMA}\mathbf{Y}\mathbf{N}b_j - \mathbf{Y}_A\mathbf{N}b_j$$

a, dans le repère des axes de l'analyse, pour coordonnées :

$$\begin{aligned}\mathbf{Y}_D\mathbf{N}b_j &= \mathbf{X}(\mathbf{X}'\mathbf{D}\mathbf{X})^{-1}\mathbf{X}'\mathbf{D}\mathbf{Y}\mathbf{N}b_j - \mathbf{XMA}\Delta_b\mathbf{B}'\mathbf{N}b_j \\ &= \lambda_j\mathbf{X}(\mathbf{X}'\mathbf{D}\mathbf{X})^{-1}\mathbf{a}_j - \mathbf{XMA}\Delta_b\mathbf{b}_{\delta j}\end{aligned}$$

où $\mathbf{b}_{\delta j}$ est vecteur colonne nul excepté le 1 en position j . Finalement

$$\mathbf{Y}_D\mathbf{N}b_j = \lambda_j[\mathbf{X}(\mathbf{X}'\mathbf{D}\mathbf{X})^{-1}\mathbf{a}_j - \mathbf{XMA}_j/||\mathbf{XMA}_j||^2].$$

Pour $K = 1$ quand les deux métriques sont les métriques identité, la discordance est nulle si les tableaux ont la même structure interne relativement à un axe j (alors $\mathbf{a}_j = \mathbf{b}_j$ est un axe principal de l'ACP de \mathbf{X} et on vérifie bien que $\mathbf{Y}_D\mathbf{b}_j = \mathbf{0}$). Quand $K > 1$, quand toutes les métriques sont les métriques identité, si un axe j est axe d'ACP pour chacun des tableaux respectifs, alors on vérifie facilement que la composante principale $\mathbf{Y}\mathbf{b}_j$ de \mathbf{Y} est proportionnelle à la composante principale $\mathbf{X}\mathbf{a}_j$ de \mathbf{X} .

D'après (14), on a alors

$$\mathbf{P}_{XA}\mathbf{Y}\mathbf{b}_j = \mathbf{P}_X\mathbf{Y}\mathbf{b}_j,$$

car la projection globale est totalement et uniquement réalisée par la projection partielle $\mathbf{X}\mathbf{a}_j(\mathbf{X}\mathbf{a}_j)'\mathbf{D}$ sur la composante principale $\mathbf{Y}\mathbf{b}_j$. Par suite $\mathbf{Y}_D\mathbf{b}_j = \mathbf{0}$.

Avec $K = 1$ et \mathbf{M} la métrique de Mahalanobis, $||\mathbf{XMA}_j|| = 1$, et la discordance est encore nulle (Lafosse, 1997). La différence d'une structure interne avec une autre structure ne peut s'établir quand l'information sur cette dernière a été supprimée.

Quand $\mathbf{P}_{XMA}\mathbf{Y} \neq \mathbf{P}_X\mathbf{Y}$ (ce qui se produit donc en général si $K \neq 1$) quand les métriques \mathbf{M}_j sont celles respectives de Mahalanobis, la discordance peut être non nulle. Il ne suffit donc pas cette fois de supprimer chacune des structures internes aux tableaux \mathbf{X}_j pour trouver une discordance \mathbf{Y}_D nulle. La suppression de chaque structure interne n'équivaut pas ici à la suppression de toute notion de structure sur \mathbf{X} . Cette discordance provient d'une différence entre les sous espaces vectoriels

engendrés pour chaque i par les variables respectives $\{X_i M_i a_{ij} \mid j = 1, \dots, S\}$. Cela se produit quand les dimensions de ces sous-espaces sont différents.

3.3. Régressions des tableaux X_i sur Y

Dans le cas de deux tableaux Y et X_1 ($K = 1$), la participation de Y à la dépendance entre Y et X_1 se définit comme se définit la participation de X_1 . Cette répétition de définition n'est utile que pour juger de situations particulières créées par l'emploi de métriques particulières, et n'est pas utile dans le contexte général de métriques quelconques.

Une situation nouvelle est créée dans le cas de la concordance entre un tableau Y et plusieurs tableaux X_i ($K > 1$), l'apport spécifique de chaque tableau X_i à la dépendance simultanée devant être ici précisé.

3.3.1. Images concordantes synthétiques et partielles

On continue de développer l'analyse CONCOR en se fondant sur les couples monogames (XMA_j, YNB_j) , pour obtenir cette fois les contributions de chaque tableau X_i aux trois notions de concordance, discordance et anti-image.

Les régressions simples simultanées maintenant considérées permettent de définir une matrice synthétisant les K images concordantes relatives aux K tableaux X_i .

On note D_a la matrice diagonale des S coefficients de régression de la régression simple simultanée de XMA sur YNB , c'est-à-dire des valeurs :

$$\text{cov}(YNB_j, XMA_j) / \text{var}(YNB_j). \quad (19)$$

Vu qu'une composante XMA_j constitue un compromis entre les composantes partenaires du $K + 1$ uple j , on se sert de la régression simple de XMA_j sur YNB_j pour synthétiser sur la dimension j la participation à la concordance des tableaux X_i .

On note L la matrice ayant en colonnes les composantes YNB_j , après les avoir cependant réduites. Les S régressions simultanément réalisées définissent ainsi une matrice synthétique $n \times m$:

$$X_A = YNB D_a A' = LL' D X M A A'. \quad (20)$$

Cette matrice est la synthèse d'images partielles X_{A_i} traduisant les contributions de chaque tableau à la concordance avec Y . Ces images partielles sont les K sous-blocs de X_A de dimensions respectives $n \times m_i$. La propriété suivante montre que ces images partielles peuvent être obtenues à partir de régressions simples simultanées, sans passer par l'image synthétique.

On note A_j la matrice $m_i \times S$ ayant pour colonnes les vecteurs normés a_{ij} , $j = 1, \dots, S$.

Propriété 3.3.1. — *Les K matrices obtenues par régression simple simultanée de $X_i M_i A_i$ sur YNB , $i = 1, \dots, K$, sont les images concordantes partielles X_{A_i} constituant l'image concordante synthétique X_A .*

On note par \mathbf{P}_i la matrice diagonale contenant sur la diagonale les normes $p_{i,j}$ des sous-vecteurs $p_{i,j}\mathbf{a}_{ij}$ de $\mathbf{a}_j, j = 1, \dots, S$. D'après (20), l'image partielle exprimant la concordance de \mathbf{X}_i avec \mathbf{Y} est définie par :

$$\mathbf{X}_{Ai} = \mathbf{Y}\mathbf{N}\mathbf{B}\mathbf{D}_a\mathbf{P}_i\mathbf{A}'_i = \mathbf{Y}\mathbf{N}\mathbf{B}\mathbf{D}_{ai}\mathbf{A}'_i. \quad (21)$$

Or la matrice diagonale $\mathbf{D}_{ai} = \mathbf{D}_a\mathbf{P}_i$ contient les valeurs :

$$[\text{cov}(\mathbf{Y}\mathbf{N}\mathbf{b}_j, \mathbf{X}\mathbf{M}\mathbf{a}_j)/\text{var}(\mathbf{Y}\mathbf{N}\mathbf{b}_j)] \times [\text{cov}(\mathbf{X}_i\mathbf{M}_i\mathbf{a}_{ij}, \mathbf{Y}\mathbf{N}\mathbf{b}_j)/\text{cov}(\mathbf{Y}\mathbf{N}\mathbf{b}_j, \mathbf{X}\mathbf{M}\mathbf{a}_j)]. \quad (22)$$

Ce sont donc les coefficients de régression linéaire simple de $\mathbf{X}_i\mathbf{M}_i\mathbf{a}_{ij}$ sur $\mathbf{Y}\mathbf{N}\mathbf{b}_j$:

$$\text{cov}(\mathbf{X}_i\mathbf{M}_i\mathbf{a}_{ij}, \mathbf{Y}\mathbf{N}\mathbf{b}_j)/\text{var}(\mathbf{Y}\mathbf{N}\mathbf{b}_j). \quad (23)$$

Donc tout se passe comme si on avait fait ce calcul de l'image partielle en partant d'une régression simple simultanée basée sur les S couples $(\mathbf{X}_i\mathbf{M}_i\mathbf{a}_{ij}, \mathbf{Y}\mathbf{N}\mathbf{b}_j)$.

Mais ce faisant directement, notre démarche n'aurait pu être justifiée, car n'étant pas monogames, ces couples ne peuvent être considérés séparément. En effet, $\mathbf{Y}\mathbf{N}\mathbf{b}_j$ est corrélée aux composantes $\mathbf{X}_i\mathbf{M}_i\mathbf{a}_{ih}$, pour $h < j$. Par suite, on retrouve ci-après pour les images concordantes partielles les propriétés usuelles des images concordantes.

3.3.2 Décomposition des contributions de \mathbf{X}_i

On pose $\mathbf{X}_{Ti} = \mathbf{X}_i - \mathbf{X}_{Ai}$ et on considère la décomposition :

$$\mathbf{X}_i = \mathbf{X}_{Ai} + \mathbf{X}_{Ti}. \quad (24)$$

Propriété 3.3.2. — *L'inertie totale de \mathbf{X}_i se décompose en somme de deux contributions totales, celle de \mathbf{X}_{Ai} et celle de \mathbf{X}_{Ti} . La somme des contributions relatives à la concordance est égale à la somme des variances expliquées de la régression simple simultanée qui définit l'image partielle concordante \mathbf{X}_{Ai} .*

D'après (21), pour chaque i on a :

$$\text{tr}(\mathbf{X}'_i\mathbf{D}\mathbf{X}_{Ai}\mathbf{M}_i) = \text{tr}(\mathbf{X}'_i\mathbf{D}\mathbf{Y}\mathbf{N}\mathbf{B}\mathbf{D}_{ai}\mathbf{A}'_i\mathbf{M}_i) = \text{tr}[(\mathbf{X}_i\mathbf{M}_i\mathbf{A}_i)'\mathbf{D}\mathbf{Y}\mathbf{N}\mathbf{B}]\mathbf{D}_{ai}. \quad (25)$$

Soit encore, d'après (22) et (23) :

$$\text{tr}(\mathbf{X}'_i\mathbf{D}\mathbf{X}_{Ai}\mathbf{M}_i) = \sum_{j=1}^S \text{cov}^2(\mathbf{X}_i\mathbf{M}_i\mathbf{a}_{ij}, \mathbf{Y}\mathbf{N}\mathbf{b}_j)/\text{var}(\mathbf{Y}\mathbf{N}\mathbf{b}_j). \quad (26)$$

On a donc :

$$\text{tr}(\mathbf{X}'_i\mathbf{D}\mathbf{X}_{Ai}\mathbf{M}_i) = \sum_{j=1}^S \rho_{ij}^2 \text{var}(\mathbf{X}_i\mathbf{M}_i\mathbf{a}_{ij}), \quad (27)$$

et aussi :

$$\begin{aligned}\text{tr}(\mathbf{X}'_{Ai} \mathbf{D} \mathbf{X}_{Ai} \mathbf{M}_i) &= \text{tr}[(\mathbf{Y} \mathbf{N} \mathbf{B} \mathbf{D}_{ai} \mathbf{A}'_i)' \mathbf{D} \mathbf{Y} \mathbf{N} \mathbf{B} \mathbf{D}_{ai} \mathbf{A}'_i \mathbf{M}_i] \\ &= \text{tr}[\mathbf{D}_{ai} (\mathbf{Y} \mathbf{N} \mathbf{B})' \mathbf{D} \mathbf{Y} \mathbf{N} \mathbf{B} \mathbf{D}_{ai}],\end{aligned}\quad (28)$$

Et ainsi on a :

$$\text{tr}(\mathbf{X}'_{Ai} \mathbf{D} \mathbf{X}_{Ai} \mathbf{M}_i) = \sum_{j=1}^S \rho_{ij}^2 \text{var}(\mathbf{X}_i \mathbf{M}_i \mathbf{a}_{ij}). \quad (29)$$

De (27) et (29) on déduit alors :

$$\text{tr}(\mathbf{X}'_i \mathbf{D} \mathbf{X}_{Ai} \mathbf{M}_i) = \text{tr}(\mathbf{X}'_{Ai} \mathbf{D} \mathbf{X}_{Ai} \mathbf{M}_i). \quad (30)$$

Comme d'après (24) on a :

$$\mathbf{X}'_{Ai} \mathbf{D} \mathbf{X}_i \mathbf{M}_i = \mathbf{X}'_{Ai} \mathbf{D} \mathbf{X}_{Ai} \mathbf{M}_i + \mathbf{X}'_{Ai} \mathbf{D} \mathbf{X}_{Ti} \mathbf{M}_i, \quad (31)$$

on obtient :

$$\text{tr}(\mathbf{X}'_{Ai} \mathbf{D} \mathbf{X}_{Ti} \mathbf{M}_i) = 0. \quad (32)$$

En usant encore de (24), on obtient alors :

$$\text{tr}(\mathbf{X}'_{Ti} \mathbf{D} \mathbf{X}_{Ti} \mathbf{M}_i) = \text{tr}(\mathbf{X}'_i \mathbf{D} \mathbf{X}_{Ti} \mathbf{M}_i). \quad (33)$$

Comme d'après (24) on a :

$$\mathbf{X}'_i \mathbf{D} \mathbf{X}_i \mathbf{M}_i = \mathbf{X}'_i \mathbf{D} \mathbf{X}_{Ai} \mathbf{M}_i + \mathbf{X}'_i \mathbf{D} \mathbf{X}_{Ti} \mathbf{M}_i. \quad (34)$$

Finalement pour chaque i , l'égalité (24) est relative à une décomposition de l'inertie

$$\text{tr}(\mathbf{X}'_i \mathbf{D} \mathbf{X}_i \mathbf{M}_i) = \text{tr}(\mathbf{X}'_{Ai} \mathbf{D} \mathbf{X}_{Ai} \mathbf{M}_i) + \text{tr}(\mathbf{X}'_{Ti} \mathbf{D} \mathbf{X}_{Ti} \mathbf{M}_i). \quad (35)$$

Si tout semble donc se passer comme en concordance d'un tableau avec un autre tableau, on remarque cependant qu'ici les matrices $\mathbf{X}'_i \mathbf{D} \mathbf{X}_{Ai}$ et $\mathbf{X}'_i \mathbf{D} \mathbf{X}_{Ti}$ ne sont pas symétriques.

En fait, beaucoup de décompositions en somme de type (24) mènent à un découpage d'inertie de type (35). Mais celui proposé a été justifié depuis le découpage initial de relation fondé sur les couples monogames de l'analyse CONCOR.

Les participations globales à la concordance simultanée des variables de \mathbf{X} s'évaluent par les termes diagonaux :

$$\text{diag}(\mathbf{X}' \mathbf{D} \mathbf{X}_A) = \text{diag}(\mathbf{X}' \mathbf{D} \mathbf{L} \mathbf{L}' \mathbf{D} \mathbf{X} \mathbf{M} \mathbf{A} \mathbf{A}'). \quad (36)$$

Quand $\mathbf{X} = \mathbf{XMAA}'$, les participations se décomposent par axe, leurs valeurs positives étant pour un axe j égales à (Lafosse, 1997) :

$$\text{diag}(\mathbf{AA}'\mathbf{MX}'\mathbf{DLL}'\mathbf{DXMa}_j\mathbf{a}'_j) = \rho_j^2 \text{var}(\mathbf{XMa}_j) \text{diag}(\mathbf{a}_j\mathbf{a}'_j). \quad (37)$$

Mais quand $\mathbf{X} \neq \mathbf{XMAA}'$, des valeurs négatives apparaissent, correspondant à quelques covariances négatives entre les colonnes de \mathbf{X} et celles de \mathbf{X}_A . Cela traduit une difficulté à être capable d'établir aussi clairement que l'on voudrait la dépendance simultanée de \mathbf{Y} avec chaque tableau \mathbf{X}_j : les signes des vecteurs \mathbf{a}_{ij} sont bien, par construction, fonctions du signe du vecteur \mathbf{b}_j de façon à retrouver les individus des tableaux respectifs projetés sur les axes respectifs dans des dispositions analogues. Mais cela étant, le découpage correspondant de la variabilité de \mathbf{X} ne peut être toujours parfait comme pour $K = 1$. Cependant les valeurs négatives sont souvent assez faibles parce que souvent associées à des axes porteurs d'une faible inertie et parce que les dimensions supplémentaires aux dimensions \mathbf{a}_{ij} leur étant orthogonales, la correction qu'elles induisent est nulle en moyenne pour chaque axe j . On tient compte de ces signes négatifs sur les biplots en modifiant alors les valeurs des signes correspondants des cosinus directeurs des flèches (\mathbf{a}_j devenant ainsi \mathbf{d}_j) sans changer les coordonnées des individus de sorte que les participations deviennent égales à

$$\text{diag}(\mathbf{X}'\mathbf{DLL}'\mathbf{DXMa}_j\mathbf{d}'_j) = \text{abs}[\text{diag}(\mathbf{X}'\mathbf{DLL}'\mathbf{DXMa}_j\mathbf{a}'_j)].$$

Cela induit une redondance des participations par rapport aux inerties, qui traduit l'idée qu'il est nécessaire de « se répéter » pour arriver à prendre en compte toutes les relations qu'ont les variables avec les contributions des individus.

3.3.3. Discordances et anti-images. Indices.

L'anti-image relative à \mathbf{X}_j définit la part de contributions de \mathbf{X}_j n'ayant aucune influence dans l'analyse parce que non corrélée à \mathbf{Y} .

Notant $\mathbf{P}_Y = \mathbf{Y}(\mathbf{Y}'\mathbf{D}\mathbf{Y})^{-1}\mathbf{Y}'\mathbf{D}$ le projecteur sur le sous-espace engendré par toutes les colonnes de \mathbf{Y} , on définit ainsi les anti-images partielles (relatives aux tableaux \mathbf{X}_j) par :

$$\mathbf{X}_{Ni} = \mathbf{X}_i - \mathbf{P}_Y\mathbf{X}_i. \quad (38)$$

Ces anti-images décomposent par blocs l'anti-image synthétique

$$\mathbf{X}_N = \mathbf{X} - \mathbf{P}_Y\mathbf{X}. \quad (39)$$

Les images concordantes partielles obtenues dans l'analyse de la concordance du tableau \mathbf{Y} avec les tableaux $\mathbf{P}_Y\mathbf{X}_i$ sont encore les images \mathbf{X}_{Ai} précédemment définies. Il suffit de remplacer les \mathbf{X}_i par les tableaux $\mathbf{P}_Y\mathbf{X}_i$ pour le constater, les calculs faits avec la matrice donnée en (2) restant alors inchangés.

Les images partielles discordantes sont alors définies par $\mathbf{X}_{Di} = \mathbf{P}_Y\mathbf{X}_i - \mathbf{X}_{Ai}$ et elles décomposent par blocs l'image discordante synthétique :

$$\mathbf{X}_D = \mathbf{P}_Y\mathbf{X} - \mathbf{X}_A \quad (40)$$

Finalement et comme en concordance d'un tableau avec un autre, on décompose chaque tableau \mathbf{X}_i de dimension $n \times m_i$ en trois tableaux de même dimension :

$$\mathbf{X}_i = \mathbf{X}_{Ai} + \mathbf{X}_{Di} + \mathbf{X}_{Ni}, \quad (41)$$

Cette décomposition correspond bien à un découpage des contributions de \mathbf{X}_i en trois ensembles de contributions, des contributions partielles des individus de \mathbf{X}_i à la dépendance simultanée avec \mathbf{Y} , à la discordance et aux anti-images :

$$\text{tr}(\mathbf{X}'_i \mathbf{D} \mathbf{X}_i \mathbf{M}_i) = \text{tr}(\mathbf{X}'_{Ai} \mathbf{D} \mathbf{X}_{Ai} \mathbf{M}_i) + \text{tr}(\mathbf{X}'_{Di} \mathbf{X}_{Di} \mathbf{M}_i) + \text{tr}(\mathbf{X}'_{Ni} \mathbf{D} \mathbf{X}_{Ni} \mathbf{M}_i). \quad (42)$$

On a finalement :

$$\mathbf{X} = [\mathbf{X}_{A1} \mathbf{X}_{A2} \dots \mathbf{X}_{Ai} \dots \mathbf{X}_{AK}] + [\mathbf{X}_{D1} \mathbf{X}_{D2} \dots \mathbf{X}_{Di} \dots \mathbf{X}_{DK}] + [\mathbf{X}_{N1} \mathbf{X}_{N2} \dots \mathbf{X}_{Ni} \dots \mathbf{X}_{NK}]. \quad (43)$$

Indices partiels et synthétiques

D'après (35) on peut proposer pour indice partiel LAIP de mesure de l'intensité de la participation partielle du tableau $(\mathbf{X}_i, \mathbf{M}_i)$ à la concordance avec (\mathbf{Y}, \mathbf{N}) :

$$\text{LAI}_P[(\mathbf{X}_i, \mathbf{M}_i)/(\mathbf{Y}, \mathbf{N})] = \text{tr}(\mathbf{X}'_{Ai} \mathbf{D} \mathbf{X}_{Ai} \mathbf{M}_i) / \text{tr}(\mathbf{X}'_i \mathbf{D} \mathbf{X}_i \mathbf{M}_i). \quad (44)$$

Le numérateur étant décomposable en somme de variances expliquées (29), rangeant ces variances par ordre décroissant, on peut fonder sur cet ordre l'établissement de résumés concernant la participation du tableau \mathbf{X}_i à la concordance.

D'après (30), cet indice est interprétable comme un cosinus mesurant la proximité entre \mathbf{X}_i et l'image concordante partielle de \mathbf{X}_i avec \mathbf{Y} :

$$\begin{aligned} \cos^2_{(D, M_i)}(\mathbf{X}_i, \mathbf{X}_{Ai}) &= \text{tr}^2(\mathbf{X}'_i \mathbf{D} \mathbf{X}_{Ai} \mathbf{M}_i) / \text{tr}(\mathbf{X}'_i \mathbf{D} \mathbf{X}_i \mathbf{M}_i) \text{tr}(\mathbf{X}'_{Ai} \mathbf{D} \mathbf{X}_{Ai} \mathbf{M}_i) \\ &= \text{LAI}_P[(\mathbf{X}_i, \mathbf{M}_i)/(\mathbf{Y}, \mathbf{N})]. \end{aligned}$$

De la même façon que par (44) et pour chaque tableau, peuvent se définir un indice partiel de discordance et un indice partiel de l'anti-image.

Les indices partiels (44) indiquent les contributions relatives de chaque tableau \mathbf{X}_i dans leur relation de concordance avec \mathbf{Y} .

La moyenne pondérée suivante de ces indices constitue alors une mesure globale de l'importance de l'image concordante synthétique \mathbf{X}_A , puisque l'on a :

$$\begin{aligned} \sum_{i=1}^K \frac{\text{tr}(\mathbf{X}'_i \mathbf{D} \mathbf{X}_i \mathbf{M}_i)}{\text{tr}(\mathbf{X}' \mathbf{D} \mathbf{X} \mathbf{M})} \text{LAI}_P[(\mathbf{X}_i, \mathbf{M}_i)/(\mathbf{Y}, \mathbf{N})] &= \frac{\text{tr}(\mathbf{X}'_A \mathbf{D} \mathbf{X}_A \mathbf{M})}{\text{tr}(\mathbf{X}' \mathbf{D} \mathbf{X} \mathbf{M})} \\ &= \text{LAI}[\mathbf{X} \mathbf{M} \mathbf{A} \mathbf{A}' / \mathbf{M}) / (\mathbf{Y}, \mathbf{N})]. \quad (45) \end{aligned}$$

De même, les moyennes pondérées par les mêmes poids des indices partiels de discordance (resp. des indices partiels de l'anti-image) conduisent à une mesure globale de l'importance de la discordante synthétique (resp. à une mesure globale de l'importance de l'anti-image synthétique).

3.4. Graphiques juxtaposés

Quand plusieurs matrices sont analysées simultanément, les possibilités de graphiques juxtaposés deviennent nombreuses. L'objectif est ici d'apporter des arguments en faveur d'un nombre assez limité d'entre elles.

Par exemple, la représentation des contributions des individus-lignes de \mathbf{X}_{Di} par sous-espace pourrait se faire en réalisant l'ACP du triplet statistique $(\mathbf{X}_{Di}, \mathbf{M}_i, \mathbf{D})$.

Mais il est en général judicieux de juxtaposer les graphiques des trois images dans les repères des axes de concordance, cette juxtaposition étant justifiée par la propriété 3.4 suivante, qui décompose axe par axe la relation (42). L'interprétation des axes est celle qui aura été établie en concordance avec les biplots.

Propriété 3.4. — *L'inertie du nuage $(\mathbf{X}_i, \mathbf{M}_i)$ projeté sur un axe \mathbf{a}_{ij} est égale à la somme des inerties des projetés sur cet axe de chacun des trois nuages $(\mathbf{X}_{Ai}, \mathbf{M}_i)$, $(\mathbf{X}_{Di}, \mathbf{M}_i)$ et $(\mathbf{X}_{Ni}, \mathbf{M}_i)$.*

Le coefficient $\rho_{i,j}$ est la corrélation linéaire entre $\mathbf{X}_{Ai}\mathbf{M}_i\mathbf{a}_{ij}$ et $\mathbf{Y}\mathbf{N}\mathbf{b}_j$. On remarque d'abord que l'on a, d'après (21), (23) et sachant l'orthonormalité des \mathbf{a}_{ij} pour i fixé :

$$\begin{aligned} \mathbf{a}'_{ij}\mathbf{M}_i\mathbf{X}'_i\mathbf{D}\mathbf{X}_{Ai}\mathbf{M}_i\mathbf{a}_{ij} &= \mathbf{a}'_{ij}\mathbf{M}_i\mathbf{X}'_i\mathbf{D}\mathbf{Y}\mathbf{N}\mathbf{b}_j\text{cov}(\mathbf{X}_i\mathbf{M}_i\mathbf{a}_{ij}, \mathbf{Y}\mathbf{N}\mathbf{b}_j)/\text{var}(\mathbf{Y}\mathbf{N}\mathbf{b}_j) \\ &= \text{cov}^2(\mathbf{X}_i\mathbf{M}_i\mathbf{a}_{ij}, \mathbf{Y}\mathbf{N}\mathbf{b}_j)/\text{var}(\mathbf{Y}\mathbf{N}\mathbf{b}_j) \\ &= \rho_{i,j}^2 \text{var}(\mathbf{X}_i\mathbf{M}_i\mathbf{a}_{ij}) \\ &= \text{var}(\mathbf{X}_{Ai}\mathbf{M}_i\mathbf{a}_{ij}) \end{aligned}$$

Soit finalement :

$$\mathbf{a}'_{ij}\mathbf{M}_i\mathbf{X}'_i\mathbf{D}\mathbf{X}_{Ai}\mathbf{M}_i\mathbf{a}_{ij} = \mathbf{a}'_{ij}\mathbf{M}_i\mathbf{X}'_{Ai}\mathbf{D}\mathbf{X}_{Ai}\mathbf{M}_i\mathbf{a}_{ij}. \quad (46)$$

Or d'après (34) et (41) et puisque $\mathbf{X}'_{Ni}\mathbf{D}\mathbf{X}_{Ai}$ est une matrice nulle, on a :

$$\mathbf{a}'_{ij}\mathbf{M}_i\mathbf{X}'_i\mathbf{D}\mathbf{X}_{Ai}\mathbf{M}_i\mathbf{a}_{ij} = \mathbf{a}'_{ij}\mathbf{M}_i\mathbf{X}'_{Ai}\mathbf{D}\mathbf{X}_{Ai}\mathbf{M}_i\mathbf{a}_{ij} + \mathbf{a}'_{ij}\mathbf{M}_i\mathbf{X}'_{Di}\mathbf{D}\mathbf{X}_{Ai}\mathbf{M}_i\mathbf{a}_{ij}. \quad (47)$$

On en déduit avec (46) que :

$$\mathbf{a}'_{ij}\mathbf{M}_i\mathbf{X}'_{Di}\mathbf{D}\mathbf{X}_{Ai}\mathbf{M}_i\mathbf{a}_{ij} = \mathbf{0}. \quad (48)$$

Par suite, puisque $\mathbf{X}'_{Ni}\mathbf{D}\mathbf{X}_{Di} = \mathbf{0}$, on a aussi :

$$\mathbf{a}'_{ij}\mathbf{M}_i\mathbf{X}'_i\mathbf{D}\mathbf{X}_{Di}\mathbf{M}_i\mathbf{a}_{ij} = \mathbf{a}'_{ij}\mathbf{M}_i\mathbf{X}'_{Di}\mathbf{D}\mathbf{X}_{Di}\mathbf{M}_i\mathbf{a}_{ij}. \quad (49)$$

Finalement, comme :

$$a'_{ij}M_iX'_iDX_iM_ia_{ij} = a'_{ij}M_iX'_iDX_{A_i}M_ia_{ij} \\ + a'_{ij}M_iX'_iDX_{D_i}M_ia_{ij} + a'_{ij}M_iX'_iDX_{N_i}M_ia_{ij},$$

on a bien :

$$a'_{ij}M_iX'_iDX_iM_ia_{ij} = a'_{ij}M_iX'_{A_i}DX_{A_i}M_ia_{ij} \\ + a'_{ij}M_iX'_{D_i}DX_{D_i}M_ia_{ij} + a'_{ij}M_iX'_{N_i}DX_{N_i}M_ia_{ij},$$

Chacun des trois graphiques à juxtaposer est donc obtenu par projection des nuages respectifs (X_{A_i}, M_i) , (X_{D_i}, M_i) et (X_{N_i}, M_i) dans un sous-espace commun, engendré par un sous-ensemble de deux ou trois vecteurs du système d'axes $\{a_{ij}\}_{j=1, \dots, S}$.

Cette complémentarité entre inerties pour les graphiques juxtaposés peut être complétée, comme en analyse de concordance d'un tableau avec un autre, par une complémentarité de mesures partielles de l'intensité par axe de chacune des trois notions de concordance, discordance et anti- image.

Deux types de graphiques peuvent donc être préconisés : un premier type juxtapose les $K + 1$ biplots de concordance dans les repères qui se correspondent et qui permettent d'interpréter les axes (Lafosse, 1999a), ou juxtapose dans ces repères les individus des $K + 1$ images discordantes, ou encore des $K + 1$ anti-images; avec le deuxième type, images concordantes, discordance et anti-image sont juxtaposés en se servant du même repère.

4. Application

4.1. Présentation des données

On reprend ici des données d'écologie marine largement décrites par diverses approches multivariées dans Fromentin *et al.* (1993). On pourra donc s'y référer pour une connaissance plus complète. Les 68 relevés successifs effectués en s'éloignant de la côte, sont ici répartis en 5 zones, reprenant en partie une subdivision adoptée par Fromentin. Ces relevés constituent les 68 lignes des deux tableaux de départ.

Zone 0 : relevés de 1 à 16, zone côtière.

zone 1 : relevés de 17 à 25, zone de divergence D1.

Zone 2 : relevés de 26 à 30, zone de convergence.

Zone 3 : relevés de 31 à 41, zone frontale.

Zone 4 : relevés de 42 à 46, zone de divergence D2.

Zone 5 : relevés de 47 à 68, zone centrale.

En fait, la zone côtière 0 a été supprimée, suite à une étude préliminaire. En effet, cette zone s'est révélée très prépondérante par rapport à toutes les autres quant à la dépendance créée entre les deux tableaux. Cette première étude pouvant donc préciser surtout le rôle spécifique de la zone 0, nous passons ici directement à une deuxième étude relative aux autres zones. On a cependant gardé les numéros initiaux des relevés.

Quatre paramètres hydrologiques ont été mesurés pour chaque relevé : température, salinité, fluorescence et densité. Ils définissent un tableau de dimension 52×4 qui a été centré et réduit pour donner un tableau Y .

Dans le même temps, le dénombrement de 24 espèces est réalisé à partir d'une même masse d'eau et du plancton chaque fois recueilli, constituant les données biologiques. Le logarithme (+1) de ces effectifs conduit à un tableau X centré et réduit de dimension 52×24 . Les espèces sont :

Acartia (Acar)	Copepodits 5 (Cop5)	Nauplii (Naup)	Gasteropods larvae (Gast)
Adults of Calanus (Cala)	Clausocalanus A (ClauA)	Oithona (Oith)	Eggs of crustaceans (Eggs)
Copepodits 1 (Cop1)	Clausocalanus B (ClauB)	Acanthaires (Acan)	Ostracods (Ostr)
Copepodits 2 (Cop2)	Clausocalanus C (ClauC)	Cladocerans (Clad)	Pteropods (Pter)
Copepodits 3 (Cop3)	Adults of Centropages (CenA)	Echinoderms larvae (Echi)	Siphonophores (Siph)
Copepodits 4 (Cop4)	Juvenils of Centropages (CenJ)	Decapods larvae (Deca)	Bells of calyophore(Bell)

Ici, les 24 espèces biologiques, qui ont toutes le même poids sachant la réduction des variables, ont été partitionnées en 3 groupes définis d'après l'importance de densité maximale observée sur les 68 relevés. Le groupe I est ainsi constitué de 8 espèces qui en nombre ne dépasse jamais l'effectif 200 par relevé (par exemple sans jamais dépasser 8 pour Pteropods). Le groupe III est constitué de 8 espèces pouvant aller jusqu'à l'effectif 10 000 sur un seul relevé et le groupe II est constitué de 8 espèces intermédiaires de ce point de vue. Les 3 groupes définissant les 3 tableaux X_j sont alors :

G I	Ostr	Acan	Siph	Clad	Echi	Pter	Gast	Bell
G II	Cop1	Cop2	Cop3	Cop4	Cop5	Naup	Eggs	Deca
G III	Acar	Cala	ClauA	ClauB	ClauC	CenA	CenJ	Oith

4.2. Une analyse CONCOR

Une des analyses menées par Fromentin est l'analyse canonique des 2 tableaux initiaux formés des 68 relevés, permettant l'étude de dépendance sans prise en compte des interactions entre espèces ni des corrélations entre paramètres hydrologiques.

La dépendance simultanée des quatre paramètres hydrologiques avec ces 3 groupes biologiques est ici étudiée en prenant en compte les interactions possibles que les espèces peuvent avoir quant à leurs présence-absences à l'intérieur de chaque groupe $X_j, j = 1, \dots, 3$. C'est dire que les trois métriques correspondantes M_j choisies sont les métriques identité.

De même, notre problème n'étant pas de définir des combinaisons linéaires des paramètres hydrologiques susceptibles de servir au mieux d'explication pour la présence des espèces, mais de détailler la dépendance entre tableaux, les raisons de ne pas tenir compte des fortes corrélations qui existent entre ces paramètres semblent alors manquer. La métrique \mathbf{N} est ainsi choisie égale à \mathbf{I} .

TABLEAU 1
Contributions par axes.

axe	hydrologiques			Biologique GI			Biologique GII			Biologique GIII		
	conc	disc	anti	conc	disc	anti	conc	disc	anti	conc	disc	anti
1	582	53	159	49	1	17	74	1	26	168	3	58
2	97	4	61	59	1	39	110	1	72	84	1	55
3	12	1	31	3	1	7	11	3	27	15	5	37
4	0	0	0	5	2	35	3	1	20	1	0	3
	691	58	252	117	4	98	198	7	145	268	9	154
	1000			219			350			431		

Les contributions sont exprimées en millièmes. La somme des contributions des 3 groupes vaut 1000. Ces mesures sont explicitées en section 3.2 pour les mesures hydrologiques (les 3 colonnes correspondant respectivement aux 3 tableaux de coordonnées (14), (18) et (16)) et en section 3.3 pour les mesures biologiques (tableaux de coordonnées (20), (40) et (39)).

On constate d'après la table 1 des contributions que les axes 3 et 4 apportent une information relative négligeable en ce qui concerne la dépendance entre tableaux (colonnes **conc** comme concordance). En fait il ne peut en être tellement autrement, sachant que les 4 dimensions de départ du tableau \mathbf{Y} pourraient être ramenées à 2, avec une ACP par exemple, sans perte sensible d'information.

Les espèces les plus dépendantes des paramètres hydrologiques sont les espèces les plus denses si on se réfère aux pourcentages 43.1 %, 35 % et 21.9 %.

Les faibles discordances calculées sur les 3 groupes biologiques indiquent que l'analyse serait peu changée si on supprimait l'information sur les corrélations entre paramètres hydrologiques (ce qui induirait la nullité des discordances sur les 3 groupes). Si on se réfère au graphique correspondant de la figure 1 on voit en effet que la structure des 4 variables est assez bien résumable par 2 variables réduites très peu corrélées, comme température et densité. Ainsi donc l'introduction de la métrique de Mahalanobis sur \mathbf{Y} ne changerait pas tellement ce résumé.

Par ailleurs, le pourcentage de 5.8 % est la mesure globale de perte de dépendance entre tableaux due à la prise en compte des interactions entre espèces à l'intérieur de chacun des 3 groupes. Il apparaît plutôt faible aussi, et on pourrait penser que cela est partiellement dû à notre façon de constituer les trois groupes, dont le but n'était pas de rassembler entre elles les espèces pouvant interagir. Or l'analyse faite en rassemblant

en un seul groupe toutes les espèces ensemble donne un pourcentage de 5.9%. Tout semble indiquer alors que notre répartition des espèces en 3 groupes n'a pas détruit l'essentiel des interactions entre espèces.

C'est pour les relevés 20, 32 et 49 que l'interaction entre espèces vient le plus contrarier la dépendance espèces - paramètres hydrologiques (contributions à la discordance observée plus élevées pour les 3 lignes correspondantes de Y_D). L'analyse faite avec un seul groupe réunissant les 3 groupes désigne aussi d'abord ces 3 relevés.

Les cumuls de contributions à la discordance Y_D (table 2) précisent que la zone 2 est celle pour laquelle l'interaction entre les espèces n'a pas d'influence sur la dépendance. Les mesures ci-après sont relativisées par le nombre de relevés dans chaque zone.

La figure 1 détaille la dépendance dans les dimensions 1 et 2. Les 4 facteurs hydrologiques ont un rôle également important quant à leur influence sur les espèces puisque les longueurs de flèches sont sensiblement égales (biplot Y en haut à gauche). Sur les 3 autres biplots les plus grandes longueurs de flèches indiquent les espèces les plus dépendantes des facteurs hydrologiques, la même échelle étant utilisée pour les trois graphiques. Par construction, les directions et sens des flèches sur ces trois biplots sont associables aux mêmes directions considérées sur le biplot Y (et par là associables entre elles). Ainsi la présence des espèces *Clausocalanus C* et *Decapods* est fortement liée à une densité de l'eau élevée. Ainsi les relevés de zone 5 associés à une température et une salinité élevées se caractérisent par une absence relative de la plupart des espèces. L'interprétation des axes est donnée par les flèches plutôt que par les points. Les 4 nuages de points, codés par l'appartenance des relevés à leur zone, se ressemblent. Ils seraient identiques si la dépendance était totale (avec des corrélations calculées égales à 1). Les relevés les plus éloignés de l'origine sont ceux qui contribuent le plus à la dépendance. Les cercles tracés servent à faciliter la comparaison de tels éloignements. L'échelle, seulement indiquée pour le premier groupe (subset1), est celle aussi adoptée pour les 2 autres biplots des 2 autres groupes (subset 2 et 3). Le nuage le plus ramassé autour de l'origine désigne ainsi le Groupe I comme celui participant le moins à la dépendance.

L'analyse de dépendance avec Y faite en réunissant les trois groupes en un même groupe présente une forte analogie : on a vérifié que le graphe unique où toutes les espèces sont alors représentées correspond assez à la superposition des trois graphes associés aux trois groupes ici obtenus. En fait cela se produit parce que nos données sont bien particulières : on a déjà remarqué que la dépendance entre espèces et les paramètres hydrologiques était captée pour l'essentiel par des espaces de dimension 2, sachant la structure des paramètres ; donc rien n'est changé à ce sujet quand on divise les espèces en trois groupes, la dépendance des espèces s'établissant sensiblement avec le même sous-espace de dimension 2 dans les deux approches. L'intérêt de la partition est ici moins exploitable à cause de cette pauvreté en dimensions.

La figure 2 décrit les contributions des relevés à la discordance, en usant des mêmes repères respectifs utilisés en concordance. L'interprétation des axes est donc inchangée : c'est celle donnée par les flèches sur les biplots de concordance. Les quatre graphiques peuvent être respectivement associés aux 4 biplots précédents, explicitant des contributions complémentaires.

TABLEAU 2

Sensibilité relative de la dépendance étudiée aux interactions entre espèces.

zone 1	zone 2	zone 3	zone 4	zone 5
.24	.07	.25	.18	.25

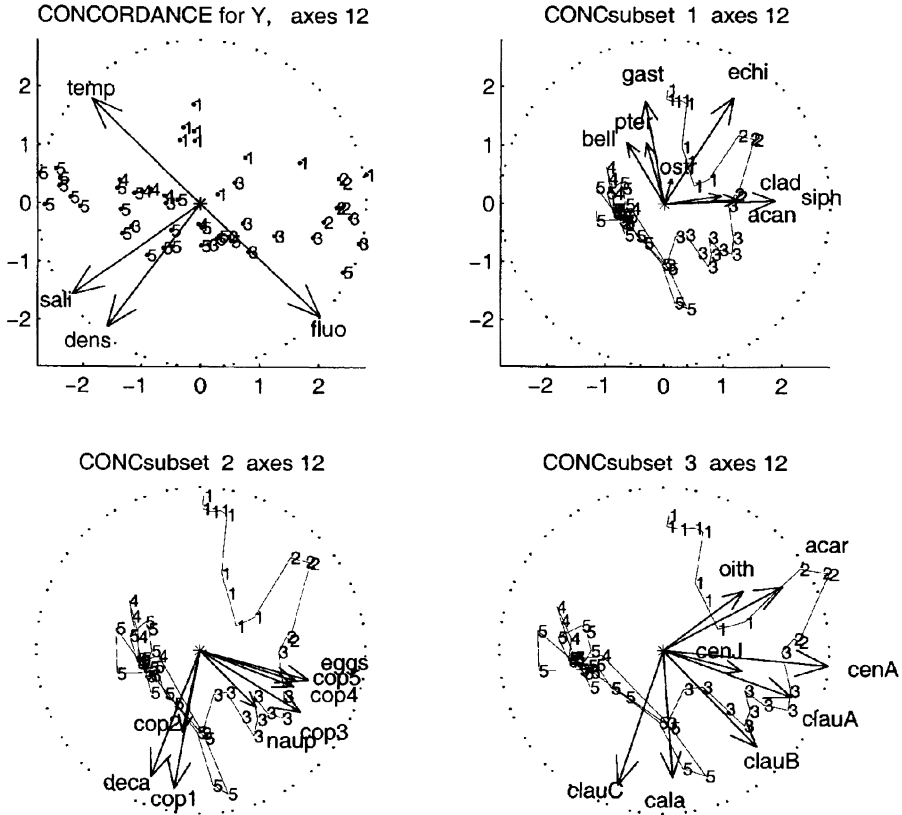


FIGURE 1

Les longueurs de flèches sont indiquées en section 3.2.1 pour le biplot Y et formule (37) pour les 3 autres biplots. Les vecteurs b_1 et b_2 pour Y, a_{i1} et a_{i2} , avec $i = 1, 2$ et 3 respectivement pour les 3 autres biplots, contiennent les cosinus directeurs des flèches. Les points correspondent aux coordonnées des lignes de Y_A projetées dans le repère (b_1, b_2) pour le biplot Y et aux coordonnées des lignes de X_{A_i} dans les repères respectifs (a_{i1}, a_{i2}) , $i = 1, 2$ et 3, pour les trois autres biplots.

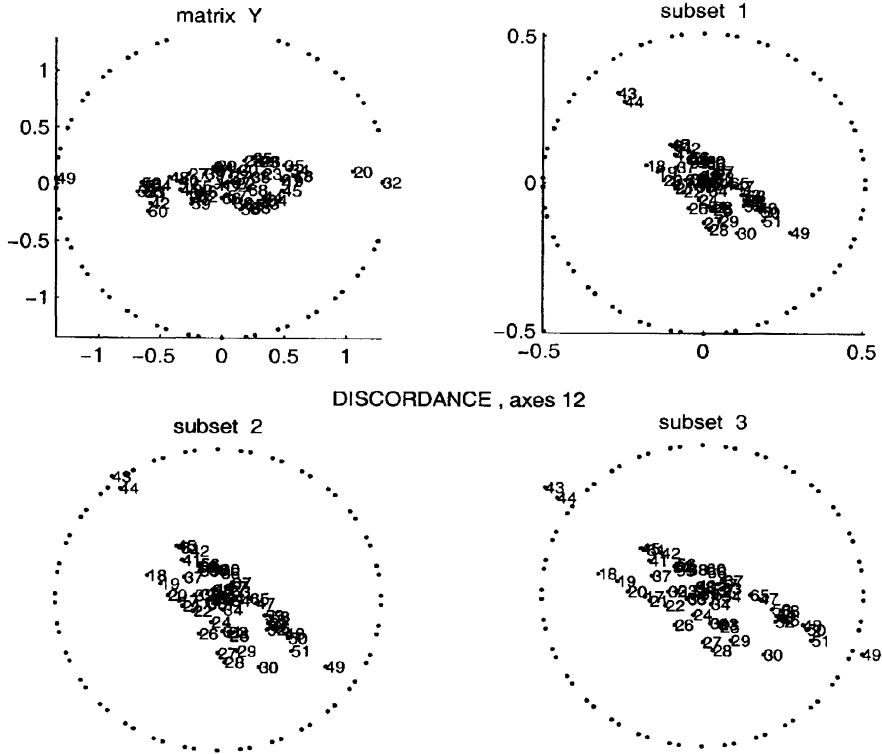


FIGURE 2

Coordonnées des lignes de Y_D dans le repère (b_1, b_2) pour le plot Y et coordonnées des lignes de X_{D_i} dans les repères respectifs (a_{i1}, a_{i2}) , $i = 1, 2$ et 3 , pour les trois autres plots.

La discordance peut paraître assez négligeable. Cependant l'allongement des nuages traduit une cohérence trouvée sur l'ensemble des relevés, qui peut être prise en compte. L'allongement horizontal observé sur le graphique Y indique que les espèces moins dépendantes des paramètres hydrologiques pour la raison qu'elles tiendraient compte de la présence ou de l'absence d'autres espèces seraient plutôt celles constitutives de l'axe 1, comme centropages adultes. Les autres trois nuages sont tous obliquement allongés dans la même direction. Cette répétition est naturelle, la discordance caractérisant quand $K > 1$ et pour des axes relatifs à des corrélations élevées une façon commune pour les trois nuages de se différencier de Y . C'est l'interaction entre température et fluorescence qui est ici désignée comme un facteur de réduction de la dépendance entre paramètres hydrologiques et les trois groupes.

La figure 3 révèle les contributions à l'indépendance entre espèces et paramètres hydrologiques pour les axes 1 et 2. Ces axes sont encore les mêmes axes respectifs précédents et les contributions respectives sont à nouveau complémentaires. Le fait que de mêmes relevés (comme 17 ou 38) apparaissent loin de l'origine sur les 4

graphiques signifient que la contribution à l'indépendance de ces relevés n'est pas liée à l'appartenance à un groupe d'espèces particulier. Plus globalement, l'ensemble des mesures caractérisent une façon originale pour chaque groupe d'être assez indépendant des paramètres biologiques, et cela pourrait déjà suffire à justifier la partition. Celle du groupe 1 se fait plutôt selon l'axe 2, grâce à la présence de Gastéropodes, Echinoderms et Bells dans les zones 1-2 et de leur absence dans la zone 3. Celle du groupe 2 se fait aussi plus selon l'axe 2, grâce notamment la présence de Decapods et Copepodits 1 dans les zones 3 et 5 et leurs absences des zones 1 et 2. Celle du groupe 3 est plus diffuse, associée aux deux axes simultanément, des espèces étant présentes presque dans toutes les zones.

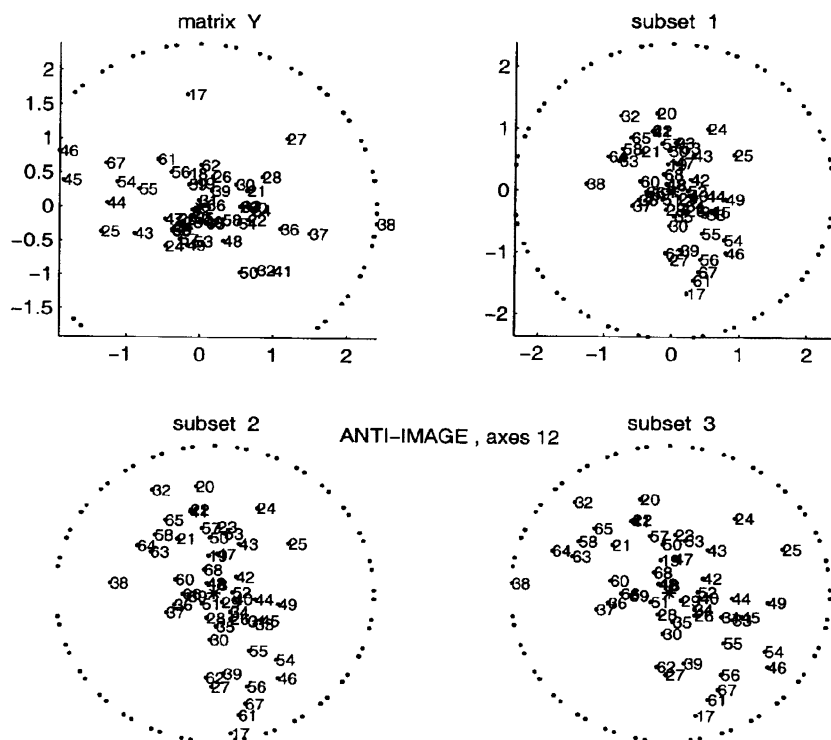


FIGURE 3

Coordonnées des lignes de Y_N dans le repère (b_1, b_2) pour le plot Y et coordonnées des lignes de X_{N_i} dans les repères respectifs (a_{i1}, a_{i2}) , $i = 1, 2$ et 3 , pour les trois autres plots.

Références

CAZES P., BAUMERDER A., BONNEFOUS S. & PAGES J.P. (1977), Codage et analyse des tableaux logiques. Introduction à la pratique des variables qualitatives. *Cahiers du BUR0*, 27, Paris VI.
 CLIFF N. (1966), Orthogonal rotation to congruence. *Psychometrika*, 31, 33-42.

- FROMENTIN J.M., IBANEZ F. & LEGENDRE P. (1993), A Phytosociological Method for interpreting Plankton data. *Mar.Ecol.Progr.Series* 93 : 285-306.
- GREEN B.F.(1952), The Orthogonal Approximation of an Oblique Structure in factor analysis. *Psychometrika*, **17**, 429-440.
- GUTTMAN L. (1953), Image theory for the structure of quantitative variates. *Psychometrika*, **18**, 277-296.
- HANAFI M. (1997), Structure de l'ensemble des analyses multivariées des tableaux de données à trois entrées : Eléments théoriques et appliqués. *Thèse*, Univ. Claude Bernard, Lyon I.
- HORST (1961), Relation among m sets of variables. *Psychometrika*, **38**, 433-451.
- LAFOSSE R. (1985), Une nouvelle analyse procrustéenne de deux tableaux. *Data Analysis and Informatics, IV (Diday & coll. eds)*. Elsevier Science Pub. North Holland, 407-414.
- LAFOSSE R. (1997), Analyse de concordance de deux tableaux : monogamies, simultanités et découpages. *Rev. Stat. Appliquée*. vol. 45. n°3. 45-72.
- LAFOSSE R. & HANAFI M. (1997), Concordance d'un tableau avec K tableaux : définition de $K + 1$ uples synthétiques. *Rev. Stat. Appliquée*. vol 45 n°4. 111-126.
- LAFOSSE (1998), Programme de l'analyse CONCOR en langage MATLAB. version 1.3. <http://www.mathworks.com/support/ftp/statv5.shtml> ou mel à : lafosse@cict.fr.
- LAFOSSE (1999a), Analyses of some relation between arrays and graphics. *Proc. Int. Conf. on Probability and Statistics and their Applications (Hanoi)*.
- LAFOSSE (1999b), Analysis of concordance between matrices and proposals for selecting variables. *IX Int. Symp. on Applied Stochastic Models and Data Analysis (Lisbonne)*.
- TEN BERGE J.M.F. (1977), Optimizing Factorial Invariance. *Thesis Groningen*.
- TEN BERGE J.M.F. (1977), Orthogonal procrustes rotation for two or more matrices. *Psychometrika*, **42**, **2**, 267-273.
- TEN BERGE J.M.F. (1988), Generalized approaches to the maxbet problem and the maxdiff problem, with applications to canonical correlations. *Psychometrika*, **52**, **4**, 487-494.
- TUCKER L.R. (1958), An interbattery method of factor analysis. *Psychometrika*, **23**, 111-136.
- VAN DE GEER P. (1984), Linear relations among K sets of variables. *Psychometrika*, **49**,**1**, 79-94.