

# NOTES DE STATISTIQUE ET D'INFORMATIQUE

2012/2

## LA REPRÉSENTATION D'UNE MATRICE PAR BIPLLOT

R. PALM, C. CHARLES et J. J. CLAUSTRIAUX

Université de Liège – Gembloux Agro-Bio Tech  
*Unité de Statistique, Informatique et Mathématique  
appliquées à la bioingénierie*  
**GEMBLOUX**  
(Belgique)

# LA REPRÉSENTATION D'UNE MATRICE PAR BILOT

R. PALM<sup>\*</sup>, C. CHARLES<sup>†</sup> et J. J. CLAUSTRIAUX<sup>‡</sup>

## RÉSUMÉ

La factorisation d'une matrice par décomposition par valeurs singulières est examinée et utilisée pour la représentation d'une matrice par biplot. Le lien entre cette représentation et les représentations habituelles des variables et des individus réalisées lors des analyses en composantes principales est ensuite examiné et illustré par un exemple numérique.

## SUMMARY

Matrix factorization by means of singular value decomposition is examined and used to produce a graphical representation of a data matrix called biplot. The link between this biplot and the plots of the variables and of the individuals usually given in principal component analysis is discussed and applied to an example.

## 1. INTRODUCTION

Un *biplot* est une représentation graphique d'une matrice, généralement dans un espace à deux dimensions, chaque ligne et chaque colonne étant représentée par un point.

Le préfixe *bi* fait allusion aux deux types de points et non à la dimension de l'espace servant à la représentation. D'un point de vue théorique, la dimension de l'espace peut être augmentée sans difficulté, mais au-delà de la dimension 3, la visualisation ne peut se faire que par des projections des points dans des sous-espaces.

La matrice qui fait l'objet de la représentation n'est généralement pas la matrice des observations brutes. Les variables sont en effet le plus souvent centrées, par soustraction de la moyenne, et standardisées d'une manière ou d'une autre.

---

<sup>\*</sup>Professeur à l'Université de Liège, Gembloux Agro-Bio Tech.

<sup>†</sup>Chargée de cours à l'Université de Liège, Gembloux Agro-Bio Tech.

<sup>‡</sup>Professeur ordinaire à l'Université de Liège, Gembloux Agro-Bio Tech.

Le biplot repose sur la factorisation des matrices. La technique générale est présentée au paragraphe 2 et une factorisation particulière, la décomposition par valeurs propres, est décrite au paragraphe 3. Les propriétés de trois formes de biplot, basées sur cette décomposition par valeurs singulières, sont examinées au paragraphe 4. La relation entre les biplots et l'analyse en composantes principales fait l'objet du paragraphe 5. Enfin, quelques informations complémentaires sont données au paragraphe 6.

La représentation par biplot a été proposée par GABRIEL [1971]. Le lecteur y trouvera une description approfondie de la méthode. Des informations peuvent également être trouvées dans plusieurs livres consacrés à l'analyse statistique multivariée, en particulier dans les ouvrages de JACKSON [1991], KHATTREE et NAIK [1995], RENCHER [2002] et SEBER [1984], ainsi que dans les monographies consacrées entièrement aux biplots de GOWER et HAND [1996] et de GOWER *et al.* [2011].

## 2. FACTORISATION, POINTS-LIGNES ET POINTS-COLONNES

### 2.1. Factorisation d'une matrice

Soit  $\mathbf{Y}$  une matrice de dimensions  $n \times p$  et de rang  $r$ . Il pourrait s'agir, par exemple, d'une matrice d'observations réalisées sur  $n$  individus pour  $p$  variables. Cette matrice  $\mathbf{Y}$  peut être exprimée sous la forme d'un produit de deux matrices :

$$\mathbf{Y} = \mathbf{A} \mathbf{B}',$$

$\mathbf{A}$  étant une matrice de dimensions  $n \times r$  et de rang  $r$  et  $\mathbf{B}$  une matrice de dimensions  $p \times r$  et de rang  $r$  également. Ces matrices  $\mathbf{A}$  et  $\mathbf{B}$  qui interviennent dans la factorisation de  $\mathbf{Y}$  ne sont pas uniques et nous examinerons, au paragraphe 3, comment elles peuvent être déterminées.

Chaque élément  $y_{ij}$  de  $\mathbf{Y}$  est égal au produit de deux vecteurs :

$$y_{ij} = \mathbf{a}_i \mathbf{b}_j' \quad (i = 1, \dots, n; j = 1, \dots, p),$$

$\mathbf{a}_i$  étant la  $i^{\text{ème}}$  ligne de  $\mathbf{A}$  et  $\mathbf{b}_j$  la  $j^{\text{ème}}$  ligne de  $\mathbf{B}$ .

Cette relation peut encore s'écrire :

$$y_{ij} = a_{i1} b_{j1} + a_{i2} b_{j2} + \dots + a_{ir} b_{jr}.$$

A titre d'illustration, considérons la matrice quelconque suivante :

$$\mathbf{Y} = \begin{bmatrix} 15 & 5 & 11 \\ 8 & 1 & 6 \\ 10 & 6 & 8 \\ 19 & 14 & 14 \end{bmatrix},$$

qui peut être factorisée par les matrices suivantes :

$$\mathbf{A} = \begin{bmatrix} 3,11 & -1,31 & 0,21 \\ 1,56 & -1,37 & -0,05 \\ 2,31 & 0,10 & -0,63 \\ 4,46 & 1,34 & 0,20 \end{bmatrix} \text{ et } \mathbf{B} = \begin{bmatrix} 4,47 & -0,76 & 0,41 \\ 2,50 & 2,12 & 0,00 \\ 3,33 & -0,57 & -0,56 \end{bmatrix}.$$

On notera que les éléments des matrices  $\mathbf{A}$  et  $\mathbf{B}$  ont été arrondis à deux décimales. Les calculs ultérieurs ont toutefois été réalisés avec un plus grand nombre de décimales. Le lecteur qui referait les calculs de façon manuelle pourrait donc obtenir des résultats légèrement différents de ceux reproduits ici.

On peut vérifier que, par exemple, l'élément  $y_{32}$  est bien égal aux produits des vecteurs  $\mathbf{a}_3$  et  $\mathbf{b}'_2$ , avec :

$$\mathbf{a}_3 = (2,31 \quad 0,10 \quad -0,63) \quad \text{et} \quad \mathbf{b}_2 = (2,50 \quad 2,12 \quad 0,00).$$

En effet :

$$y_{32} = (2,31)(2,50) + (0,10)(2,12) - (0,63)(0,00) \simeq 6.$$

Des calculs similaires pourraient être réalisés pour les douze éléments de  $\mathbf{Y}$ .

## 2.2. Points-lignes et points-colonnes pour une matrice de rang deux

Nous considérons d'abord la représentation graphique d'une matrice de rang deux. Les matrices  $\mathbf{A}$  et  $\mathbf{B}$  n'ont alors que deux colonnes et les vecteurs  $\mathbf{a}_i$  et  $\mathbf{b}_j$  n'ont que deux éléments. Les  $n$  vecteurs  $\mathbf{a}_i$  et les  $p$  vecteurs  $\mathbf{b}_j$  peuvent être représentés dans un espace de dimension 2. Et pour distinguer les vecteurs  $\mathbf{a}_i$  des vecteurs  $\mathbf{b}_j$  dans ce graphique, on représente une série de vecteurs, par exemple les  $\mathbf{a}_i$ , par des points situant l'extrémité des vecteurs et l'autre série de vecteurs, les  $\mathbf{b}_j$ , par des flèches rejoignant l'origine des axes aux extrémités des vecteurs. Ces vecteurs seront appelés par la suite, respectivement les points-lignes et les points-colonnes, car les lignes de  $\mathbf{A}$  sont associées aux lignes de  $\mathbf{Y}$  et les lignes de  $\mathbf{B}$  sont associées aux colonnes de  $\mathbf{Y}$ .

Les vecteurs  $\mathbf{a}_i$  et  $\mathbf{b}_j$  présentent trois propriétés qui seront utilisées par la suite et que nous énonçons ci-dessous.

1. La longueur d'un vecteur (longueur du segment rejoignant l'origine au point pour un vecteur  $\mathbf{a}_i$  ou longueur de la flèche pour un vecteur  $\mathbf{b}_j$ ), appelée aussi norme, est égale à la racine carrée du produit du vecteur par sa transposée :

$$\|\mathbf{a}_i\| = (\mathbf{a}_i \mathbf{a}'_i)^{1/2} = \sqrt{a_{i1}^2 + a_{i2}^2}$$

et

$$\|\mathbf{b}_j\| = (\mathbf{b}_j \mathbf{b}'_j)^{1/2} = \sqrt{b_{j1}^2 + b_{j2}^2}.$$

2. Le produit scalaire (ou produit interne) des vecteurs  $\mathbf{a}_i$  et  $\mathbf{b}_j$  est égal à  $y_{ij}$  :

$$\langle \mathbf{a}_i, \mathbf{b}_j \rangle = \|\mathbf{a}_i\| \|\mathbf{b}_j\| \cos \theta_{ij} = a_{i1} b_{j1} + a_{i2} b_{j2} = y_{ij},$$

$\theta_{ij}$  étant l'angle formé par les deux vecteurs. L'observation  $y_{ij}$  est donc égale à la longueur du vecteur  $\mathbf{b}_j$ , multipliée par la longueur de la projection du vecteur  $\mathbf{a}_i$  sur le vecteur  $\mathbf{b}_j$  ou encore à la longueur du vecteur  $\mathbf{a}_i$ , multipliée par la longueur de la projection du vecteur  $\mathbf{b}_j$  sur le vecteur  $\mathbf{a}_i$ . Il en résulte que l'angle  $\theta_{ij}$  entre les deux vecteurs  $\mathbf{a}_i$  et  $\mathbf{b}_j$  est donné par la relation :

$$\theta_{ij} = \arccos \left( \frac{\langle \mathbf{a}_i, \mathbf{b}_j \rangle}{\|\mathbf{a}_i\| \|\mathbf{b}_j\|} \right) = \arccos \left( \frac{\sum_{s=1}^2 a_{is} b_{js}}{\sqrt{\sum_{s=1}^2 a_{is}^2} \sqrt{\sum_{s=1}^2 b_{js}^2}} \right).$$

Deux vecteurs  $\mathbf{a}_i$  et  $\mathbf{b}_j$  forment par conséquent un angle droit lorsque  $y_{ij} = 0$ , puisque, dans ce cas, le produit scalaire est nul. Ils forment un angle aigu lorsque  $y_{ij} > 0$  et ils forment un angle obtus si  $y_{ij} < 0$ .

3. L'angle formé par les deux vecteurs  $\mathbf{a}_i$  et  $\mathbf{a}_{i'}$  ( $i' \neq i$ ) ou  $\mathbf{b}_j$  et  $\mathbf{b}_{j'}$  ( $j' \neq j$ ) est égal à :

$$\theta_{ii'} = \arccos \left( \frac{\sum_{s=1}^2 a_{is} a_{i's}}{\sqrt{\sum_{s=1}^2 a_{is}^2} \sqrt{\sum_{s=1}^2 a_{i's}^2}} \right) \quad \text{et} \quad \theta_{jj'} = \arccos \left( \frac{\sum_{s=1}^2 b_{js} b_{j's}}{\sqrt{\sum_{s=1}^2 b_{js}^2} \sqrt{\sum_{s=1}^2 b_{j's}^2}} \right).$$

### 2.3. Points-lignes et points-colonnes pour une matrice de rang quelconque

Si la matrice est de rang  $r$  ( $r > 2$ ), les vecteurs  $\mathbf{a}_i$  et  $\mathbf{b}_j$  doivent être représentés dans un espace de dimension  $r$  et les propriétés vues ci-dessus restent valables lorsque les différentes sommes qui interviennent dans les relations sont étendues aux  $r$  termes.

Si une telle extension ne pose pas de problèmes théoriques, elle ne permet cependant plus une représentation graphique simple. Pour se limiter, par la suite, à une représentation graphique à deux dimensions, il peut se justifier de ne retenir que deux colonnes de  $\mathbf{A}$  et de  $\mathbf{B}$  et de négliger les  $r - 2$  autres colonnes. Les propriétés relatives aux vecteurs rappelées ci-dessus restent valables, sauf en ce qui concerne le produit scalaire des vecteurs  $\mathbf{a}_i$  et  $\mathbf{b}_j$ . Celui-ci n'est, en effet, plus égal à  $y_{ij}$ , mais à une approximation de  $y_{ij}$ ,  $r - 2$  termes ayant été négligés. Comme nous le verrons ci-dessous, dans la pratique on retiendra souvent les deux premières colonnes de  $\mathbf{A}$  et de  $\mathbf{B}$  et on négligera les  $r - 2$  dernières colonnes, mais la factorisation de  $\mathbf{Y}$  sera optimisée, de manière à ce que l'approximation des  $y_{ij}$  soit aussi bonne que possible.

Pour ne pas alourdir inutilement les notations, nous ne ferons plus, par la suite la distinction entre le cas où  $\mathbf{Y}$  est de rang 2 et le cas où  $\mathbf{Y}$  est de

rang supérieur à 2. Nous donnerons les formules générales pour les matrices de rang  $r$ , étant entendu que, dans la pratique, on se limitera le plus souvent aux deux premières colonnes de  $\mathbf{A}$  et de  $\mathbf{B}$ , ce qui revient à dire qu'on n'étudie pas la matrice  $\mathbf{Y}$ , mais une approximation de celle-ci, notée  $\hat{\mathbf{Y}}_{(2)}$  ou encore, plus simplement,  $\hat{\mathbf{Y}}$ , qui est une matrice de rang 2.

Si, pour l'exemple examiné précédemment, on élimine la dernière colonne des matrices  $\mathbf{A}$  et  $\mathbf{B}$ , le produit matriciel ne redonne pas la matrice  $\mathbf{Y}$  mais une approximation de  $\mathbf{Y}$  :

$$\hat{\mathbf{Y}} = \begin{bmatrix} 3,11 & -1,31 \\ 1,56 & -1,37 \\ 2,31 & 0,10 \\ 4,46 & 1,34 \end{bmatrix} \begin{bmatrix} 4,47 & 2,50 & 3,33 \\ -0,76 & 2,12 & -0,57 \end{bmatrix} = \begin{bmatrix} 14,92 & 5,00 & 11,11 \\ 8,02 & 1,00 & 5,97 \\ 10,26 & 6,00 & 7,65 \\ 18,92 & 14,00 & 14,11 \end{bmatrix},$$

qui, compte tenu du choix de  $\mathbf{A}$  et  $\mathbf{B}$ , est excellente, la somme des carrés des écarts entre les valeurs initiales et les valeurs approchées étant égale à :

$$\sum_{i=1}^4 \sum_{j=1}^3 (y_{ij} - \hat{y}_{ij})^2 = 0,23.$$

La figure 1 donne la représentation des lignes de  $\mathbf{A}$  et de  $\mathbf{B}$  après suppression de la dernière colonne de chacune de ces matrices. Les points représentant les lignes de  $\mathbf{A}$  (points-lignes) sont désignés par les symboles  $L_1$  à  $L_4$  et les points représentant les lignes de  $\mathbf{B}$  (points-colonnes) sont représentés par les symboles  $C_1$  à  $C_3$ .

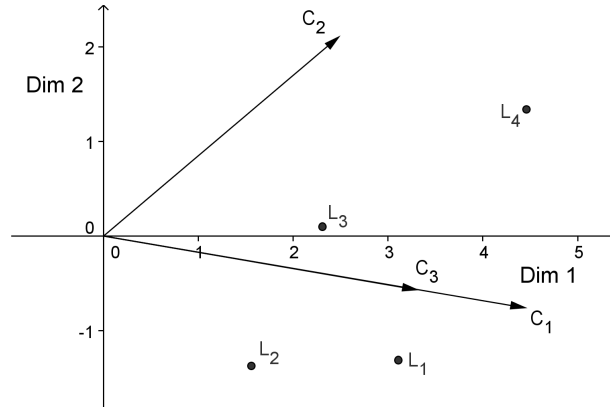


Figure 1 – Biplot de la matrice  $\hat{\mathbf{Y}}$  ( $L_1$  à  $L_4$  : points-lignes ;  $C_1$  à  $C_3$  : points-colonnes).

Ce graphique montre immédiatement que les vecteurs  $C_1$  et  $C_3$  ont pratiquement la même direction : la première et la troisième ligne de  $\mathbf{B}$  (limitées aux deux premières colonnes) sont donc à peu près proportionnelles. On constate aussi que le produit scalaire de  $L_4$  et  $C_1$  est le plus grand, ce qui signifie que

l'élément  $\hat{y}_{41}$  de  $\hat{\mathbf{Y}}$  est le plus grand. On note encore que tous les produits scalaires d'un point-ligne avec un point-colonne sont positifs, chaque angle étant inférieur à  $90^\circ$ ; la matrice  $\hat{\mathbf{Y}}$  ne contient donc que des éléments positifs. L'angle formé par les vecteurs  $\mathbf{L}_2$  et  $\mathbf{C}_2$  est proche de  $90^\circ$ ; la valeur  $\hat{y}_{22}$  est donc très faible par rapport aux autres valeurs. Dans la mesure où  $\hat{\mathbf{Y}}$  est proche de  $\mathbf{Y}$ , les commentaires faits à propos des valeurs  $\hat{y}_{ij}$  restent valables pour les  $y_{ij}$ .

### 3. DÉCOMPOSITION PAR VALEURS SINGULIÈRES

#### 3.1. Optimisation de la factorisation

Au paragraphe 2.1, nous avons vu qu'une matrice de rang 2 peut être représentée sous la forme d'un biplot, qui est une représentation simultanée d'un "effet ligne" et d'un "effet colonne". Ce biplot est réalisé à partir de la factorisation de  $\mathbf{Y}$  sous la forme d'un produit de deux matrices. Lorsque la matrice est de rang  $r$  ( $r > 2$ ), le biplot construit à partir des deux premières colonnes de  $\mathbf{A}$  et de  $\mathbf{B}$  est la représentation d'une matrice de rang 2, notée  $\hat{\mathbf{Y}}_{(2)}$  ou  $\hat{\mathbf{Y}}$ , qui est une approximation de  $\mathbf{Y}$ . Puisque la factorisation de  $\mathbf{Y}$  par le produit de deux matrices  $\mathbf{A}$  et  $\mathbf{B}'$  n'est pas unique, il y a évidemment intérêt à choisir, parmi l'ensemble des solutions possibles, celle qui, lorsqu'on ne retient que les deux premières colonnes de  $\mathbf{A}$  et de  $\mathbf{B}$ , donne la meilleure approximation de  $\mathbf{Y}$ . La décomposition par valeurs singulières de la matrice  $\mathbf{Y}$  donne la solution à ce problème.

La *décomposition par valeurs singulières*<sup>1</sup> d'une matrice  $\mathbf{Y}$ , de dimension  $n \times p$  et de rang  $r$  consiste en la factorisation suivante [ECKART et YOUNG, 1936] :

$$\mathbf{Y} = \mathbf{U}\mathbf{L}\mathbf{V}',$$

où  $\mathbf{U}$  est une matrice de dimensions  $n \times r$ ,  $\mathbf{L}$  une matrice de dimensions  $r \times r$  et  $\mathbf{V}$  une matrice de dimensions  $p \times r$ .  $\mathbf{U}$  et  $\mathbf{V}$  sont des matrices orthonormées, c'est-à-dire telles que :

$$\mathbf{U}'\mathbf{U} = \mathbf{I} \text{ et } \mathbf{V}'\mathbf{V} = \mathbf{I}.$$

De plus,  $\mathbf{U}$  est la matrice dont les colonnes sont les vecteurs propres, normés à l'unité, associés aux valeurs propres non nulles de  $\mathbf{Y}\mathbf{Y}'$ ;  $\mathbf{V}$  est la matrice dont les colonnes sont les vecteurs propres normés à l'unité, associés aux valeurs propres non nulles de  $\mathbf{Y}'\mathbf{Y}$ ;  $\mathbf{L}$  est la matrice diagonale dont les éléments diagonaux sont les valeurs singulières de  $\mathbf{Y}$ , c'est-à-dire les racines carrées des valeurs propres non nulles de  $\mathbf{Y}'\mathbf{Y}$  ou de  $\mathbf{Y}\mathbf{Y}'$ , ces deux matrices ayant des valeurs propres identiques.

La décomposition par valeurs singulières factorise  $\mathbf{Y}$  sous la forme d'un produit de trois matrices et nous verrons, au paragraphe suivant, comment celles-ci peuvent être utilisées pour une représentation par biplot.

---

1. En anglais : singular-value decomposition.

Pour la matrice  $\mathbf{Y}$  présentée au paragraphe précédent, la décomposition conduit aux résultats suivants :

$$\mathbf{U} = \begin{bmatrix} 0,51 & -0,56 & 0,30 \\ 0,26 & -0,59 & -0,07 \\ 0,38 & 0,04 & -0,91 \\ 0,73 & 0,58 & 0,29 \end{bmatrix}, \quad \mathbf{V} = \begin{bmatrix} 0,73 & -0,33 & 0,60 \\ 0,41 & 0,91 & 0,00 \\ 0,55 & -0,24 & -0,80 \end{bmatrix}$$

$$\mathbf{L} = \begin{bmatrix} 37,36 & 0 & 0 \\ 0 & 5,40 & 0 \\ 0 & 0 & 0,48 \end{bmatrix}.$$

Les matrices  $\mathbf{U}$  et  $\mathbf{V}$  sont liées par les relations suivantes :

$$\mathbf{U} = \mathbf{Y}\mathbf{V}\mathbf{L}^{-1} \quad \text{et} \quad \mathbf{V} = \mathbf{Y}'\mathbf{U}\mathbf{L}^{-1},$$

ou encore, pour une colonne  $k$  particulière :

$$\mathbf{u}_k = \frac{1}{l_k} \mathbf{Y} \mathbf{v}_k \quad \text{et} \quad \mathbf{v}_k = \frac{1}{l_k} \mathbf{Y}' \mathbf{u}_k \quad (k = 1, \dots, r).$$

Les matrices  $\mathbf{U}$ ,  $\mathbf{L}$  et  $\mathbf{V}$  ont encore les propriétés suivantes :

$$(\mathbf{U}\mathbf{L})(\mathbf{U}\mathbf{L})' = \mathbf{U}\mathbf{L}\mathbf{L}'\mathbf{U}' = \mathbf{U}\mathbf{L}^2\mathbf{U}' = \mathbf{Y}\mathbf{Y}'$$

$$(\mathbf{V}\mathbf{L})(\mathbf{V}\mathbf{L})' = \mathbf{V}\mathbf{L}\mathbf{L}'\mathbf{V}' = \mathbf{V}\mathbf{L}^2\mathbf{V}' = \mathbf{Y}'\mathbf{Y},$$

$\mathbf{L}^2$  étant la matrice diagonale des valeurs propres non nulles de  $\mathbf{Y}'\mathbf{Y}$  et de  $\mathbf{Y}\mathbf{Y}'$ .

### 3.2. Reconstitution de $\mathbf{Y}$

La factorisation de  $\mathbf{Y}$  peut s'écrire sous la forme d'une somme de produits de vecteurs :

$$\mathbf{Y} = \mathbf{U}\mathbf{L}\mathbf{V}' = \sum_{k=1}^r l_k \mathbf{u}_k \mathbf{v}_k' = \mathbf{Y}_1 + \mathbf{Y}_2 + \dots + \mathbf{Y}_r.$$

Ainsi, pour l'exemple, on a :

$$\mathbf{Y} = \begin{bmatrix} 13,91 & 7,78 & 10,37 \\ 6,97 & 3,90 & 5,20 \\ 10,34 & 5,78 & 7,71 \\ 19,94 & 11,16 & 14,87 \end{bmatrix} + \begin{bmatrix} 1,00 & -2,78 & 0,74 \\ 1,05 & -2,90 & 0,77 \\ -0,08 & 0,22 & -0,06 \\ -1,03 & 2,84 & -0,76 \end{bmatrix}$$

$$+ \begin{bmatrix} 0,09 & 0,00 & -0,11 \\ -0,02 & 0,00 & 0,03 \\ -0,26 & 0,00 & 0,35 \\ 0,08 & 0,00 & -0,11 \end{bmatrix}.$$



Chaque matrice  $\mathbf{Y}_k$  ( $k = 1, \dots, r$ ) est une matrice de rang 1 et la somme de  $q$  ( $q \leq r$ ) de ces matrices donne une matrice de rang  $q$ . L'importance de chacune de ces matrices dans la reconstitution de  $\mathbf{Y}$  est décroissante :  $\mathbf{Y}_1$  est la matrice de rang 1 qui donne la meilleure approximation de  $\mathbf{Y}$ ;  $\mathbf{Y}_1 + \mathbf{Y}_2$  est la matrice de rang 2 qui donne la meilleure approximation de  $\mathbf{Y}$  et ainsi de suite, la qualité de l'approximation se mesurant par la somme des carrés des écarts entre les  $y_{ij}$  et les  $\hat{y}_{ij}$  approchés, notés  $\hat{y}_{ij}$ , pour l'ensemble des couples  $ij$ .

La qualité relative de la reconstitution de  $\mathbf{Y}$  à partir de  $\mathbf{Y}_1 + \mathbf{Y}_2 + \dots + \mathbf{Y}_q$  est mesurée par le rapport :

$$(l_1^2 + l_2^2 + \dots + l_q^2) / (l_1^2 + l_2^2 + \dots + l_r^2),$$

qui est égal à :

$$\sum_{i=1}^n \sum_{j=1}^p \hat{y}_{ij}^2 / \sum_{i=1}^n \sum_{j=1}^p y_{ij}^2.$$

Plus ce rapport est proche de l'unité, plus l'approximation de  $\mathbf{Y}$  est bonne.

Ainsi, pour l'exemple ci-dessus, si on néglige la dernière colonne des matrices  $\mathbf{U}$  et  $\mathbf{V}$ , on a :

$$\hat{\mathbf{Y}}_{(2)} = \mathbf{Y}_1 + \mathbf{Y}_2,$$

et l'approximation de  $\mathbf{Y}$  par  $\hat{\mathbf{Y}}_{(2)}$  est excellente car le rapport :

$$(37, 36^2 + 5, 40^2) / (37, 36^2 + 5, 40^2 + 0, 48^2) = 0,9998$$

est pratiquement égal à l'unité.

## 4. COORDONNÉES DES POINTS DANS LE BIPLLOT

### 4.1. Différentes solutions

La décomposition par valeurs singulières factorise la matrice  $\mathbf{Y}$  en un produit de trois matrices.

D'autre part, nous avons vu, au paragraphe 2, que le biplot repose sur la factorisation de  $\mathbf{Y}$  en deux matrices  $\mathbf{A}$  et  $\mathbf{B}$ . Ces deux matrices vont être définies à partir de la décomposition par valeurs singulières.

On a :

$$\mathbf{Y} = \mathbf{U}\mathbf{L}\mathbf{V}' = \mathbf{A}\mathbf{B}',$$

et la détermination de  $\mathbf{A}$  et  $\mathbf{B}$  à partir de  $\mathbf{U}$ ,  $\mathbf{L}$  et  $\mathbf{V}$  peut se faire de plusieurs manières. En fonction d'un paramètre  $\alpha$ , on peut écrire la décomposition sous la forme suivante :

$$\mathbf{Y} = (\mathbf{U}\mathbf{L}^\alpha)(\mathbf{L}^{1-\alpha}\mathbf{V})' = \mathbf{A}\mathbf{B}' \quad (0 < \alpha < 1),$$

$\mathbf{L}^\alpha$  étant la matrice diagonale dont les éléments sont les valeurs  $\lambda_k^\alpha$ .

Dans la pratique, on considère le plus souvent l'une des trois solutions suivantes, qui correspondent respectivement à  $\alpha = 1$ ,  $\alpha = 0$  et  $\alpha = 1/2$  :

$$(UL)(V)' = A_1 B_1' \quad \text{avec} \quad A_1 = UL \text{ et } B_1 = V$$

$$(U)(LV)' = A_2 B_2' \quad \text{avec} \quad A_2 = U \text{ et } B_2 = VL$$

et

$$(UL^{1/2})(L^{1/2}V)' = A_3 B_3' \quad \text{avec} \quad A_3 = UL^{1/2} \text{ et } B_3 = VL^{1/2}.$$

#### 4.2. Respect des distances et des angles

Considérons d'abord le cas d'une matrice  $\mathbf{Y}$  de rang 2 et examinons la factorisation pour  $\alpha = 1$ . D'après les propriétés données au paragraphe 3.1, on a :

$$A_1 A_1' = UL(UL)' = YY'.$$

Les éléments diagonaux des matrices  $A_1 A_1'$  et  $YY'$  correspondent aux sommes des carrés des éléments des lignes des matrices  $A_1$  et  $Y$ . On a donc :

$$\mathbf{y}_i \mathbf{y}_i' = \mathbf{a}_i \mathbf{a}_i' = \sum_{j=1}^p y_{ij}^2.$$

La longueur du vecteur correspondant à un point-ligne du biplot est donc égale à la longueur du vecteur correspondant à la même ligne de la matrice  $Y$ , si on représente les lignes de  $Y$  dans l'espace des colonnes de  $Y$ .

D'autre part, pour un élément hors de la diagonale de  $A_1 A_1'$  et de  $YY'$ , on a aussi :

$$\mathbf{y}_i \mathbf{y}_{i'}' = \mathbf{a}_i \mathbf{a}_{i'}' \quad (i \neq i').$$

Il en résulte que l'angle formé par deux points-lignes du biplot est égal à l'angle formé par les deux vecteurs  $\mathbf{y}_i$  et  $\mathbf{y}_{i'}$  si on représente ces deux vecteurs dans l'espace des  $p$  colonnes :

$$\theta_{ii'} = \arccos \left( \frac{\langle \mathbf{y}_i, \mathbf{y}_{i'} \rangle}{\|\mathbf{y}_i\| \|\mathbf{y}_{i'}\|} \right) = \arccos \left( \frac{\langle \mathbf{a}_i, \mathbf{a}_{i'} \rangle}{\|\mathbf{a}_i\| \|\mathbf{a}_{i'}\|} \right),$$

ou encore :

$$\theta_{ii'} = \arccos \left( \frac{\sum_{j=1}^p y_{ij} y_{i'j}}{\sqrt{\sum_{j=1}^p y_{ij}^2} \sqrt{\sum_{j=1}^p y_{i'j}^2}} \right) = \arccos \left( \frac{\sum_{s=1}^2 a_{is} a_{i's}}{\sqrt{\sum_{s=1}^2 a_{is}^2} \sqrt{\sum_{s=1}^2 a_{i's}^2}} \right).$$

Le respect des longueurs et des angles entraîne le respect des distances euclidiennes entre les points-lignes du biplot par rapport aux distances euclidiennes entre les lignes de  $Y$ .

Cette propriété du respect des distances euclidiennes et des angles n'est, par contre, pas vérifiée pour les points-colonnes, pour lesquels :

$$\mathbf{B}'_1 \mathbf{B}_1 = \mathbf{I} \neq \mathbf{Y}'\mathbf{Y}.$$

Par symétrie, la factorisation avec  $\alpha = 0$  conduit au résultat suivant :

$$\mathbf{B}_2 \mathbf{B}'_2 = (\mathbf{V}\mathbf{L})(\mathbf{V}\mathbf{L})' = \mathbf{Y}'\mathbf{Y}.$$

Pour cette factorisation, le biplot conserve les angles et les longueurs des vecteurs relatifs aux points-colonnes, mais pas des vecteurs relatifs aux points-lignes.

Enfin, la factorisation avec  $\alpha = 1/2$  ne respecte les longueurs et les angles, ni pour les vecteurs relatifs aux points-lignes, ni pour les vecteurs relatifs aux points-colonnes.

En conclusion, la factorisation avec  $\alpha = 1$  privilégie les lignes, la factorisation avec  $\alpha = 0$  privilégie les colonnes et la factorisation avec  $\alpha = 1/2$  donne un poids égal aux lignes et aux colonnes. Ces factorisations sont parfois appelées factorisation JK ou RMP (*row metric preserving*), GH ou CMP (*column metric preserving*) et SYM (*Symmetric*).

Lorsque  $\mathbf{Y}$  n'est pas une matrice de rang 2, nous avons vu que les biplots ne sont plus des représentations de  $\mathbf{Y}$  mais des représentations d'une approximation de rang 2 de  $\mathbf{Y}$ . Dans ces conditions, les longueurs et les angles ne sont plus exactement conservés pour les points-lignes de  $\mathbf{Y}$  lorsque  $\alpha = 1$ , ni pour les points-colonnes de  $\mathbf{Y}$  lorsque  $\alpha = 0$ .

Pour l'exemple considéré précédemment, et en négligeant la troisième valeur propre, on obtient pour  $\alpha = 1$  :

$$\mathbf{A}_1 = \begin{bmatrix} 19,02 & -3,05 \\ 9,53 & -3,18 \\ 14,13 & 0,24 \\ 27,26 & 3,11 \end{bmatrix} \quad \text{et} \quad \mathbf{B}_1 = \begin{bmatrix} 0,73 & -0,33 \\ 0,41 & 0,91 \\ 0,55 & -0,24 \end{bmatrix}$$

et pour  $\alpha = 0$  :

$$\mathbf{A}_2 = \begin{bmatrix} 0,51 & -0,56 \\ 0,26 & -0,59 \\ 0,38 & 0,04 \\ 0,73 & 0,58 \end{bmatrix} \quad \text{et} \quad \mathbf{B}_2 = \begin{bmatrix} 27,33 & -1,78 \\ 15,29 & 4,93 \\ 20,37 & -1,31 \end{bmatrix}.$$

Pour  $\alpha = 1/2$ , on obtient la factorisation donnée au paragraphe 2.1.

Les représentations sous forme de biplots pour les factorisations avec  $\alpha = 1$  et  $\alpha = 0$  sont données aux figures 2 et 3 et, pour  $\alpha = 1/2$ , à la figure 1, examinée au paragraphe 2.3.

On constate que, dans la figure 2, qui privilégie la représentation des lignes, les points-colonnes ont des coordonnées très faibles, ce qui rend le graphique peu

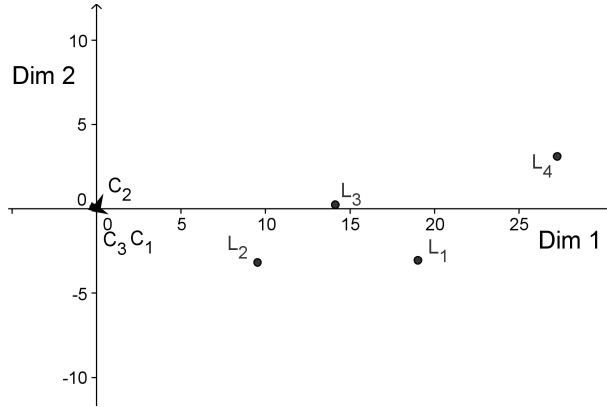


Figure 2 – Biplot de la matrice  $\hat{Y}$ , factorisation avec  $\alpha = 1$  ( $L_1$  à  $L_4$  : points-lignes ;  $C_1$  à  $C_3$  : points-colonnes).

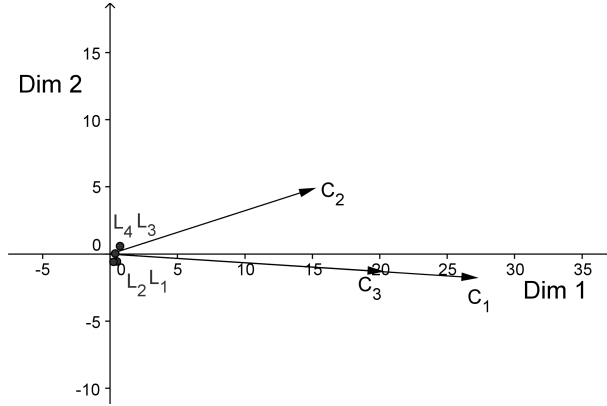


Figure 3 – Biplot de la matrice  $\hat{Y}$ , factorisation avec  $\alpha = 0$  ( $L_1$  à  $L_4$  : points-lignes ;  $C_1$  à  $C_3$  : points-colonnes) .

lisible pour les points-colonnes. On a, par contre, la situation inverse pour la figure 3, les points-lignes étant alors peu visibles. Enfin, dans la figure 1, qui consiste en un compromis, les deux types de points sont visibles simultanément. Cette discordance d'échelle pour les points-lignes et les points-colonnes lorsque  $\alpha = 1$  ou  $\alpha = 0$  sera systématique, dès que les valeurs singulières de  $Y$  seront différentes de l'unité.

L'examen de la figure 2, qui respecte les longueurs et les angles pour les points-lignes, montre par exemple que le vecteur pour  $L_4$  est à peu près deux fois plus long que le vecteur pour  $L_3$ . Cela signifie donc que la moyenne quadratique des données de la ligne 4 de  $\hat{Y}$  est approximativement le double de la moyenne quadratique des données de la ligne 3 de  $\hat{Y}$ .

Effectivement, les racines carrées des sommes des carrés des éléments des

lignes 3 et 4 de  $\hat{\mathbf{Y}}$  sont respectivement égales à (paragraphe 2.3) :

$$\sqrt{10,26^2 + 6,00^2 + 7,65^2} = 14,13$$

et

$$\sqrt{18,92^2 + 14,00^2 + 14,11^2} = 27,44,$$

ces valeurs étant égales aux longueurs des vecteurs des points-lignes, pour les lignes 3 et 4 de la matrice  $\mathbf{A}_1$  :

$$\sqrt{14,13^2 + 0,24^2} = 14,13$$

et

$$\sqrt{27,26^2 + 3,11^2} = 27,44.$$

On voit immédiatement aussi que les moyennes quadratiques des observations des lignes s'ordonnent de la manière croissante suivante : 2, 3, 1 et 4.

La figure 3, qui respecte les longueurs et les distances pour les points-colonnes, montre par exemple que la colonne 2 présente des valeurs globalement plus faibles que la colonne 1 et que la colonne 3 présente des valeurs intermédiaires, les carrés des longueurs des vecteurs dans la figure 3 étant égaux aux sommes des carrés des éléments des colonnes de  $\hat{\mathbf{Y}}$ . Par exemple, pour la première colonne de  $\hat{\mathbf{Y}}$  (paragraphe 2.3) et la première ligne de  $\mathbf{B}_2$ , on a :

$$14,92^2 + 8,02^2 + 10,26^2 + 18,92^2 = 27,33^2 + (-1,78^2) = 750,1.$$

Cette valeur est aussi très proche de la somme des carrés de la première colonne de  $\mathbf{Y}$ , puisque  $\hat{\mathbf{Y}}$  est très proche de  $\mathbf{Y}$  (paragraphe 3.2).

## 5. BIPLLOT ET ANALYSE EN COMPOSANTES PRINCIPALES

### 5.1. Présentation classique de l'analyse en composantes principales

L'analyse en composantes principales repose sur le calcul des valeurs et des vecteurs propres de la matrice des variances et covariances  $\mathbf{S}$  ou de la matrice de corrélation  $\mathbf{R}$  des variables. Les deux méthodes se différencient par la standardisation préliminaire des variables de départ. Ainsi, pour l'analyse en composantes principales basée sur la matrice des variances et covariances, on considère que la matrice de départ, que nous notons  $\mathbf{X}$ , est la matrice des observations centrées obtenues en retranchant de chaque observation d'une colonne, c'est-à-dire d'une variable, la moyenne de la colonne. Pour une analyse sur la matrice de corrélation, on considère que la matrice  $\mathbf{X}$  est constituée de variables centrées et réduites : on retranche des observations d'une colonne la moyenne de la colonne et on divise le résultat par l'écart-type de la colonne. Par la suite, nous considérons essentiellement le cas de la matrice de corrélation, qui correspond à la situation la plus fréquente en pratique.

Soit  $\mu_k^2$  et  $\mathbf{c}_k$  les valeurs propres et les vecteurs propres de la matrice  $\mathbf{R}$ . Soit  $\mathbf{C}$  la matrice obtenue par la juxtaposition des vecteurs propres. Les composantes principales sont obtenues par la relation :

$$\mathbf{Z} = \mathbf{XC},$$

et le graphique des individus dans les différents plans factoriels correspond au diagramme de dispersion de  $z_{ik}$  et  $z_{ik'}$  ( $k \neq k'$ ).

Pour la représentation des variables dans les cercles de corrélation, on calcule la matrice de corrélation  $\mathbf{W}$  entre les variables initiales et les composantes principales, en multipliant les vecteurs propres  $\mathbf{c}_k$  par la racine carrée de la valeur propre correspondante  $\mu_k$ , soit :

$$\mathbf{W} = \mathbf{CM},$$

$\mathbf{M}$  étant la matrice diagonale dont les éléments diagonaux sont les  $\mu_k$ . Dans le plan factoriel  $(k, k')$ , la variable  $j$  a comme coordonnées les valeurs  $w_{jk}$  et  $w_{jk'}$ .

A titre d'illustration, nous reprenons l'exemple proposé par HARTIGAN (1975), concernant les teneurs en protéines, graisse et lactose du lait de 16 mammifères. L'analyse en composantes principales de ces données a été détaillée antérieurement (PALM, 1998). Le lecteur y trouvera les données et l'ensemble des résultats de l'analyse, ainsi que des éléments d'interprétation. Nous n'envisageons ici que les aspects plus directement en relation avec les biplots.

Le tableau 1 reprend les données centrées et réduites notées  $x_{i1}$ ,  $x_{i2}$  et  $x_{i3}$ . La juxtaposition de ces trois colonnes donne la matrice  $\mathbf{X}$ . Le tableau reprend également les valeurs des composantes principales, notées  $z_{i1}$ ,  $z_{i2}$  et  $z_{i3}$ .

La matrice de corrélation des trois variables est égale à :

$$\mathbf{R} = \begin{bmatrix} 1,000 & 0,897 & -0,938 \\ 0,897 & 1,000 & -0,865 \\ -0,938 & -0,865 & 1,000 \end{bmatrix}.$$

Les valeurs propres de cette matrice valent :

$$\mu_1^2 = 2,800 \quad \mu_2^2 = 0,142 \quad \text{et} \quad \mu_3^2 = 0,058,$$

et la matrice des vecteurs propres s'écrit :

$$\mathbf{C} = \begin{bmatrix} -0,585 & -0,233 & 0,777 \\ -0,569 & 0,801 & -0,188 \\ 0,578 & 0,552 & 0,601 \end{bmatrix}.$$

Les valeurs des composantes principales du tableau 1 s'obtiennent par le produit  $\mathbf{XC}$  et le graphique des individus dans le premier plan factoriel est simplement le diagramme de dispersion de  $z_{i1}$  et  $z_{i2}$ .

Tableau 1 – Teneurs en protéines, graisse et lactose du lait de 16 mammifères : données centrées et réduites et valeurs des composantes principales.

Code	Nom	$x_{i1}$	$x_{i2}$	$x_{i3}$	$z_{i1}$	$z_{i2}$	$z_{i3}$
a	ânesse	-1,354	-1,024	1,263	2,105	0,193	-0,100
b	baleine	1,332	1,858	-1,615	-2,770	0,285	-0,285
c	biche	1,132	1,639	-0,990	-2,167	0,502	-0,023
d	brebis	-0,239	-0,297	0,325	0,496	-0,002	0,065
e	buffle	-0,154	-0,078	0,325	0,322	0,152	0,090
f	chamelle	-0,839	-0,733	0,387	1,132	-0,177	-0,282
g	cobaye	0,275	-0,180	-0,927	-0,594	-0,720	-0,310
h	jument	-1,096	-1,083	1,701	2,241	0,328	0,374
i	lama	-0,725	-0,762	0,888	1,371	0,049	0,113
j	lapine	1,675	0,679	-1,428	-2,191	-0,636	0,316
k	mule	-1,268	-0,966	0,825	1,768	-0,022	-0,308
l	rate	0,789	0,606	-0,551	-1,125	-0,004	0,168
m	renarde	0,046	-0,369	0,450	0,443	-0,058	0,376
n	renne	1,218	1,727	-1,052	-2,303	0,517	-0,010
o	truie	0,189	-0,486	-0,301	-0,008	-0,599	0,057
p	zèbre	-0,982	-0,529	0,700	1,280	0,192	-0,243

En multipliant  $\mathbf{C}$  par la matrice diagonale dont les éléments sont les  $\mu_k$ , on obtient les corrélations des variables initiales et des composantes principales :

$$\mathbf{W} = \begin{bmatrix} -0,585 & -0,233 & 0,777 \\ -0,569 & 0,801 & -0,188 \\ 0,578 & 0,552 & 0,601 \end{bmatrix} \begin{bmatrix} \sqrt{2,800} & 0 & 0 \\ 0 & \sqrt{0,142} & 0 \\ 0 & 0 & \sqrt{0,058} \end{bmatrix}$$

$$= \begin{bmatrix} -0,979 & -0,088 & 0,186 \\ -0,952 & 0,301 & -0,045 \\ 0,968 & 0,208 & 0,144 \end{bmatrix}.$$

Dans le premier plan factoriel, les variables protéines, graisse et lactose auront donc comme coordonnées les valeurs suivantes :

$$(-0,979, -0,088), (-0,952, 0,301) \text{ et } (0,968, 0,208).$$

## 5.2. Relations avec les représentations par biplots

Selon que  $\mathbf{X}$  est la matrice des variables centrées ou des variables centrées réduites, la quantité :

$$(\mathbf{X}'\mathbf{X})/(n-1)$$

représente la matrice de variances et covariances ou la matrice de corrélation. C'est donc la matrice dont on calcule les valeurs propres et les vecteurs propres.

Si on divise par  $\sqrt{n-1}$  tous les éléments de  $\mathbf{X}$  :

$$\mathbf{Y} = \mathbf{X}/\sqrt{n-1},$$

alors la matrice  $\mathbf{Y}'\mathbf{Y}$  est directement la matrice de variances et covariances ou la matrice de corrélation, selon la définition de  $\mathbf{X}$ .

Dès lors, si on effectue la décomposition par valeurs singulières de  $\mathbf{Y}$  :

$$\mathbf{Y} = \mathbf{ULV}',$$

la matrice  $\mathbf{L}$  est identique à la matrice  $\mathbf{M}$  et la matrice  $\mathbf{V}$  est identique à la matrice  $\mathbf{C}$ .

Il en résulte que le produit  $\mathbf{YV}$  est égal, à une constante près, à la matrice  $\mathbf{Z}$  :

$$\mathbf{YV} = \mathbf{Z}/\sqrt{n-1}.$$

En remplaçant  $\mathbf{Y}$  par sa décomposition par valeurs singulières et en tenant compte du fait que  $\mathbf{V}$  est une matrice orthonormée, on a :

$$\mathbf{Z}/\sqrt{n-1} = \mathbf{ULV}'\mathbf{V} = \mathbf{UL}.$$

Les coordonnées des individus dans les graphiques établis en analyse en composantes principales sont donc identiques, à une constante près, égale pour toutes les dimensions, aux coordonnées du biplot lors de la factorisation de  $\mathbf{Y}$  avec  $\alpha = 1$  (paragraphe 4.1).

Du fait de l'égalité des matrices  $\mathbf{C}$  et  $\mathbf{V}$  d'une part, et  $\mathbf{M}$  et  $\mathbf{L}$  d'autre part, la matrice  $\mathbf{W}$  qui donne les coordonnées des variables est égale à :

$$\mathbf{W} = \mathbf{VL} = \mathbf{CM}.$$

Les coordonnées des variables dans l'analyse en composantes principales sont donc égales aux coordonnées des points-colonnes du biplot obtenu par factorisation de  $\mathbf{Y}$  avec  $\alpha = 0$  (paragraphe 4.1).

Pour les données relatives aux mammifères, la matrice  $\mathbf{Y}$  s'écrit :

$$\mathbf{Y} = \begin{bmatrix} -0,350 & -0,264 & 0,326 \\ 0,344 & 0,480 & -0,417 \\ \dots & \dots & \dots \\ -0,254 & -0,137 & 0,181 \end{bmatrix}.$$

Elle comporte 16 lignes, mais seules les deux premières et la dernière lignes ont été reprises.

La décomposition par valeurs singulières de cette matrice donne les résultats suivants :

$$\mathbf{U} = \begin{bmatrix} 0,325 & -0,133 & 0,108 \\ -0,427 & -0,195 & 0,307 \\ \dots & \dots & \dots \\ 0,198 & -0,132 & 0,262 \end{bmatrix},$$



$$\mathbf{L} = \begin{bmatrix} 1,674 & 0 & 0 \\ 0 & 0,376 & 0 \\ 0 & 0 & 0,239 \end{bmatrix}$$

et

$$\mathbf{V} = \begin{bmatrix} -0,585 & 0,233 & -0,777 \\ -0,569 & -0,801 & 0,188 \\ 0,578 & -0,552 & -0,601 \end{bmatrix}.$$

On constate tout d'abord que la matrice  $\mathbf{V}$  est bien égale à la matrice des vecteurs propres de  $\mathbf{R}$ , notée  $\mathbf{C}$  au paragraphe 5.1, à condition toutefois de changer le signe des deux dernières colonnes. Cette inversion de signe résulte du fait qu'en analyse en composantes principales le signe des vecteurs propres est arbitraire et on peut toujours multiplier par  $-1$  tous les éléments d'un vecteur propre. Si les signes n'ont pas été changés au paragraphe 5.1, c'est pour respecter les résultats donnés dans l'étude antérieure [PALM, 1998]. On constate aussi que si on élève au carré les éléments diagonaux de  $\mathbf{L}$  on retrouve les valeurs propres  $\mu_k^2$  de la matrice  $\mathbf{R}$  (paragraphe 5.1).

D'autre part, si on multiplie la première colonne de  $\mathbf{U}$  par 1,674, la deuxième colonne de  $\mathbf{U}$  par 0,376 et la troisième colonne de  $\mathbf{U}$  par 0,239, on retrouve, au facteur  $\sqrt{15}$  près, les valeurs des composantes principales. Ainsi, pour le premier animal, on a :

$$\begin{aligned} (0,325)(1,674)(\sqrt{15}) &= 2,107 \\ (-0,133)(0,376)(\sqrt{15}) &= -0,194 \\ (0,108)(0,239)(\sqrt{15}) &= 0,100. \end{aligned}$$

Ces valeurs sont, aux erreurs d'arrondis près, égales aux valeurs des composantes principales données dans le tableau 1, si on modifie le signe des deux dernières composantes.

Enfin, si on multiplie la matrice  $\mathbf{V}$  par  $\mathbf{L}$ , c'est-à-dire si on multiplie la première colonne de  $\mathbf{V}$  par 1,674, la deuxième colonne de  $\mathbf{V}$  par 0,376 et la troisième colonne de  $\mathbf{V}$  par 0,239, on retrouve la matrice de corrélation  $\mathbf{W}$  du paragraphe 5.1, si on inverse les signes des deux dernières colonnes.

### 5.3. Distances euclidiennes et distances de MAHALANOBIS

Nous avons vu, au paragraphe 4.2, que la factorisation de  $\mathbf{Y}$  après décomposition par valeurs singulières avec  $\alpha = 1$  conduit à une représentation graphique qui respecte les longueurs pour les points-lignes et les distances euclidiennes entre les points-lignes. Ainsi, pour le premier individu, la somme des carrés des éléments de la première ligne de  $\mathbf{Y}$  est égale à :

$$(-0,350)^2 + (-0,264)^2 + 0,326^2 = 0,298,$$

et la somme des différences entre les deux premières lignes vaut :

$$(-0,350 - 0,344)^2 + (-0,264 - 0,480)^2 + (0,326 + 0,417)^2 = 1,587.$$

La première valeur est égale à la somme des carrés des éléments de la première ligne de  $\mathbf{UL}$ , soit :

$$[(0, 325)(1, 674)]^2 + [(-0, 133)(0, 376)]^2 + [(0, 108)(0, 239)]^2 = 0, 299.$$

La deuxième valeur est égale à la somme des carrés des différences entre les deux premières lignes de  $\mathbf{UL}$  :

$$[(0, 325 + 0, 427)(1, 674)]^2 + [(-0, 133 + 0, 195)(0, 376)]^2 + [(0, 108 - 0, 307)(0, 239)]^2 = 1, 588.$$

Ces résultats montrent bien que, lorsque toutes les composantes principales sont retenues, les longueurs et les distances euclidiennes entre individus sont préservées. Bien entendu, si on néglige une ou plusieurs composantes, ce sont les distances pour le tableau approché  $\mathbf{Y}$  qui sont conservées.

Si, pour la représentation des individus, on avait retenu la factorisation avec  $\alpha = 0$ , les longueurs et les distances euclidiennes entre individus n'auraient pas été respectées. Par contre, on peut montrer que les carrés des longueurs des vecteurs relatifs aux points-lignes sont, au facteur  $n - 1$  près, les carrés des distances au sens de MAHALANOBIS des individus au centre de gravité. On peut démontrer en effet que, pour un individu  $i$  donné :

$$\mathbf{u}_i \mathbf{u}_i' = \frac{1}{n-1} \mathbf{y}_i \mathbf{S}^{-1} \mathbf{y}_i',$$

$\mathbf{S}$  étant la matrice de variances et covariances des données. Contrairement aux distances euclidiennes, les distances de MAHALANOBIS tiennent compte des variances et covariances entre variables. Elles sont largement utilisées en analyse multivariée, en relation avec l'hypothèse de multinormalité des populations.

De même, les distances euclidiennes entre deux points-lignes  $\mathbf{u}_i$  et  $\mathbf{u}_{i'}$  dans le biplot établi pour  $\alpha = 0$  sont, à une constante près, égales au carré des distances au sens de MAHALANOBIS entre les individus  $\mathbf{y}_i$  et  $\mathbf{y}_{i'}$  :

$$(\mathbf{u}_i - \mathbf{u}_{i'}) (\mathbf{u}_i - \mathbf{u}_{i'})' = \frac{1}{n-1} (\mathbf{y}_i - \mathbf{y}_{i'}) \mathbf{S}^{-1} (\mathbf{y}_i - \mathbf{y}_{i'})'.$$

Pour les données relatives aux mammifères, le carré de la longueur du vecteur-ligne pour le premier individu est égal à :

$$0, 325^2 + (-0, 133)^2 + (0, 108)^2 = 0, 135,$$

et le carré de la distance euclidienne entre les deux premiers individus vaut :

$$(0, 325 + 0, 427)^2 + (-0, 133 + 0, 195)^2 + (0, 108 - 0, 307)^2 = 0, 609.$$

Si on calcule le carré de la distance de MAHALANOBIS pour le premier individu et le carré de la distance de MAHALANOBIS entre les deux premiers

individus, on trouve respectivement 2,020 et 9,137. En divisant par 15 ces deux valeurs, on retrouve bien 0,135 et 0,609.

Considérons maintenant la représentation des variables réalisée lors de l'analyse en composantes principales. Nous avons vu que cette représentation correspond à la représentation des points-colonnes du biplot pour  $\alpha = 0$  (paragraphe 5.2). Dans l'espace factoriel complet, le carré de la distance d'un point-colonne  $j$  à l'origine des axes correspond à l'élément  $jj$  de la matrice (paragraphe 4.2) :

$$\mathbf{W}\mathbf{W}' = \mathbf{V}\mathbf{L}(\mathbf{V}\mathbf{L})' = \mathbf{Y}'\mathbf{Y}.$$

Si l'analyse en composantes principales est réalisée sur la matrice de corrélation, alors (paragraphe 5.2) :

$$\mathbf{Y}'\mathbf{Y} = \mathbf{R},$$

et les éléments diagonaux sont égaux à l'unité. Dans l'espace factoriel complet, les points représentant les variables se trouvent donc sur une hypersphère de rayon unitaire. Dans un sous-espace factoriel de dimension 2, ils se trouvent sur ou à l'intérieur du cercle de rayon unitaire, appelé cercle de corrélation.

Si on examine l'angle formé par deux vecteurs-colonnes  $j$  et  $j'$ , on a (paragraphe 2) :

$$\cos \theta_{jj'} = \frac{\sum_{s=1}^r w_{js} w_{j's}}{\sqrt{\sum_{s=1}^r w_{js}^2 \sum_{s=1}^r w_{j's}^2}} = \sum_{s=1}^r w_{js} w_{j's}.$$

puisque, comme on vient de le signaler, la longueur des vecteurs relatifs aux points-colonnes est égale à l'unité. Il en résulte que  $\cos \theta_{jj'}$  est l'élément  $jj'$  de la matrice  $\mathbf{W}\mathbf{W}'$  et donc de  $\mathbf{R}$ , c'est-à-dire aussi la corrélation des variables  $j$  et  $j'$ . Deux points sont confondus si les variables correspondantes ont une corrélation égale à l'unité : ils forment un angle droit si la corrélation est nulle et ils forment un angle de 180 degrés si la corrélation est égale à  $-1$ . L'égalité entre le cosinus de l'angle et la corrélation des variables n'est cependant vérifiée que pour l'espace factoriel complet.

Les propriétés énoncées ci-dessus peuvent être appliquées aux données relatives aux mammifères. Pour les deux premières variables, c'est-à-dire pour les deux premières lignes de  $\mathbf{V}\mathbf{L}$  prises comme exemple, on a bien :

$$(-0,979)^2 + 0,088^2 + (-0,186)^2 \simeq 1,$$

$$(-0,952)^2 + (-0,301)^2 + 0,045^2 \simeq 1,$$

et

$$(-0,979)(-0,952) + (0,088)(-0,301) + (-0,186)(0,045) = 0,897,$$

cette dernière valeur étant la corrélation entre les deux premières variables, c'est-à-dire entre les teneurs en protéines et en graisse.

#### 5.4. Représentation simultanée des variables et des individus.

En analyse en composantes principales, les individus et les variables font le plus souvent l'objet de représentations graphiques séparées. Et nous avons vu, au paragraphe 5.2, que le graphique des individus est équivalent à la représentation des lignes dans le biplot avec  $\alpha = 1$  et que le cercle de corrélation correspond à la représentation des colonnes dans le biplot avec  $\alpha = 0$ .

Une représentation simultanée des individus et des variables par le biplot avec  $\alpha = 1$  peut cependant présenter un intérêt pratique, si on calibre les axes correspondant aux variables.

La figure 4 donne le biplot pour l'exemple des mammifères. Les individus sont identifiés par leur code (tableau 1) et les variables par les symboles P pour protéines, G pour graisse et L pour lactose.

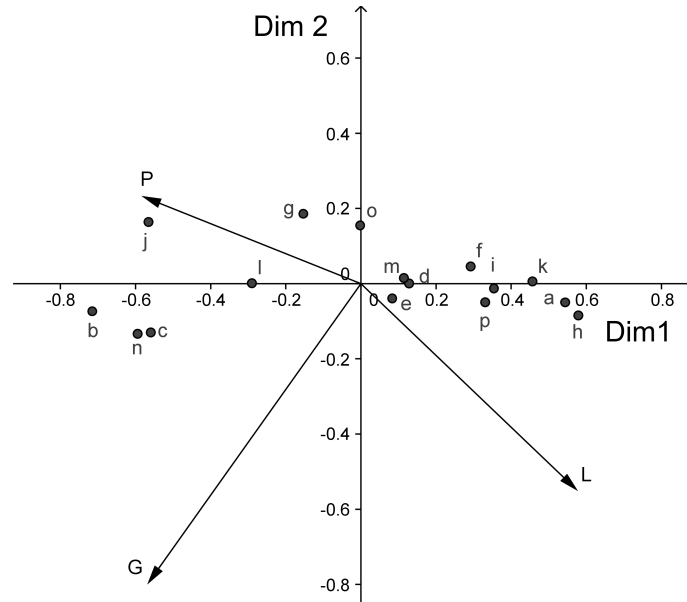


Figure 4 – Biplot de la matrice des données relatives aux mammifères, factorisation de  $\mathbf{Y}$  avec  $\alpha = 1$  ( $a, b, \dots, p$  : code des animaux ; P : protéines, G : graisse et L : lactose).

A la figure 5, les vecteurs représentant les variables ont été remplacés par des axes gradués et, pour ne pas alourdir le graphique, seule l'ânesse, prise à titre d'exemple, a été identifiée par le symbole  $a$ . Si on projette un individu perpendiculairement sur un de ces axes, on obtient, comme coordonnée sur cet axe, la valeur approchée de la variable pour cet individu, que nous avons notée précédemment  $\hat{y}_{ij}$ . Ainsi par exemple, pour l'ânesse, on a respectivement pour l'axe "Protéines", "Graisse" et "Lactose", les valeurs suivantes :

$$\hat{y}_{11} = -0,329, \quad \hat{y}_{12} = -0,269 \quad \text{et} \quad \hat{y}_{13} = 0,342.$$

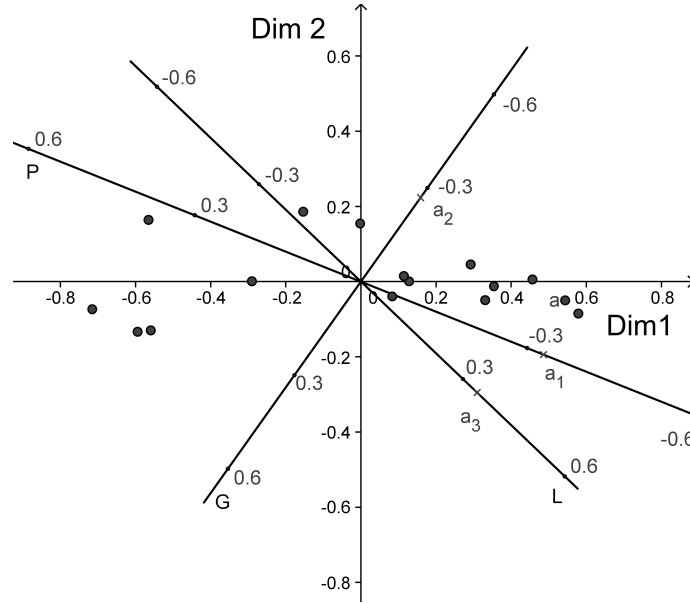


Figure 5 – Biplot de la matrice des données relatives aux mammifères, avec graduation des axes relatifs aux variables ( $a$  : ânesse;  $a_1$  : projection de  $a$  sur  $P$ ;  $a_2$  : projection de  $a$  sur  $G$ ;  $a_3$  : projection de  $a$  sur  $L$ ).

Ces projections ont été matérialisées par les symboles  $a_1$ ,  $a_2$  et  $a_3$ . Les valeurs des projections sont très proches des valeurs de la première ligne de la matrice  $\mathbf{Y}$  donnée au paragraphe 5.2 :

$$y_{11} = -0,350, \quad y_{12} = -0,264 \quad \text{et} \quad y_{13} = 0,326,$$

car, pour cet exemple, les deux premières composantes expliquent 98 % de la variabilité et la matrice  $\mathbf{Y}$  est proche de  $\hat{\mathbf{Y}}_{(2)}$ .

La figure 5 peut donc être interprétée comme une extension multivariée d'un diagramme de dispersion : au lieu de deux axes perpendiculaires, on a plusieurs axes, nécessairement non perpendiculaires. Un tel graphique n'est cependant utile que si la matrice  $\hat{\mathbf{Y}}_{(2)}$  est suffisamment proche de  $\mathbf{Y}$ , comme c'est le cas pour cet exemple.

Des informations concernant la détermination des graduations sur les axes sont données par GOWER et *al.* (2011).

## 6. INFORMATIONS COMPLÉMENTAIRES

Dans cette note, nous avons présenté un outil permettant la visualisation graphique d'une matrice, généralement dans un espace de dimension réduite. La particularité de l'outil est de donner une représentation simultanée des lignes et des colonnes de la matrice, ce qui justifie le nom de biplot donné au graphique. La méthode repose sur la factorisation de la matrice sous la forme de deux matrices.

Cette factorisation est optimisée, de manière à ce que le sous-espace de dimension réduite donne la meilleure approximation de la matrice. L'objectif est atteint par l'utilisation de la décomposition par valeurs singulières :

$$\mathbf{Y} = \mathbf{ULV}'.$$

Les lignes et les colonnes de la matrice initiale  $\mathbf{Y}$  sont alors représentées par des points dont les coordonnées sont les éléments des lignes et des matrices :

$$\mathbf{UL}^\alpha \text{ et } \mathbf{VL}^\beta.$$

Pour que le produit scalaire d'un vecteur correspondant à un point-ligne  $i$  et à un point-colonne  $j$  soit égal à l'élément  $ij$  de  $\mathbf{Y}$ , il faut que la somme des exposants de  $\mathbf{L}$  dans les deux matrices ci-dessus soit égal à l'unité, c'est-à-dire que  $\beta$  soit égal à  $1 - \alpha$ .

Dans la pratique cependant, des représentations graphiques à partir de matrices pour lesquelles  $\alpha + \beta \neq 1$  sont aussi utilisées, comme nous l'avons vu pour l'analyse en composantes principales, où  $\alpha = \beta = 1$ .

Une généralisation du biplot aux situations pour lesquelles  $\alpha + \beta \neq 1$  est étudiée par GOWER [2004] et la qualité de la représentation du biplot pour différents choix de  $\alpha$  et  $\beta$  a été examinée par GABRIEL [2002].

Dans cette note, nous nous sommes volontairement limités à une introduction aux biplots et à leur relation avec l'analyse en composantes principales. Il existe cependant aussi des relations entre les représentations par biplots et d'autres méthodes statistiques multivariées. Le lecteur trouvera des informations complémentaires dans les livres de GOWER et HAND [1996] et de GOWER *et al.* [2011].

## BIBLIOGRAPHIE

- ECKART C., YOUNG G. [1936]. The approximation of a matrix by another of lower rank. *Psychometrika* 1, 211-218.
- GABRIEL K. R. [1971]. The biplot graphic display of matrices with application to principal component analysis. *Biometrika*, 58(3), 453-467.
- GABRIEL K. R. [2002]. Goodness of fit of biplots and correspondance analysis. *Biometrika* 89(2), 423-436.
- GOWER J. C., HAND D. J. [1996]. *Biplots*. New York, Chapman and Hall, 277 p.
- GOWER J. C., LUBBE S., ROUX N. [2011]. *Understanding biplots*. New York, Wiley, 463 p.
- HARTIGAN J. A. [1975]. *Clustering algorithms*. New York, Wiley, 351 p.
- JACKSON J. E. [1991]. *A user's guide to principal components*. New York, Wiley, 569 p.
- KHATTREE R., NAIK D. N. [1995]. *Applied multivariate statistics with SAS software*. Cary N. C., SAS Institute Inc., 396 p.

- PALM R. [1998]. L'analyse en composantes principales : principes et applications.  
*Notes Stat. Inform.* Gembloux, 98/2, 31 p.
- RENCHEA A. C. [2002]. *Methods of multivariate analysis*. New York, Wiley,  
708 p.
- SEBER G. A. F. [1984]. *Multivariate observations*. New York, Wiley, 686 p.

La collection

### ***NOTES DE STATISTIQUE ET D'INFORMATIQUE***

réunit divers travaux (documents didactiques, notes techniques, rapports de recherche, publications, etc.) émanant de l'Unité de Statistique, Informatique et Mathématique appliquées à la bioingénierie de l'Université de Liège – Gembloux Agro-Bio Tech et de l'Unité Systèmes agraires, Territoire et Technologies de l'Information du Centre wallon de Recherches agronomiques (Gembloux - Belgique).

La liste des notes disponibles peut être obtenue sur simple demande à l'adresse ci-dessous :

*Université de Liège – Gembloux Agro-Bio Tech  
Unité de Statistique, Informatique et Mathématique appliquées à la bioingénierie  
Avenue de la Faculté d'Agronomie, 8  
B-5030 GEMBLoux (Belgique)  
E-mail : [sima.gembloux@ulg.ac.be](mailto:sima.gembloux@ulg.ac.be)*

Plusieurs notes sont directement accessibles à l'adresse Web suivante, section Publications :

*<http://www.gembloux.ulg.ac.be/si/>*

En relation avec certaines notes, des programmes spécifiques sont également disponibles à la même adresse, section Macros.

Quelques titres récents sont cités ci-après :

- PALM R., BROSTAUX Y. [2009]. Etude des séries chronologiques par les méthodes de décomposition. *Notes Stat. Inform.* (Gembloux) 2009/1, 17 p.
- CHARLES C. [2011]. Introduction aux ondelettes. *Notes Stat. Inform.* (Gembloux) 2011/1, 22 p.
- CHARLES C. [2011]. Introduction aux applications des ondelettes. *Notes Stat. Inform.* (Gembloux) 2011/2, 35 p.
- PALM R., BROSTAUX Y. et CLAUSTRIAUX J. J. [2011]. Macros Minitab pour le choix d'une transformation pour la normalisation de variables. *Notes Stat. Inform.* (Gembloux) 2011/3, 15 p.
- PALM R., BROSTAUX Y. [2011]. La régression logistique avec Minitab. *Notes Stat. Inform.* (Gembloux) 2011/4, 15 p.
- PALM R., BROSTAUX Y., CLAUSTRIAUX J. J. [2011]. Inférence statistique et critères de qualité de l'ajustement en régression logistique binaire. *Notes Stat. Inform.* (Gembloux) 2011/5, 32 p.
- CLAUSTRIAUX J. J., PALM R., FERRANDIS-VALLTERRA S., BROSTAUX Y. et PLANCHON V. [2012]. Tables de contingence à trois dimensions : aspects théoriques, applications et analogie avec l'analyse de la variance à trois critères de classification. *Notes Stat. Inform.* (Gembloux) 2012/1, 19 p.