

# ANALYSE en COMPOSANTES PRINCIPALES

1. Cas général de l'ACP
2. ACP bipondérée
3. ACP simple (covariances)
4. ACP standard (corrélations)

# 1 ACP : Cas général

Protocole multivarié de  $K$  variables sur  $(J, n_J) : (x^J k)_{k \in K}$   
 Remplacer les  $K$  variables par  $L'$  nouvelles variables

## 1.1 Du protocole multivarié au nuage de points

Nuage de points  $(M^J, n_J)$  dans un espace affn  $\mathcal{U}$  de dimension  $K$   
 muni du repère cartésien  $(O, (\delta_k)_{k \in K})$ .

Le profil  $x^{jK} = (x^{jk})_{k \in K}$  de  $j$  est représenté, par le point pondéré  $(M^j, n_j) : \overline{OM^j} = \sum_{k \in K} x^{jk} \delta_k$ .

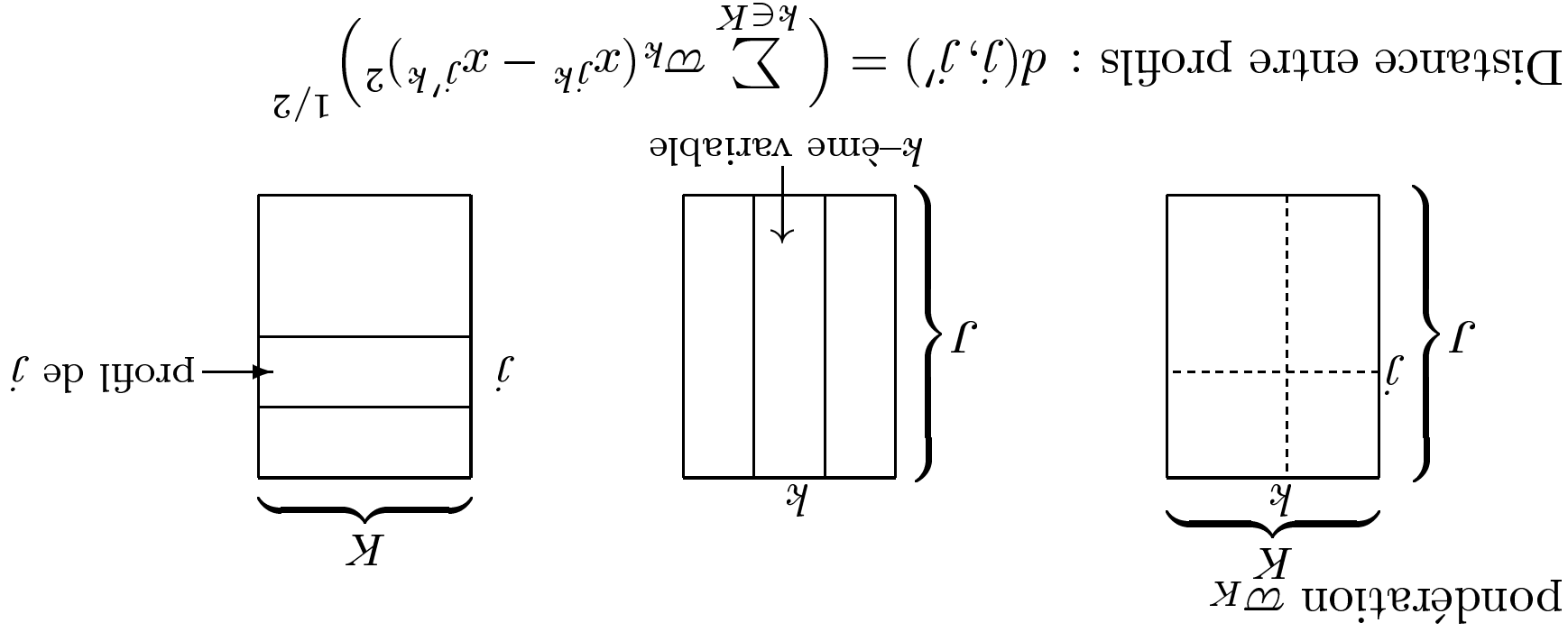
Nuage euclidien :  $(M^j M^{j'})_2 = \sum_{k \in K} \sum_{k' \in K} q_{kk'} (x^{jk} - x^{j'k'} - x^{j'k} + x^{jk'})$

avec  $q_{kk'} = \langle \delta_k^K | \delta_{k'}^K \rangle$

Effectuer l'*analyse en composantes principales* du protocole  
multivarie  $x_{JK}$  c'est déterminer les variables principales du nuage  
euclidien  $(M_J, n_J)$ , selon les méthodes présentées au chapitre  
NUAGE.

## 2 ACP bipondérée

Protocole de notes  $x^{JK}$ , muni de la mesure-effectifs  $n_j$  et de la



Espace affiné  $\mathcal{U}$  de dimension  $K$  muni du repère cartésien  $(O, (\delta_k)_{k \in K})$

$$\overrightarrow{OM_j} = \sum_{j \in J} x_{jk} \delta_k$$

Point moyen  $G = \sum_{j \in J} f_j M_j$ , d'où  $\overrightarrow{OG} = \sum_{i \in I} f_i \overrightarrow{OM_i} = \sum_{k \in K} x_k \delta_k$

$$(\overrightarrow{M_j M_{j'}})_2 = \sum_k \varpi_k (x_{jk} - x_{j'k})^2$$

Pour  $k \neq k'$ , on a  $\delta_k \perp \delta_{k'}$  et  $\|\delta_k\| = \sqrt{\varpi_k}$ .

La matrice  $\mathbf{Q}$  des produits scalaires est diagonale, notée  $\mathbf{\Omega}_K$ , avec  $q_{kk} = \varpi_k$  et  $q_{kk'} = 0$  pour  $k \neq k'$ .

Variance du nuage :  $\sum_{k \in K} w_k \text{Var } x_{J_k}$ .

Contribution de la variable  $x_{J_k}$  à la variance du nuage :  $\text{Ct}_{ak} = w_k \text{Var } x_{J_k}$ .

Nous appellerons l'ACP d'un tel protocole : ACP *pondérée*.

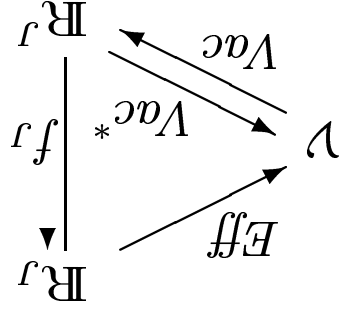
## 2.1 Directions et variables principales

*Formules de passage*

$$\vec{\delta}_k : \delta_k \mapsto \left( \langle \vec{\delta}_k | \vec{GM}^j \rangle \right)_{j \in J} = \varpi_k x_0^{jk}$$

$$\text{Vac}^* : \delta_j^j \mapsto \overrightarrow{f_j \text{GM}^j} = \sum_{k \in K} ((f_j x_0^{jk}) \delta_j^j)$$

$$\text{Som} = \text{Vac}^* \circ \text{Vac} : \vec{\delta}_k \mapsto \sum_{k' \in K} \varpi_k \varpi_{k'} \vec{\delta}_{k'}$$



Variable  $a_K = \sum_{k \in K} a_k \delta_k^k \mapsto$  vecteur géométrique  $\vec{\alpha} = \sum_{k \in K} a_k \delta_k^k$ .

Equation aux directions et valeurs propres :  $\overrightarrow{\text{Som}}(\vec{\alpha}) = \vec{\lambda \alpha}$  soit  $\sum_{k' \in K} \varpi_{k'} a_{k'} = \lambda a_k$  avec :  $a_k = a_k \varpi_k$



Variable principale calibrée (formule de passage) :

$$y_{\ell}^j = \sum_{k \in K} a_{k\ell} x_{0k}^j, \text{ avec } \sum_{k \in K} w_k (a_{k\ell}^2) = \sum_{k \in K} \frac{w_k}{a_{k\ell}^2} = 1$$

$y_{\ell}^j$  est la coordonnée du point  $M^j$  sur l'axe principal  $\ell$   
 Moy  $y_{\ell}^j = 0$  et  $\text{Var } y_{\ell}^j = \lambda_{\ell}$

## 2.2 Formules de reconstruction

Reconstruction d'ordre  $L'$  du tableau des variables centrées :

$$\tilde{x}_{L'}^{jk} = \sum_{\ell=1}^{L'} a_{\ell}^k y_{\ell}^j \quad \text{avec} \quad \sum_{k \in K} w_k (a_{\ell}^k)^2 = 1$$

Reconstruction des distances :

$$(GM^j)_2 = \sum_{k \in K} w_k (x^{jk} - \bar{x}^k)^2 = \sum_{\ell=1}^{L'} (y_{\ell}^j)^2$$

Variance du nuage :  $\text{Var } M^j = \sum_{k \in K} w_k \text{Var } x^{jk} = \sum_{\ell=1}^{L'} \lambda_{\ell}$

## 2.3 Espace des variables : caractérisation statistique

En termes statistiques, le problème de l'ACP pondérée se formule :  
*Problème* : parmi les variables  $\sum_{k \in K} a_k x_{Jk}$ , combinaisons linéaires des

variables initiales  $(x_{Jk})_{k \in K}$ , chercher celle(s) qui vérifie(nt) :  

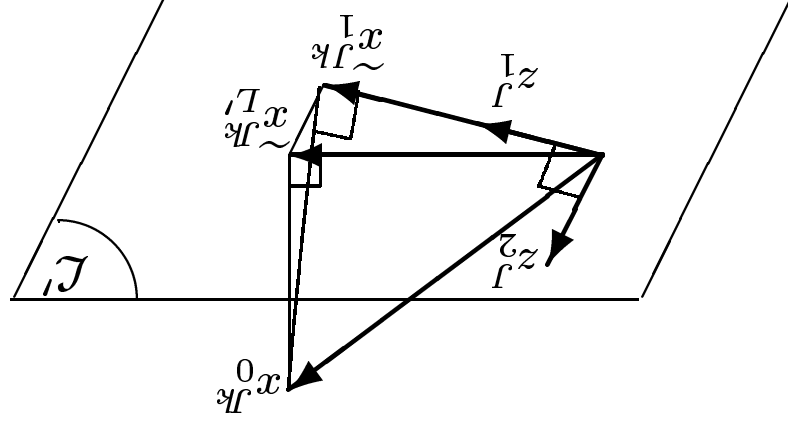
$$\text{Var} \left( \sum_{k \in K} a_k x_{Jk} \right) / \left( \sum_{k \in K} a_k^2 / \varpi_k \right)$$
 maximum.

Chercher  $\text{Var} \left( \sum_{k \in K} a_k x_{Jk} \right) / \left( \sum_{k \in K} a_k^2 / \varpi_k \right)$  maximum revient à

déterminer la première variable principale, puisque cette statistique est la variance du nuage dans la direction  $\vec{\alpha}$ .

## 2.4 Espace des variables

Régression des variables initiales sur les variables principales.



Formule de reconstitution :  $x_{jk}^0 = \sum_{\ell=1}^L b_{k\ell}^j z_{\ell}^j$  avec  $b_{k\ell}^j = a_{k\ell}^j \xi_{\ell}$

Le coefficient de régression partielle de  $x_{jk}^0$  est le coefficient de régression simple  $b_{k\ell}^j = \text{Cov}(x_{jk}^0 | z_{\ell}^j)$ , avec  $\text{Var } x_{jk}^0 = \sum_{\ell=1}^L (b_{k\ell}^j)^2$

**Corrélation.**

$$r_k^\ell = \text{Corr}(x_{Jk} | z_\ell^J) = b_k^\ell / \text{Ety } x_{Jk}$$

$$\sum_{\ell=1}^L (r_k^\ell)^2 = 1$$

$r_k^\ell$  = au cosinus de l'angle entre les variables (vecteurs)  $x_{Jk}^0$  et  $z_\ell^J$

Corrélation multiple  $R$  entre  $x_{Jk}$  et les  $L'$  premières variables

principales :  $R^2 = \sum_{\ell=1}^{L'} (r_k^\ell)^2 = \sum_{\ell=1}^{L'} (b_k^\ell)^2 / \text{Var } x_{Jk}$  ;

## 2.4.1 Variable moyenne

Posons  $\varpi = \sum_{k \in K} \varpi_k$  et notons  $\hat{x}^j$  la  $\varpi_k$ -moyenne des notes de

l'individu  $j$ , d'où la variable des moyennes, ou *variable moyenne*  
 $\hat{x}^j = (x^j)_{j \in J}$  avec :  $\hat{x}^j = \sum_{k \in K} \varpi_k x^{jk} / \varpi$ .

On a : Moy  $\hat{x}^j = \sum_{k \in K} \varpi_k x^k / \varpi$ , et  $\text{Var } \hat{x}^j = \sum_{k \in K} \sum_{k' \in K} \varpi_k \varpi_{k'} / \varpi^2$ .

La variable moyenne (centrée) est *principale* si et seulement si  $\text{Cov}(x^{jk} | x^j)$  ne dépend pas de  $k$ .

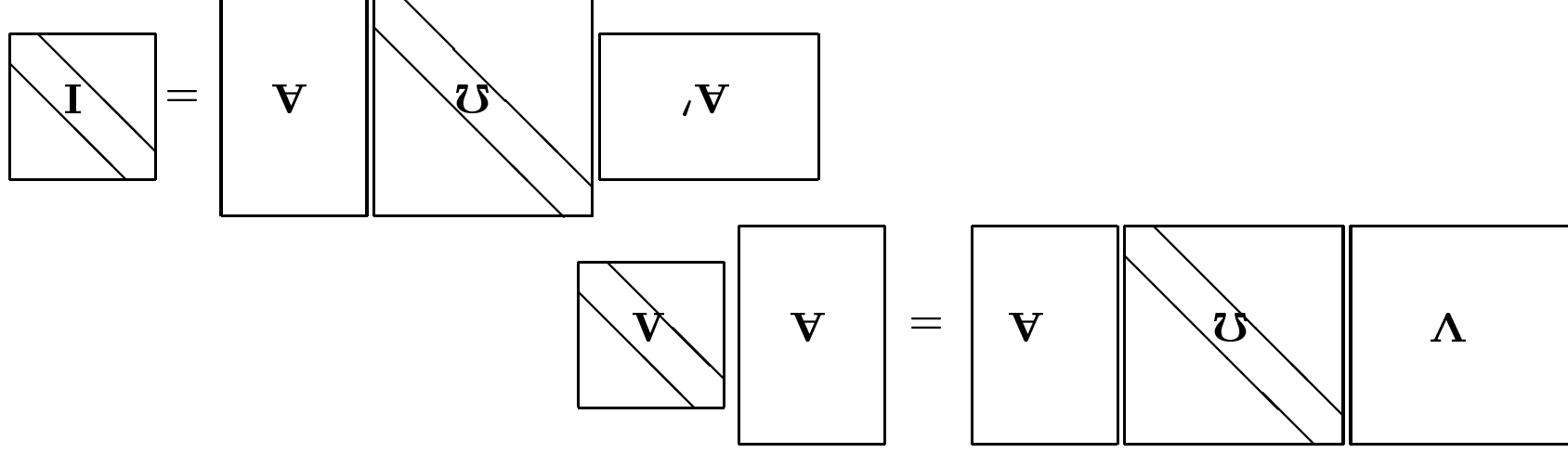
La variable principale calibrée est alors  $y^j = \sqrt{\varpi} \hat{x}^j$ , associée à la valeur propre  $\lambda_\ell = \varpi \text{Var } \hat{x}^j$ .

Si  $\hat{x}_0^j$  est variable principale associée à la valeur propre  $\lambda_\ell$ , toute variable principale associée à  $\lambda_{\ell'} \neq \lambda_\ell$  est l'*effet d'un contraste* sur  $K$  appliqué aux  $K$  variables initiales.

## 2.5 Formulaire : diagrammes résumés des formules principales

Equation aux directions et valeurs propres :

$$\sum_{k \in K} w_k (a_k^\lambda)^2 = 1 \quad \text{avec} \quad \sum_{k' \in K} v_{k'} (w_{k'} a_{k'}^\lambda) = \lambda a_k^\lambda$$



$$\begin{array}{c}
 \boxed{A'} \quad \boxed{\text{diag}(A)} \quad \boxed{A} = \boxed{V} \\
 \sum_{\ell=1}^L \lambda_{k\ell} a_{k\ell} = v_{kk'}
 \end{array}$$

$$\begin{array}{c}
 \boxed{\Xi} \quad \boxed{A} = \boxed{B} \\
 b_k^\ell = \xi_{k\ell} a_k^\ell
 \end{array}$$

$$\begin{array}{c}
 \boxed{A'} \quad \boxed{Y} = \boxed{X_0} \\
 \sum_{\ell=1}^L a_{k\ell} y_{j\ell} = x_{jk}^0
 \end{array}$$

$$\begin{array}{c}
 \boxed{A} \quad \boxed{\text{diag}(U)} \quad \boxed{X_0} = \boxed{Y} \\
 y_j^\ell = \sum_{k \in K} x_{jk}^0 (u_{k\ell})
 \end{array}$$

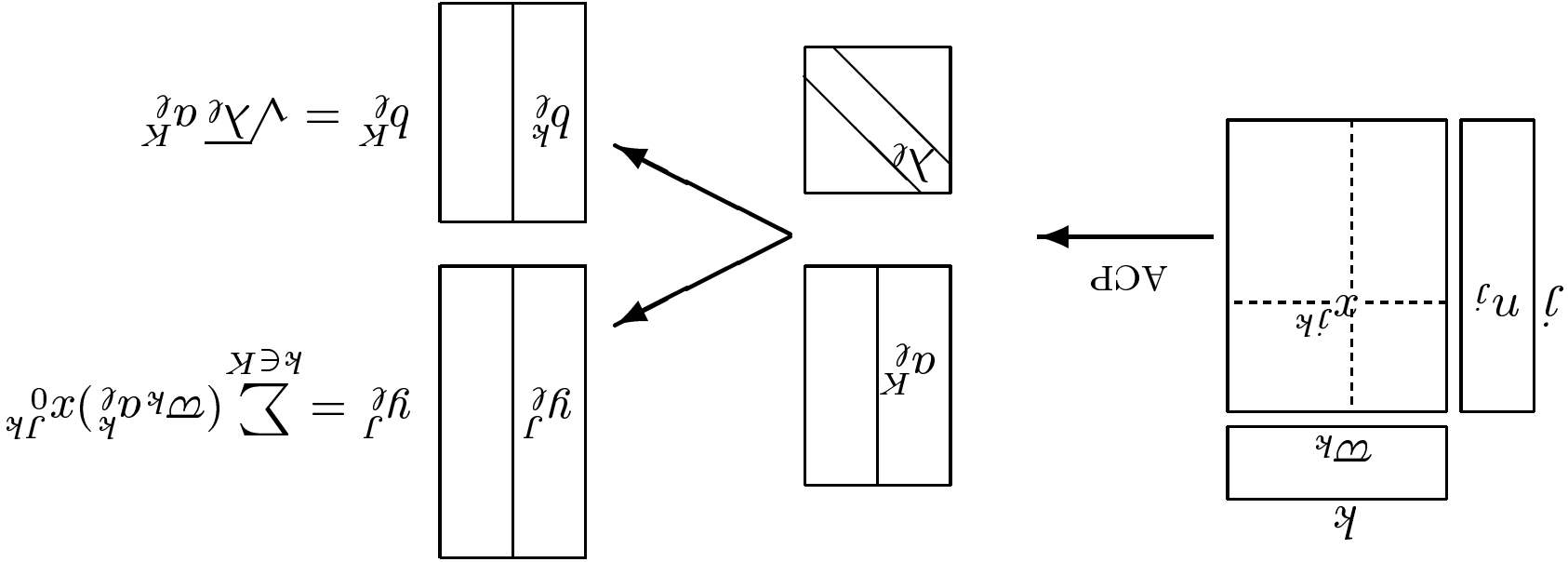
Formules de reconstitution

Coordonnées principales



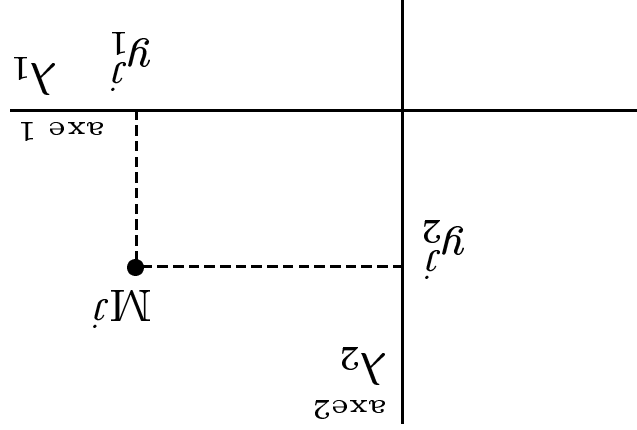
## 2.6 Principaux résultats

- vecteurs et valeurs propres  $\lambda_\ell$  ;
- coordonnées  $(y_\ell^j)$  des  $J$  / points du nuage sur les  $L'$  ( $L' \leq L$ ) axes principaux retenus pour l'interprétation ;
- coefficients de régression  $(b_k^\ell)$  des  $K$  / variables initiales sur les  $L'$  premières variables principales réduites.



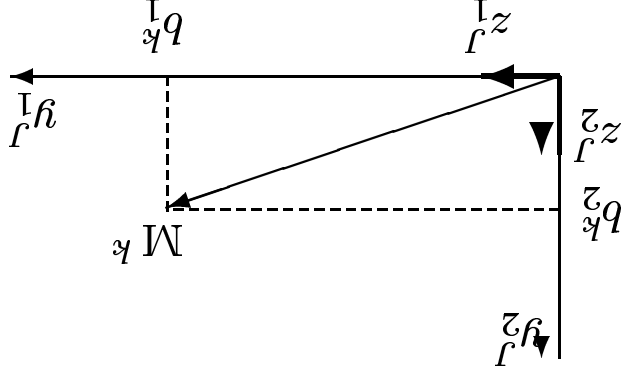
## 2.7 Représentations graphiques

*Espace d'observables*, projection du nuage des individus sur les premiers axes 1, 2, etc., ou dans les plans principaux 1-2, 1-3, 2-3, etc.  $M^j$  est représenté par le point d'abscisse  $y_1^j$  et d'ordonnée  $y_2^j$ , etc.



*Espace des variables*

Représentation des régressions des  $K$  variables initiales centrées variables principales 1, 2, 3, etc. Coordonnées :  $(b_1^k, b_2^k)$  par rapport à  $(z_1^j, z_2^j)$ .



**Procédure de calcul.**

1. On calcule la matrice  $\mathbf{S}$  de terme général  $s_{kk'} = \sqrt{w_k w_{k'}} v_{kk'}$ .
2. On diagonalise  $\mathbf{S}$ , d'où :  $(\lambda_\ell, (c_{k\ell})_{k \in K})_{\ell=1, \dots, L}$ , avec  $\sum_{k \in K} c_{k\ell}^2 = 1$ .
3. Pour  $\ell = 1, \dots, L'$  ( $L' \leq L$ ) et  $j \in J$ , on calcule :  $y_\ell^j = \sum_{k \in K} \sqrt{w_k} c_{k\ell} x_0^{jk}$ .
4. Pour  $\ell = 1, \dots, L'$  et  $k \in K$ , on calcule :  $b_\ell^k = \sqrt{\lambda_\ell} c_{k\ell} / \sqrt{w_k}$ .<sup>a</sup>

---

<sup>a</sup>Un programme d'ACP pondérée, réalisée avec P. Bonnet, est disponible auprès des auteurs.

**Formulations matricielles.**  $\mathbf{c}_\ell = K$ -colonne  $(c_{k\ell})_{k \in K}$ ,

$\mathbf{C}$  : matrice  $K \times L'$  des  $K$ -colonnes  $(\mathbf{c}_\ell)_{\ell=1, \dots, L'}$  ;

$y_\ell$  la  $J$ -colonne des coordonnées principales  $(y_\ell^j)_{j \in J}$  ;

$\mathbf{b}_\ell$  la  $K$ -colonne des coefficients de régression  $(b_k^\ell)_{k \in K}$  ;

$\mathbf{\Omega}_{1/2}^K$  la matrice diagonale  $(\sqrt{w_k})_k \in K$  ; on a :

$$1. \mathbf{S} = \mathbf{\Omega}_{1/2}^K \mathbf{V} \mathbf{\Omega}_{1/2}^K .$$

$$2. \mathbf{SC} = \mathbf{CA} \text{ avec } \mathbf{C}'\mathbf{C} = \mathbf{I} .$$

$$3. \text{Pour } \ell = 1, \dots, L' : \mathbf{y}_\ell = \mathbf{X}_0 \mathbf{\Omega}_{1/2}^K \mathbf{c}_\ell .$$

$$4. \text{Pour } \ell = 1, \dots, L' : \mathbf{b}_\ell = \sqrt{\lambda_\ell} \mathbf{\Omega}_{1/2}^K \mathbf{c}_\ell .$$

### 3 ACP simple ou des covariances

Toutes les variables sont affectées de poids  $w_k$  égaux à 1 ( $M_j M_j'$ )<sub>2</sub> =  $\sum_{k \in K} (x^{jk} - x^{j'k})^2$  (métrique euclidienne)

élémentaire). La base  $(\delta_k)_{k \in K}$  est orthonormée ( $\mathbf{Q} = \mathbf{I}_K$ ).

Variance du nuage = somme des variances des variables initiales  
 $(\text{Var } M_J = \sum_{k \in K} \text{Var } x^{Jk})$

Equation aux directions et valeurs propres :  $\sum_{k' \in K} v^{kk'} a_{k'\ell} = \lambda_\ell a_\ell^k$

Dans l'ACP simple, directions principales et valeurs propres s'obtiennent par diagonalisation de la matrice  $\mathbf{V}$  des variances et covariances entre les variables initiales : l'ACP simple est l'analyse des covariances.

### 3.1 Formulaire de l'ACP simple

Equation aux directions et valeurs propres :

$$\sum_{k' \in K} v_{kk'} a_{k'\ell} = \lambda_{\ell} a_{k\ell} \quad \text{avec} \quad a_{k\ell}^2 = a_{k\ell} \quad \text{et} \quad \sum_{k \in K} (a_{k\ell}^2) = 1$$

The diagram shows a square matrix with a diagonal line from the top-left to the bottom-right. The diagonal elements are labeled with the Greek letter lambda (λ). This matrix is equal to the product of two square matrices, both labeled with the letter A.

The diagram shows a square matrix with a diagonal line from the top-left to the bottom-right. The diagonal elements are labeled with the letter I. This matrix is equal to the product of two square matrices, both labeled with the letter A'.

Coordonnées principales

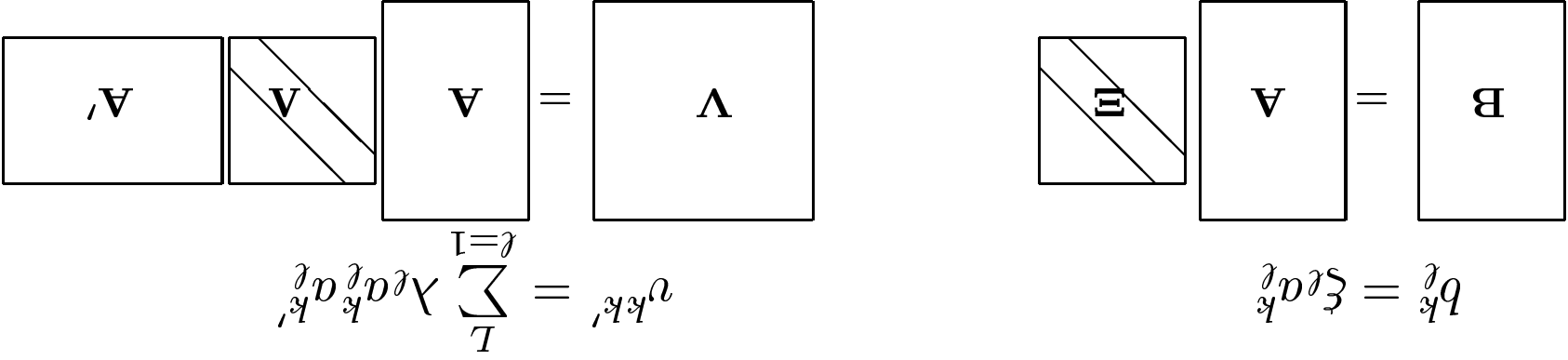
The diagram shows a square matrix labeled X<sub>0</sub> equal to the product of a square matrix labeled A and a square matrix labeled Y.

$$y_j^{\ell} = \sum_{k \in K} a_{k\ell} x_{jk}^0$$

Formules de reconstitution

The diagram shows a square matrix labeled X<sub>0</sub> equal to the product of a square matrix labeled Y' and a square matrix labeled A.

$$x_{jk}^0 = \sum_{\ell=1}^L a_{k\ell}^{\ell} y_j^{\ell}$$





## 3.2 Procédure de calcul

1. On calcule la matrice  $\mathbf{V}$  des covariances.
2. On diagonalise  $\mathbf{V}$ , d'où  $(\lambda_\ell, (a_{k\ell})_{k \in K})_{\ell=1, \dots, L}$  avec  $\sum_{k \in K} (a_{k\ell})^2 = 1$ .
3. Pour  $\ell = 1, 2, \dots, L'$  ( $L' \leq L$ ), on calcule  $\forall j \in J : y_\ell^j = \sum_{k \in K} a_{k\ell} x_0^{jk}$ .
4. Pour  $\ell = 1, 2, \dots, L'$  ( $L' \leq L$ ), on calcule  $\forall k \in K : b_\ell^k = \sqrt{\lambda_\ell} a_{k\ell}$ .

## 4 ACP standard ou des corrélations

Variables hétérogènes : ramener à une même échelle par une procédure de solidarisation

L'ACP standard revient donc d'une part à *solidariser par réduction* les échelles des variables, d'autre part à donner à ces variables des *poids égaux à 1* ; les variables ont donc la même contribution à la variance du nuage.

Dans l'ACP standard, directions principales et valeurs propres s'obtiennent par diagonalisation de la matrice des corrélations entre les variables initiales : l'ACP standard est l'*analyse des corrélations*.

## 4.1 Propriétés de l'ACP standard

- La variance du nuage = somme des valeurs propres = nombre de variables :

- Les coefficients de régression  $b_k^\ell = \xi^\ell a_k^\ell$  sont les coefficients de corrélation entre variables initiales et variables principales :  $b_k^\ell = r_k^\ell$ .

Les coefficients  $a_k^\ell$  vérifient donc la double propriété :

$$\forall \ell : \sum_{k \in K} (a_k^\ell)^2 = 1 \text{ et } \forall k : \sum_{\ell=1}^{\ell} \lambda_\ell (a_k^\ell)^2 = 1 (= \sum_T (b_k^\ell)^2)$$

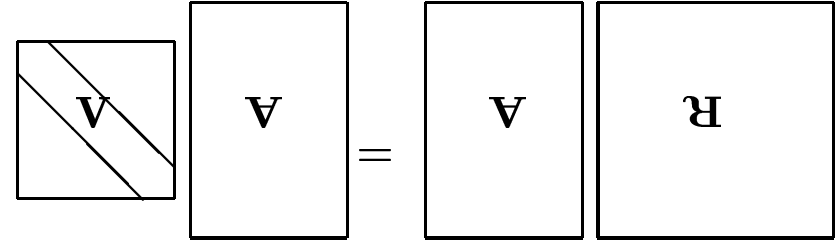
On a  $(a_k^\ell)^2 = \text{Ctr}_k^\ell$  (contribution relative de la variable  $k$  à l'axe  $\ell$ ) ; et  $(b_k^\ell)^2 = \lambda_\ell (a_k^\ell)^2 = \text{Ctr}_k^\ell$  (contribution relative de l'axe  $\ell$  à la variable  $k$ ).

- Le coefficient de corrélation multiple  $R_{L'}$  entre la variable initiale  $x_{J^k}$  et les  $L'$  premières variables principales est donné par

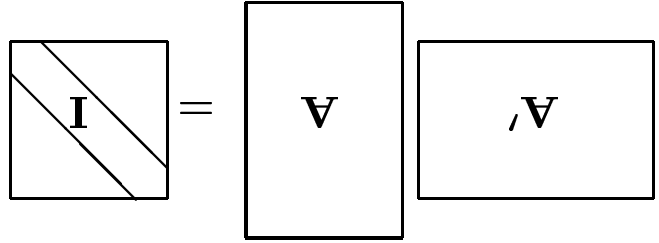
$$R_{L'}^2 = \sum_{L'}^{\ell=1} (r_k^\ell)^2 = \sum_{L'}^{\ell=1} (b_k^\ell)^2.$$

## 4.2 Formulaire de l'ACP standard

Equation aux directions et valeurs propres :

$$\sum_{k' \in K} r^{kk'} a_{k'\ell} = \lambda_\ell a_{k\ell}$$


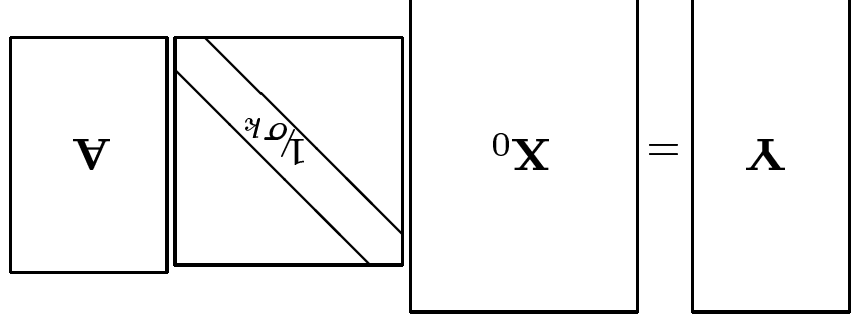
avec

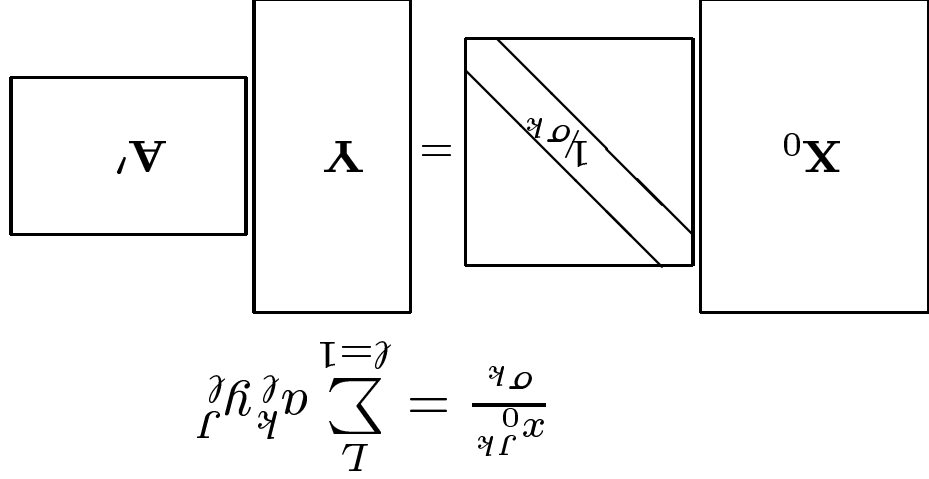
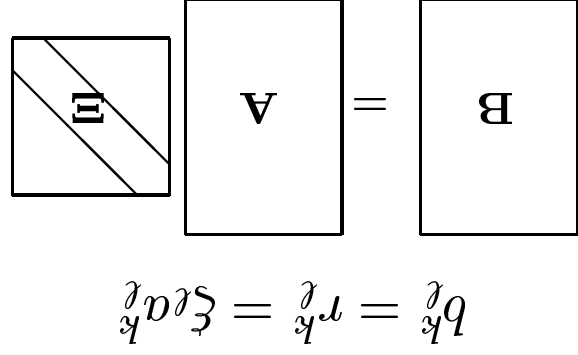
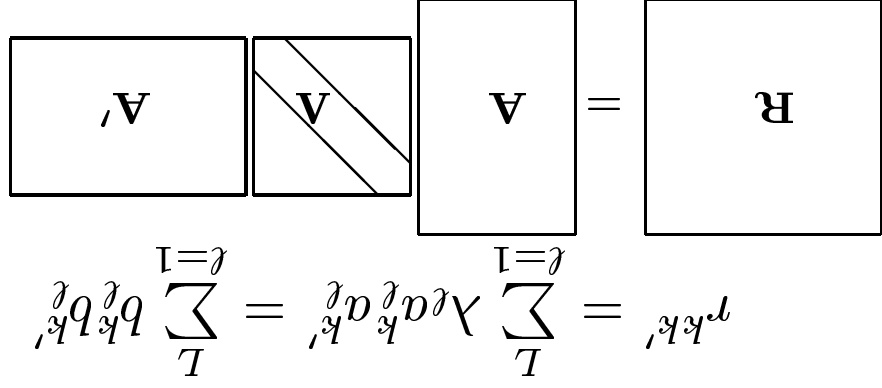
$$\sum_{k \in K} (a_{k\ell})^2 = \sum_{k \in K} (a_{k\ell})^2 = 1$$


Formules de reconstruction

Coordonnées principales

$$y_\ell^j = \sum_{k \in K} a_{k\ell} x_0^j$$





### 4.3 Procédure de calcul

1. On calcule la matrice  $\mathbf{R}$  des corrélations.
2. On diagonalise  $\mathbf{R}$ , d'où  $(\lambda_\ell, (a_{k\ell})_{k \in K})_{\ell=1, \dots, L}$  avec  $\sum_{k \in K} (a_{k\ell})^2 = 1$ .
3. Pour  $\ell = 1, \dots, L'$  ( $L' \leq L$ ) et  $\forall j \in J$ , on calcule :
 
$$y_\ell^j = \sum_{k \in K} a_{k\ell} x_0^{jk} / \sigma_k.$$
4. Pour  $\ell = 1, \dots, L'$  ( $L' \leq L$ ) et  $\forall k \in K$ , on calcule :  $b_\ell^k = \sqrt{\lambda_\ell} a_{k\ell}$ .

## 5 Méthodologie et interprétation

### 5.1 Analyse globale du nuage

- Contributions des individus, contributions des variables.
- Contributions des axes

On retient au moins tous les axes dont la contribution est supérieure à une contribution moyenne (e.g. à  $\text{Var}M^j/k$ , ou  $\text{Var}M^j/L$ ) ; on s'assure ensuite que les  $L^j$  axes retenus prennent en compte une part importante de la variance du nuage.

- *Facteur de taille* : sur le cercle des corrélations, les variables sont toutes situées d'un même côté.

Corrélation de la variable moyenne avec la variable principale importante = facteur de taille