

# Analyse Factorielle des Correspondances (AFC)

Abdelhakim Necir

*Département de Mathématiques  
Université de Biskra*

Master 1, 2021-2022

# Exemple

Le tableau suivant représente la couleur des cheveux et la couleur des yeux dans un échantillon de **370** individus.

	$V_2$			
	↓			
$V_1$	1: Brun	2: Châtain	3: Roux	4: Blond
↓				
1: Marron	68	119	26	7
2: Noisette	15	54	14	10
3: Vert	4	29	14	10

- $V_1 := (\text{marron}, \text{noisette}, \text{vert}) \rightarrow$  variables catégorielle (qualitative)
- $V_2 := (\text{brun}, \text{chatain}, \text{roux}, \text{blond}) \rightarrow$  variables catégorielle (qualitative)  
Marron  $\rightarrow$  une modalité de  $V_1$  et Blond  $\rightarrow$  une modalité de  $V_2$ , ...

Effectif de  $(1, 1) = 68$  et effectif de  $(2, 3) = 14$ , ...

# Tableau de contingence: variables qualitatives, modalités

Notons:

- $V_1 \rightarrow$  variable qualitative à  $p$  modalités
- $V_2 \rightarrow$  variable qualitative à  $q$  modalités

**Le tableau de contingence**  $N^*$  obtenu en croisant les deux variables  $V_1$  et  $V_2$ . Plus précisément, on a :

$$N^* = \begin{pmatrix} x_{11} & \cdots & x_{1q} \\ \vdots & \ddots & \vdots \\ x_{p1} & \cdots & x_{pq} \end{pmatrix},$$

où

$x_{ij} \rightarrow$  le nombre d'observations (ou effectif) pour lesquelles

$$V_1 = i \text{ et } V_2 = j.$$

# Tableau de contingence: variables qualitatives, modalités

Dans notre exemple (couleur des yeux vs couleur des cheveux), on a

$$N^* = \begin{pmatrix} 68 & 119 & 26 & 7 \\ 15 & 54 & 14 & 10 \\ 4 & 29 & 14 & 10 \end{pmatrix}.$$

# Tableaux des fréquences

On définit l'effectif total par

$$n := \sum_{i=1}^p \sum_{j=1}^q x_{ij}.$$

La fréquence observée du croisement des deux modalités  $(i, j)$ , est définie par

$$f_{ij} = \frac{x_{ij}}{n} = \mathbf{P}(V_1 = i, V_2 = j).$$

Ainsi, on définit le tableau des **fréquences observées** par

$$N := \begin{pmatrix} f_{11} & \cdots & f_{1q} \\ \vdots & \ddots & \vdots \\ f_{p1} & \cdots & f_{pq} \end{pmatrix} = \frac{1}{n} N^*.$$

# Tableaux des fréquences

On définit, respectivement, les effectifs marginaux des lignes et des colonnes, par

$$x_{i.} := \sum_{j=1}^q x_{ij} \text{ et } x_{.j} := \sum_{i=1}^p x_{ij}.$$

De même, on définit, respectivement, les fréquences marginales des lignes et des colonnes, par

$$f_{i.} := \sum_{j=1}^q f_{ij} = \mathbf{P}(V_1 = i) \text{ et } f_{.j} := \sum_{i=1}^p f_{ij} = \mathbf{P}(V_2 = j).$$

Il est important de souligner que

$$\sum_{i=1}^p \sum_{j=1}^q f_{ij} = \sum_{i=1}^p f_{i.} = \sum_{j=1}^q f_{.j} = 1.$$

On définit les fréquences conditionnelles associées aux profils-lignes par

$$f_{j/i} := \mathbf{P}(V_2 = j \mid V_1 = i) = \frac{\mathbf{P}(V_2 = j, V_1 = i)}{\mathbf{P}(V_1 = i)} = \frac{f_{ij}}{f_{i.}}.$$

En outre, on définit les fréquences conditionnelles associées aux profils-colonnes, par

$$f_{i/j} := \mathbf{P}(V_1 = i \mid V_2 = j) = \frac{\mathbf{P}(V_1 = i, V_2 = j)}{\mathbf{P}(V_2 = j)} = \frac{f_{ij}}{f_{.j}}.$$

On note aussi que

$$\sum_{j=1}^q f_{j/i} = 1, \text{ pour chaque } i = 1, \dots, p,$$

et

$$\sum_{i=1}^p f_{i/j} = 1, \text{ pour chaque } j = 1, \dots, q.$$



# Test d'indépendance du Khi-deux

- A priori nous n'avons aucune information sur la dépendance entre les deux variables  $V_1$  et  $V_2$ .
- Il se peut que ces deux dernières sont indépendantes et par conséquent l'AFC est inutile.

On peut traduire ceci, du fait que sous l'hypothèse d'indépendance,  $H_0$ , on a

$$\mathbf{P}(V_1 = i, V_2 = j) = \mathbf{P}(V_1 = i) \mathbf{P}(V_2 = j).$$

Contrairement à l'hypothèse de dépendance  $H_1$ , on a

$$\mathbf{P}(V_1 = i, V_2 = j) \neq \mathbf{P}(V_1 = i) \mathbf{P}(V_2 = j).$$

En d'autres termes, on a à tester l'hypothèse nulle

$$H_0 : f_{ij} = f_{i.} f_{.j}$$

contre l'hypothèse alternative

$$H_1 : f_{ij} \neq f_{i.} f_{.j}$$

# Test d'indépendance du Khi-deux

Donc l'hypothèse d'indépendance,  $H_0$  est une l'hypothèse théorique que nous devons la vérifier par une règle décision. Donc, la quantité

$$\tilde{f}_{ij} := f_{i.} f_{.j}$$

peut être vue comme **la fréquence théorique**. Ainsi nous définissons le tableau des fréquences théoriques par

$$\tilde{N} = \begin{pmatrix} f_{1.} f_{.1} & \cdots & f_{1.} f_{.q} \\ \vdots & \ddots & \vdots \\ f_{p.} f_{.1} & \cdots & f_{p.} f_{.q} \end{pmatrix}.$$

# Test d'indépendance du Khi-deux

Ce qui nous permet de définir **l'effectif théorique** par

$$e_{ij} := nf_{i.}f_{.j} = n \frac{x_{i.}}{n} \frac{x_{.j}}{n} = \frac{x_{i.}x_{.j}}{n}.$$

Ans le tableau les effectifs théoriques par

$$E := \begin{pmatrix} e_{11} & \cdots & e_{1q} \\ \vdots & \ddots & \vdots \\ e_{p1} & \cdots & e_{pq} \end{pmatrix} = n\tilde{N}.$$

# Test d'indépendance du Khi-deux

Notons  $\mathbf{f}_{ij}$  la variable aléatoire de la fréquence associée à la fréquence observée  $f_{ij}$ . Nous définissons **l'écart à l'indépendance** par

$$\begin{aligned}\Phi^2 &= \sum_{i=1}^p \sum_{j=1}^q \frac{(\mathbf{f}_{ij} - \mathbf{f}_{i.} \mathbf{f}_{.j})^2}{\mathbf{f}_{i.} \mathbf{f}_{.j}} \\ &= \sum_{i=1}^p \sum_{j=1}^q \frac{(\text{proba. conjointe} - \text{produit proba. marginales})^2}{\text{produit proba. marginales}}.\end{aligned}$$

La statistique du test associée est la **statistique du khi-deux** définie par

$$\chi^2 = n\Phi^2 = n \sum_{i=1}^p \sum_{j=1}^q \frac{(\mathbf{f}_{ij} - \mathbf{f}_{i.} \mathbf{f}_{.j})^2}{\mathbf{f}_{i.} \mathbf{f}_{.j}}.$$

# Test d'indépendance du Khi-deux

Nous avons, quand  $n \rightarrow \infty$

$\chi^2 \xrightarrow{D}$  la loi de Khi-deux à  $r$  degrés de liberté,

où

$$r := (p - 1) (q - 1) .$$

En pratique, nous devons s'assurer que

$n > 30$  et  $e_{ij} \geq 5$ , pour tout les  $(i, j)$ .

# Test d'indépendance du Khi-deux

Une fois que les fréquences  $f_{ij}$  sont observées, elles seront notées par  $f_{ij}$ , ainsi la statistique du Khi-deux  $\chi^2$  sera aussi observée, et on écrit

$$\chi_{obs}^2 = n \sum_{i=1}^p \sum_{j=1}^q \frac{(f_{ij} - f_{i.} f_{.j})^2}{f_{i.} f_{.j}}.$$

# Test d'indépendance du Khi-deux

Rappelons que, dans notre exemple, la matrice des données est

$$N^* = \begin{pmatrix} 68 & 119 & 26 & 7 \\ 15 & 54 & 14 & 10 \\ 4 & 29 & 14 & 10 \end{pmatrix}$$

Pour tester l'indépendance entre les deux variables  $V_1$  et  $V_2$ , il suffit d'utiliser entre autres, la fonction "chisq.test" du logiciel R. Voici les commandes à utiliser:

```
tab<-matrix(c(68,119,26,7,15,54,14,10,4,29,14,10),
             ,ncol=4,byrow=TRUE)

test=chisq.test(tab)

X-squared =34.114, df=6,p-value=6.394e-06
```

# Test d'indépendance du Khi-deux

Comme

$$p - value < \alpha := 0.05,$$

alors on rejette l'hypothèse, nulle, d'association (ou d'indépendance) des deux variables  $V_1$  et  $V_2$ .

En conclusion, il y a une liaison entre la couleur des cheveux et la couleur des yeux. L'AFC donc á un sens.



# Écart à l'indépendance

Rappelons que **l'écart à l'indépendance observé** est défini par

$$\Phi_{obs}^2 = \sum_{i=1}^p \sum_{j=1}^q \frac{(f_{ij} - f_{i.} f_{.j})^2}{f_{i.} f_{.j}} = \frac{\chi_{obs}^2}{n}.$$

Soient les événements suivants

$$A_i := \{V_1 = i\}, \quad B_j := \{V_2 = j\}$$

et

$$A_i \cap B_j := \{V_1 = i, V_2 = j\}.$$

Nous avons

$$\mathbf{P}(A_i \cap B_j) = f_{ij}, \quad \mathbf{P}(A_i) = f_{i.}, \quad \mathbf{P}(B_j) = f_{.j}.$$

Ainsi, nous avons

$$\Phi_{obs}^2 = \sum_{i=1}^p \sum_{j=1}^q \frac{(\mathbf{P}(A_i \cap B_j) - \mathbf{P}(A_i) \mathbf{P}(B_j))^2}{\mathbf{P}(A_i) \mathbf{P}(B_j)}.$$

Rappelons que

$$\begin{aligned} A_i \text{ et } B_j \text{ sont indépendants} &\Leftrightarrow \mathbf{P}(A_i \cap B_j) = \mathbf{P}(A_i) \mathbf{P}(B_j) \\ &\Leftrightarrow \Phi_{obs}^2 = 0. \end{aligned}$$

# Écart à l'indépendance

Remarquons que l'écart à l'indépendance observé peut être réécrit comme suit

$$\Phi_{obs}^2 = \sum_{i=1}^p \sum_{j=1}^q f_{i.} \frac{\left( \frac{f_{ij}}{f_{i.}} - f_{.j} \right)^2}{f_{.j}} = \sum_{i=1}^p \sum_{j=1}^q f_{i.} \frac{\underbrace{(f_{j/i} - f_{.j})^2}}{f_{.j}}.$$

Autrement dit

$$\Phi_{obs}^2 = \sum_{i=1}^p \sum_{j=1}^q \mathbf{P}(A_i) \frac{\underbrace{(\mathbf{P}(B_j/A_i) - \mathbf{P}(B_j))^2}}{\mathbf{P}(B_j)}.$$

Nous avons

$$\begin{aligned} A_i \text{ et } B_j \text{ sont indépendants} &\Leftrightarrow \mathbf{P}(B_j/A_i) = \mathbf{P}(B_j) \\ &\Leftrightarrow \Phi_{obs}^2 = 0. \end{aligned}$$

# Écart à l'indépendance

Par symétrie nous avons aussi

$$\Phi_{obs}^2 = \sum_{i=1}^p \sum_{j=1}^q f_{.j} \frac{\left( \frac{f_{ij}}{f_{.j}} - f_{i.} \right)^2}{f_{i.}} = \sum_{i=1}^p \sum_{j=1}^q f_{.j} \frac{\underbrace{(f_{i/j} - f_{i.})^2}}{f_{i.}}.$$

Autrement dit

$$\Phi_{obs}^2 = \sum_{i=1}^p \sum_{j=1}^q \mathbf{P}(B_j) \frac{\underbrace{(\mathbf{P}(A_i/B_j) - \mathbf{P}(A_i))^2}}{\mathbf{P}(A_i)}.$$

Nous avons

$$\begin{aligned} A_i \text{ et } B_j \text{ sont indépendants} &\Leftrightarrow \mathbf{P}(A_i/B_j) = \mathbf{P}(A_i) \\ &\Leftrightarrow \Phi_{obs}^2 = 0. \end{aligned}$$

# Écart à l'indépendance

En résumé, nous avons trois variantes de l'écart à l'indépendance:

$$\Phi_{obs}^2 = \sum_{i=1}^p \sum_{j=1}^q \frac{(f_{ij} - f_{i.} f_{.j})^2}{f_{i.} f_{.j}},$$

$$\Phi_{obs}^2 = \sum_{i=1}^p \sum_{j=1}^q f_{i.} \frac{\left( \frac{f_{ij}}{f_{i.}} - f_{.j} \right)^2}{f_{.j}} = \sum_{i=1}^p \sum_{j=1}^q f_{i.} \frac{(f_{j/i} - f_{.j})^2}{f_{.j}}$$

et

$$\Phi_{obs}^2 = \sum_{i=1}^p \sum_{j=1}^q f_{.j} \frac{\left( \frac{f_{ij}}{f_{.j}} - f_{i.} \right)^2}{f_{i.}} = \sum_{i=1}^p \sum_{j=1}^q f_{.j} \frac{(f_{i/j} - f_{i.})^2}{f_{i.}}$$

# Tableaux des profils-lignes et profils-colonnes

On note

$$D_r := \begin{pmatrix} f_{1.} & \cdots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \cdots & f_{p.} \end{pmatrix} \rightarrow D_r^{-1} := \begin{pmatrix} 1/f_{1.} & \cdots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \cdots & 1/f_{p.} \end{pmatrix},$$

$$N := \begin{pmatrix} f_{11} & \cdots & f_{1q} \\ \vdots & \ddots & \vdots \\ f_{p1} & \cdots & f_{pq} \end{pmatrix} \rightarrow D_r^{-1} N = \begin{pmatrix} \frac{f_{11}}{f_{1.}} & \cdots & \frac{f_{1q}}{f_{1.}} \\ \vdots & \ddots & \vdots \\ \frac{f_{p1}}{f_{p.}} & \cdots & \frac{f_{pq}}{f_{p.}} \end{pmatrix} =: X_r$$

$X_r \rightarrow$  le tableau des profils-lignes.

# Tableaux des profils-lignes et profils-colonnes

$$D_c := \begin{pmatrix} f_{.1} & \cdots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \cdots & f_{.q} \end{pmatrix} \rightarrow D_c^{-1} := \begin{pmatrix} 1/f_{.1} & \cdots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \cdots & 1/f_{.q} \end{pmatrix},$$

$$N^t := \begin{pmatrix} f_{11} & \cdots & f_{1q} \\ \vdots & \ddots & \vdots \\ f_{p1} & \cdots & f_{pq} \end{pmatrix} \rightarrow D_c^{-1} N^t = \begin{pmatrix} \frac{f_{11}}{f_{.1}} & \cdots & \frac{f_{1q}}{f_{.1}} \\ \vdots & \ddots & \vdots \\ \frac{f_{p1}}{f_{.p}} & \cdots & \frac{f_{pq}}{f_{.p}} \end{pmatrix} =: X_c$$

$X_c \rightarrow$  le tableau des profils-colonnes.

On définit, le centre de gravité de ce nuage de points des profils-lignes par

$$\begin{aligned} g_r &= (f_{.1}, \dots, f_{.q})^t \\ &= X_r^t D_r \mathbf{1}_p \\ &= (D_r^{-1} N)^t D_r \mathbf{1}_p = N^t \underbrace{D_r^{-1} D_r}_{\mathbf{I}_p} \mathbf{1}_p = N^t \mathbf{1}_p, \end{aligned}$$

où  $\mathbf{1}_p = (1, \dots, 1)^t$ , le vecteur unitaire de  $p \times 1$ .



De même, on définit, leur centre de gravité du nuage des profils-colonnes par

$$\begin{aligned} g_c &= (f_{1\cdot}, \dots, f_{p\cdot})^t = X_c^t D_c \mathbf{1}_q \\ &= (D_c^{-1} N^t)^t D_c \mathbf{1}_q = N \underbrace{D_c^{-1} D_c}_{\mathbf{1}_p} \mathbf{1}_p = N \mathbf{1}_q, \end{aligned}$$

où  $\mathbf{1}_q := (1, \dots, 1)^t$ , le vecteur unitaire de  $q \times 1$ .

# Exemple

Considérons la matrice des données de l'exemple ci-dessus

$$N^* = \begin{pmatrix} 68 & 119 & 26 & 7 \\ 15 & 54 & 14 & 10 \\ 4 & 29 & 14 & 10 \end{pmatrix}$$

Nous avons ici  $p = 3$ ,  $q = 4$  et

$$n = 68 + 15 + 4 + 119 + 54 + 29 + 26 + 14 + 14 + 7 + 10 + 10 = 370.$$

Ainsi

$$N = \begin{pmatrix} 68/370 & 119/370 & 26/370 & 7/370 \\ 15/370 & 54/370 & 14/370 & 10/370 \\ 4/370 & 29/370 & 14/370 & 10/370 \end{pmatrix}.$$

# Exemple

Les fréquences marginales-lignes sont

$$f_{1.} = \sum_{j=1}^4 f_{1j} = 68/370 + 119/370 + 26/370 + 7/370 = 220/370$$

$$f_{2.} = \sum_{j=1}^4 f_{2j} = 15/370 + 54/370 + 14/370 + 10/370 = 93/370$$

$$f_{3.} = \sum_{j=1}^4 f_{3j} = 4/370 + 29/370 + 14/370 + 10/370 = 57/370.$$

# Exemple

Les fréquences marginales-colonnes sont

$$f_{.1} = \sum_{i=1}^3 f_{i1} = 68/370 + 15/370 + 4/370 = 87/370$$

$$f_{.2} = \sum_{i=1}^3 f_{i2} = 119/370 + 54/370 + 29/370 = 202/370$$

$$f_{.3} = \sum_{i=1}^3 f_{i3} = 26/370 + 14/370 + 14/370 = 27/185$$

$$f_{.4} = \sum_{i=1}^3 f_{i4} = 7/370 + 10/370 + 10/370 = 27/370.$$

# Exemple

Les centres de gravité des profils-linges et profils-colonnes, respectivement, sont

$$g_r = (87/370, 202/370, 27/185, 27/370)^t,$$

et

$$g_c = (220/370, 93/370, 57/370)^t.$$

# Exemple

Les matrices diagonales de profils-lignes et profils-colonnes sont respectivement

$$D_r = \begin{pmatrix} 220/370 & 0 & 0 \\ 0 & 93/370 & 0 \\ 0 & 0 & 57/370 \end{pmatrix}$$

et

$$D_c = \begin{pmatrix} 87/370 & 0 & 0 & 0 \\ 0 & 202/370 & 0 & 0 \\ 0 & 0 & 27/185 & 0 \\ 0 & 0 & 0 & 27/370 \end{pmatrix}$$

# Exemple

Les matrices profils-lignes et profils-colonnes, respectivement, sont

$$X_r = D_r^{-1}N = \begin{pmatrix} \frac{17}{55} & \frac{119}{220} & \frac{13}{110} & \frac{7}{220} \\ \frac{5}{31} & \frac{18}{31} & \frac{14}{93} & \frac{10}{93} \\ \frac{4}{57} & \frac{29}{57} & \frac{14}{57} & \frac{10}{57} \end{pmatrix}$$

et

$$X_c = D_c^{-1}N^t = \begin{pmatrix} \frac{68}{87} & \frac{5}{29} & \frac{4}{87} \\ \frac{119}{202} & \frac{27}{101} & \frac{29}{202} \\ \frac{13}{27} & \frac{7}{27} & \frac{7}{27} \\ \frac{27}{27} & \frac{27}{27} & \frac{27}{27} \end{pmatrix}.$$

Rappelons que la matrice des profils-lignes est

$$X_r := \begin{pmatrix} 1 \rightarrow & \frac{f_{11}}{f_{1.}} & \dots & \frac{f_{1q}}{f_{1.}} \\ & \dots & \ddots & \vdots \\ i \rightarrow & \frac{f_{i1}}{f_{i.}} & \dots & \frac{f_{iq}}{f_{i.}} \\ & \vdots & \ddots & \vdots \\ i' \rightarrow & \frac{f_{i'1}}{f_{i'.}} & \dots & \frac{f_{i'q}}{f_{i'.}} \\ & \vdots & \ddots & \vdots \\ p \rightarrow & \frac{f_{p1}}{f_{p.}} & \dots & \frac{f_{pq}}{f_{p.}} \end{pmatrix}.$$



La distance euclidienne entre deux profils-lignes  $i$  et  $i'$  est définie par

$$d^2(i, i') = \sum_{j=1}^q \left( \frac{f_{ij}}{f_{i.}} - \frac{f_{i'j}}{f_{i'.}} \right)^2 = \|i - i'\|^2.$$

La distance de khi-deux entre deux profils-lignes  $i$  et  $i'$  est

$$d_{\chi^2}^2(i, i') = \sum_{j=1}^q \frac{1}{f_{.j}} \left( \frac{f_{ij}}{f_{i.}} - \frac{f_{i'j}}{f_{i'.}} \right)^2,$$

La pondération par  $1/f_{.j}$  à chaque carré de différence revient à donner des importances comparables aux diverses modalités  $j$  de la variable  $V_2$ . Sans cette pondération, la distance reflète surtout la différence entre les modalités de plus grandes fréquences.

# Métrique du khi-deux

Autrement dit

$$d_{\chi^2}^2(i, i') = \|i - i'\|_{M_r}^2,$$

où

$$M_r := D_c^{-1} = \begin{pmatrix} 1/f_{.1} & \cdots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \cdots & 1/f_{.q} \end{pmatrix}.$$

Celle-ci peut être vue comme la distance euclidienne, entre les deux profils-lignes  $i$  et  $i'$ , pondérée. Plus précisément

$$d_{\chi^2}^2(i, i') = (i - i')^t M_r (i - i').$$

La distance entre un profil-ligne et son centre de gravité  $g_r$  est définie par

$$d_{\chi^2}^2(i, g_r) = \sum_{j=1}^q \frac{1}{f_{.j}} \left( \frac{f_{ij}}{f_{.j}} - f_{.j} \right)^2 = \|i - g_r\|_{M_r}^2.$$

La distance entre deux profils-colonnes  $j$  et  $j'$  est

$$d_{\chi^2}^2(j, j') = \sum_{i=1}^p \frac{1}{f_{i.}} \left( \frac{f_{ij}}{f_{.j}} - \frac{f_{ij'}}{f_{.j'}} \right)^2 = \|j - j'\|_{M_c}^2,$$

où

$$M_c := D_r^{-1} = \begin{pmatrix} 1/f_{1.} & \cdots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \cdots & 1/f_{p.} \end{pmatrix}.$$

La distance entre un profil-colonne et son centre gravité  $g_c$  est définie par

$$d_{\chi^2}^2(j, g_c) = \sum_{i=1}^p \frac{1}{f_{i.}} \left( \frac{f_{ij}}{f_{.j}} - f_{i.} \right)^2 = \|j - g_c\|_{M_c}^2.$$

# Métrique du khi-deux

La matrice des profils-lignes est

$$X_r = \begin{pmatrix} 68/220 & 119/220 & 26/220 & 7/220 \\ 15/93 & 54/93 & 14/93 & 10/93 \\ 4/57 & 29/57 & 14/57 & 10/57 \end{pmatrix},$$

et les fréquences marginales des profils-colonnes sont

$$f_{.1} = 87/370, \quad f_{.2} = 202/370, \quad f_{.3} = 27/185, \quad f_{.4} = 27/370.$$

La distance euclidienne entre la première et la deuxième lignes est

$$\begin{aligned} d^2(1,2) &= \left( \frac{68}{220} - \frac{15}{93} \right)^2 + \left( \frac{119}{220} - \frac{54}{93} \right)^2 \\ &\quad + \left( \frac{26}{220} - \frac{14}{93} \right)^2 + \left( \frac{7}{220} - \frac{10}{93} \right)^2 \\ &= 3.0203 \times 10^{-2}. \end{aligned}$$

La distance de chi-deux entre la première et la deuxième lignes est

$$\begin{aligned}d_{\chi^2}^2(1, 2) &= \frac{370}{87} \left( \frac{68}{220} - \frac{15}{93} \right)^2 + \frac{370}{202} \left( \frac{119}{220} - \frac{54}{93} \right)^2 \\&\quad + \frac{370}{27} \left( \frac{26}{220} - \frac{14}{93} \right)^2 + \frac{370}{27} \left( \frac{7}{220} - \frac{10}{93} \right)^2 \\&= 0.18869.\end{aligned}$$

Les inerties totales des nuages de points profils-lignes et profils-colonnes par rapport aux centres de gravité correspondants sont définies respectivement par

$$\text{Inertie } (X_r / g_r) := \sum_{i=1}^p f_i \cdot d_{\chi^2}^2 (i, g_r) ,$$

et

$$\text{Inertie } (X_c / g_c) := \sum_{j=1}^q f_j d_{\chi^2}^2 (j, g_c) .$$

Observons que

$$\begin{aligned}\text{Inertie}(X_r/g_r) &= \sum_{i=1}^p f_{i\cdot} \left\{ d_{\chi^2}^2(i, g_r) \right\} \\&= \sum_{i=1}^p f_{i\cdot} \left\{ \sum_{j=1}^q \left( \frac{1}{f_{\cdot j}} \left( \frac{f_{ij}}{f_{i\cdot}} - f_{\cdot j} \right)^2 \right) \right\} \\&= \sum_{i=1}^p \sum_{j=1}^q f_{i\cdot} \left( \frac{1}{f_{\cdot j}} \left( \frac{f_{ij}}{f_{i\cdot}} - f_{\cdot j} \right)^2 \right) \\&= \sum_{i=1}^p \sum_{j=1}^q \frac{f_{i\cdot}}{f_{\cdot j}} \left( \frac{f_{ij}}{f_{i\cdot}} - f_{\cdot j} \right)^2 \\&= \sum_{i=1}^p \sum_{j=1}^q \frac{(f_{ij} - f_{i\cdot} f_{\cdot j})^2}{f_{i\cdot} f_{\cdot j}} = \Phi^2 = \frac{\chi^2}{n}.\end{aligned}$$

En d'autres termes, étudier l'inertie de  $X_r$  revient à étudier l'écart à l'indépendance  $\Phi^2$ .

On montre aussi que

$$\text{Inertie}(X_r/g_r) = \text{Inertie}(X_c/g_c) = \Phi^2 = \frac{\chi^2}{n}.$$



Pour notre exemple on a

$$X_r = \begin{pmatrix} 68/220 & 119/220 & 26/220 & 7/220 \\ 15/93 & 54/93 & 14/93 & 10/93 \\ 4/57 & 29/57 & 14/57 & 10/57 \end{pmatrix}.$$

Donc

$$\text{Inertie}(X_r/g_r) = \sum_{i=1}^3 f_i \cdot d_{\chi^2}^2(i, g_r) = \frac{\chi^2}{370} = \frac{34.114}{370} = 0.0922 = 9.22\%$$

On en déduit que

$$\text{Inertie}(X_c/g_c) = 9.22\%.$$

Rappelons que la matrice des profils-lignes est définie par

$$X_r = D_r^{-1} N = \begin{pmatrix} \frac{f_{11}}{f_{1\cdot}} & \cdots & \frac{f_{1q}}{f_{1\cdot}} \\ \vdots & \ddots & \vdots \\ \frac{f_{p1}}{f_{p\cdot}} & \cdots & \frac{f_{pq}}{f_{p\cdot}} \end{pmatrix},$$

et son centre de gravité est

$$g_r = (f_{\cdot 1}, \dots, f_{\cdot q})^t.$$

On définit le nuage profils-lignes centré par

$$Y_r := X_r - \mathbf{1}_p g_r^t,$$

où

$$\mathbf{1}_p = (1, \dots, 1)^t,$$

le vecteur unitaire de  $p \times 1$ . En d'autres termes

$$Y_r = \begin{pmatrix} \frac{f_{11}}{f_{1.}} - f_{.1} & \dots & \frac{f_{1q}}{f_{1.}} - f_{.q} \\ \vdots & \ddots & \vdots \\ \frac{f_{p1}}{f_{p.}} - f_{.1} & \dots & \frac{f_{pq}}{f_{p.}} - f_{.q} \end{pmatrix}.$$

Les éléments de  $Y_r$  sont

$$y_{i,j} = \frac{f_{ij}}{f_{i.}} - f_{.j}, \quad i = 1, \dots, p \text{ et } j = 1, \dots, q.$$

On désigne par

$$\mathbf{y}_i := (y_{i,1}, \dots, y_{i,q})^t, \quad i = 1, \dots, p,$$

les vecteurs lignes de la matrice  $Y_r$ .

Rappelons tout d'abord que nous avons muni l'espace  $\mathbb{R}^q$  de la métrique de  $\chi^2$ . Plus précisément

$$d_{\chi^2}^2(\mathbf{y}_i, \mathbf{y}_{i'}) = \sum_{j=1}^q \frac{1}{f_{.j}} (y_{ij} - y_{i'j})^2 = \|\mathbf{y}_i - \mathbf{y}_{i'}\|_{M_r}^2 = \mathbf{y}_i^t M_r \mathbf{y}_{i'},$$

où

$$M_r := D_c^{-1} = \begin{pmatrix} 1/f_{.1} & \cdots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \cdots & 1/f_{.q} \end{pmatrix}.$$

# ACP des deux nuages de profils

Nous allons procéder l'ACP sur la matrice des données  $Y_r$ .

On note par  $E$  l'axe principal de l'ACP et  $u$  son vecteur directeur associé de norme 1 par rapport à métrique  $M_r$ , c'est à dire

$$\|u\|_{M_r}^2 = \langle u, u \rangle_{M_r} = u^t M_r u = 1.$$

On note

la projection du point  $\mathbf{y}_i$  sur  $E^\perp := \mathbf{Proj}_{E^\perp, i}$ .

Nous définissons l'inertie du nuage  $Y_r$  par rapport à  $E^\perp$  par

$$\text{Inertie} \left( Y_r / E^\perp \right) = \sum_{i=1}^p f_i. d_{\chi^2}^2 \left( \mathbf{y}_i, \mathbf{Proj}_{E^\perp, i} \right).$$

Rappelons que, d'après la relation de Chasles, on a

$$\frac{\langle \mathbf{y}_i, u \rangle_{M_r} u}{\|u\|_{M_r}^2} + \mathbf{Proj}_{E^\perp, i} = \mathbf{y}_i,$$

ce qui implique que

$$\mathbf{Proj}_{E^\perp, i} - \mathbf{y}_i = \frac{\langle \mathbf{y}_i, u \rangle_{M_r} u}{\|u\|_{M_r}^2}.$$

Alors

$$\begin{aligned}d_{\chi^2}^2(\mathbf{y}_i, \mathbf{Proj}_{E^\perp, i}) &= \|\mathbf{y}_i - \mathbf{Proj}_{E^\perp, i}\|_{M_r}^2 \\&= \left\| \frac{\langle \mathbf{y}_i, \mathbf{u} \rangle_{M_r} \mathbf{u}}{\|\mathbf{u}\|_{M_r}^2} \right\|_{M_r}^2 = \frac{\langle \mathbf{y}_i, \mathbf{u} \rangle_{M_r}^2 \|\mathbf{u}\|_{M_r}^2}{\|\mathbf{u}\|_{M_r}^4} \\&= \langle \mathbf{y}_i, \mathbf{u} \rangle_{M_r}^2 \\&= (\mathbf{y}_i^t M_r \mathbf{u})^2 = (\mathbf{y}_i^t M_r \mathbf{u}) (\mathbf{y}_i^t M_r \mathbf{u})^t \\&= (\mathbf{y}_i^t M_r \mathbf{u}) (\mathbf{u}^t M_r \mathbf{y}_i) \\&= (\mathbf{u}^t M_r \mathbf{y}_i) (\mathbf{y}_i^t M_r \mathbf{u}) = \mathbf{u}^t M_r \mathbf{y}_i \mathbf{y}_i^t M_r \mathbf{u}.\end{aligned}$$



Par conséquent

$$\text{Inertie} \left( Y_r / E^\perp \right) = u^t M_r \left[ \sum_{i=1}^p f_i \cdot \mathbf{y}_i \mathbf{y}_i^t \right] M_r u.$$

On note que

$$\sum_{i=1}^p f_i \cdot \mathbf{y}_i \mathbf{y}_i^t = Y_r^t D_r Y_r =: \mathbf{V}_r,$$

comme étant la matrice de variance-covariance associée à la matrice  $Y_r$  affectée aux poids  $D_r$ . Ainsi

$$\text{Inertie} \left( Y_r / E^\perp \right) = u^t M_r \mathbf{V}_r M_r u.$$

Nous allons maintenant chercher le vecteur  $u$  maximisant l'inertie  $(Y_r / E^\perp)$  sous la contrainte  $u^t M_r u = 1$ .

Ce que revient, en utilisant le multiplicateur de Lagrange, à maximiser la fonction

$$u \rightarrow \eta(u) := u^t M_r \mathbf{V}_r M_r u - \lambda (u^t M_r u - 1).$$

Il est clair que la dérivée de cette fonction est

$$\eta'(u) = 2M_r \mathbf{V}_r M_r u - 2\lambda M_r u.$$

En résolvant l'équation  $\eta'(u) = 0$ , on trouve

$$M_r \mathbf{V}_r M_r u = \lambda M_r u.$$

Rappelons que la matrice  $M_r$  est inversible, alors la dernière équation se réduit à

$$\mathbf{V}_r M_r u = \lambda u.$$

Comme nous l'avons fait au premier chapitre, nous allons appliquer l'ACP à la matrice  $\mathbf{V}_r M_r$ . En d'autres termes nous cherchons les valeurs propres et les vecteurs propres associés à  $\mathbf{V}_r M_r$ , définissant les axes principaux.

## Rappel:

$X_r \rightarrow$  matrice des profils-lignes.

$Y_r := X_r - \mathbf{1}_p g_r^t \rightarrow$  matrice des profils-lignes centrée.

$\mathbf{V}_r = Y_r^t D_r Y_r = X_r^t D_r X_r - g_r g_r^t \rightarrow$  matrice de Var-Covar.

$E \rightarrow$  l'axe principal et son vecteur directeur  $u$  avec  $\|u\|_{M_r}^2 = 1$ .

Inertie  $(Y_r / E^\perp) = u^t M_r \mathbf{V}_r M_r u \rightarrow$  inertie totale ( $M_r = D_c^{-1}$ ).

$\max_{u, \|u\|_{M_r}^2=1} u^t M_r \mathbf{V}_r M_r u \rightarrow \mathbf{V}_r M_r u = \lambda u \rightarrow$  vecteur propre de  $\mathbf{V}_r M_r$ .

**Remarque.** La matrice  $\mathbf{V}_r M_r$  est  $M_r$ -symétrique, c'est à dire  $M_r \mathbf{V}_r M_r$  est symétrique.

En effet, comme  $\mathbf{V}_r$  et  $M_r = D_c^{-1}$  sont symétriques, il est alors évident que  $M_r \mathbf{V}_r M_r$  l'est aussi. En effet,

$$\begin{aligned}(M_r \mathbf{V}_r M_r)^t &= ((M_r \mathbf{V}_r) M_r)^t = M_r^t (M_r \mathbf{V}_r)^t = M_r^t \mathbf{V}_r^t M_r^t \\ &= M_r \mathbf{V}_r M_r,\end{aligned}$$

car  $\mathbf{V}_r$  et  $M_r$  sont tout les deux symétriques.

**Théorème.** Le centre de gravité du nuage des profils-lignes,  $g_r$ , est  $M_r$ -orthogonal au nuage des profils-lignes centré  $Y_r$  (de même,  $g_c$  est  $M_c$ -orthogonal au nuage des profils-colonnes centré  $Y_c$ ).

**Preuve.** Tout d'abord observons que

$$M_r g_r = D_c^{-1} g_r = \begin{pmatrix} 1/f_{.1} & \cdots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \cdots & 1/f_{.q} \end{pmatrix} \begin{pmatrix} f_{.1} \\ \vdots \\ f_{.q} \end{pmatrix} = \begin{pmatrix} 1 \\ \vdots \\ 1 \end{pmatrix} = \mathbf{1}_q$$

et

$$\langle g_r, g_r \rangle_{M_r} = \|g_r\|_{M_r}^2 = g_r^t M_r g_r = g_r^t \mathbf{1}_q = f_{.1} + \dots + f_{.q} = 1.$$

Soit  $\mathbf{y}_i := r_i - g_r$  la  $i$ -ème ligne de  $Y_r$  avec

$$r_i := (f_{i1}/f_{i.}, \dots, f_{iq}/f_{i.})^t$$

la  $i$ -ème ligne de  $X_r$ . Alors

$$\begin{aligned}\langle \mathbf{y}_i, g_r \rangle_{M_r} &= \langle r_i - g_r, g_r \rangle_{M_r} = \langle r_i, g_r \rangle_{M_r} - \langle g_r, g_r \rangle_{M_r} \\ &= r_i^t M_r g_r - \langle g_r, g_r \rangle_{M_r} \\ &= r_i^t M_r g_r - \langle g_r, g_r \rangle_{M_r} \\ &= r_i^t M_r g_r - 1 = r_i^t \mathbf{1}_q - 1 \\ &= (f_{i1}/f_{i.}, \dots, f_{iq}/f_{i.}) \mathbf{1}_q - 1 \\ &= f_{i1}/f_{i.} + \dots + f_{iq}/f_{i.} - 1 = 1 - 1 = 0.\end{aligned}$$

# ACP des deux nuages de profils

**Corollaire.** Le centre de gravité  $g_r$  est un vecteur propre de  $V_r M_r$  associé à la valeur propre  $\lambda = 0$ .

**Preuve.** On va démontrer que  $V_r M_r g_r = 0g_r$ . Nous avons

$$V_r M_r g_r = V_r \mathbf{1}_q = \begin{pmatrix} \frac{f_{11}}{f_{1.}} - f_{.1} + \dots + \frac{f_{1q}}{f_{1.}} - f_{.q} \\ \vdots \\ \frac{f_{p1}}{f_{p.}} - f_{.1} + \dots + \frac{f_{pq}}{f_{p.}} - f_{.q} \end{pmatrix}.$$

Observons que

$$\begin{aligned} \frac{f_{11}}{f_{1.}} - f_{.1} + \dots + \frac{f_{1q}}{f_{1.}} - f_{.q} &= \left( \frac{f_{11}}{f_{1.}} + \dots + \frac{f_{1q}}{f_{1.}} \right) - (f_{.1} + \dots + f_{.q}) \\ &= 1 - 1 = 0. \end{aligned}$$

De la même façon on montre aussi que les autres lignes sont nuls. Donc  $V_r M_r g_r = 0g_r$ .



**Remarque.** On rappelle que le déterminant d'une matrice carré égal au produit des valeurs propres.

Comme  $\lambda = 0$  est une valeur propre de  $V_r M_r$  alors son déterminant est nul. Ce qui implique que cette matrice n'est pas inversible. Comme le rang égal au nombre maximum des valeurs propres non nuls, donc

$$\text{rg}(V_r M_r) \leq q - 1.$$

Avec le même raisonnement, on en conclus que

$$\text{rg}(V_c M_c) \leq p - 1.$$

**Théorème.** La matrice  $V_r M_r$  a les mêmes valeurs propres que

$$X_r^t D_r X_r M_r = N^t D_r^{-1} N D_c^{-1} = X_r^t X_c^t := A_r,$$

sauf  $g_r$  qui a une valeur propre  $\lambda = 1$ .

**Preuve.** Soit  $0 \neq v \neq g_r$  un vecteur propre de matrice  $V_r M_r$ , i.e.  $V_r M_r v = \lambda v$ . Ceci implique que

$$(X_r^t D_r X_r - g_r g_r^t) M_r v = \lambda v,$$

ainsi

$$X_r^t D_r X_r M_r v - g_r g_r^t M_r v = \lambda v.$$

Observons que

$$X_r^t D_r X_r M_r v - g_r \langle g_r, v \rangle_{M_r} = \lambda v.$$

D'un autre côté  $g_r$  est  $v$  sont deux vecteurs propres,  $M_r$ —orthogonaux, de  $V_r M_r$ , donc

$$X_r^t D_r X_r M_r v = \lambda v,$$

ce qui implique que  $v$  est un vecteur propre de  $X_r^t D_r X_r M_r$  associé à la même paramètre  $\lambda$ . Par ailleurs nous avons

$$\begin{aligned} X_r^t D_r X_r M_r &= (D_r^{-1} N)^t D_r D_r^{-1} N D_c^{-1} \\ &= N^t D_r^{-1} N D_c^{-1} = X_r^t X_c^t = A_r. \end{aligned}$$

**Corollaire.** Il n'est pas nécessaire de centrer le nuage de points des profils-lignes avant de réaliser l'ACP, et on peut travailler directement avec la matrice

$$A_r = X_r^t X_c^t.$$

Dans le cas des profils-colonnes, on s'intéressera à la matrice

$$A_c = X_c^t X_r^t.$$

# Lien entre les deux ACP

Rappel:

ACP profils-lignes  $\rightarrow V_r M_r$ .

ACP profils-colonnes  $\rightarrow V_c M_c$ .

val. pro  $V_r M_r \longleftrightarrow (\text{sauf, } g_r, \lambda = 0)$



val. pro. de  $X_r^t X_c^t \longleftrightarrow (\text{sauf, } g_r, \lambda = 1)$

De même

val. pro  $V_c M_c \longleftrightarrow (\text{sauf, } g_c, \lambda = 0)$



val. pro. de  $X_c^t X_r^t \longleftrightarrow (\text{sauf, } g_c, \lambda = 1)$

**Exemple:** Supposons que  $V_r M_r \in \mathcal{M}(4 \times 4)$  et que ces valeurs propres sont

$$\lambda = \{0.5, 0.3, 0.2, 0\}.$$

Alors les valeurs propres de  $A_r$  sont

$$\lambda = \{1, 0.5, 0.3, 0.2\}.$$

Rappel:  $M_r = D_c^{-1}$  et  $M_c = D_r^{-1}$ .

**Théorème.** Si  $u$ , est un vecteur propre, associé à la valeur propre  $\lambda \neq 0$ ,  $M_r$ —norme 1, de  $A_r$ , alors

$$\tilde{u} := \frac{1}{\sqrt{\lambda}} N D_c^{-1} u$$

est un vecteur propre,  $M_c$ —norme 1, pour  $A_c$ , pour la même valeur propre  $\lambda$ .

# Lien entre les deux ACP: formule de transition

Inversement, si  $\tilde{u}$  est vecteur propre, associé à la valeur propre  $\lambda \neq 0$ ,  $M_c$ —norme 1, de  $A_c$ , alors

$$u := \frac{1}{\sqrt{\lambda}} N^t D_r^{-1} \tilde{u},$$

est un vecteur propre,  $M_r$ —norme 1, pour  $A_r$ , pour la même valeur propre  $\lambda$ .



## Lien entre les deux ACP: formule de transition

**Preuve.** Notons  $u$  un vecteur propre  $A_r = N^t D_r^{-1} N D_c^{-1}$ , de norme 1 pour la métrique  $M_r$ , et associé à la valeur propre  $\lambda$ . C'est à dire  $A_r u = \lambda u$ . En d'autres termes

$$N^t D_r^{-1} N D_c^{-1} u = \lambda u.$$

En multipliant les deux membres de l'équation par  $N D_c^{-1}$ , on obtient

$$(N D_c^{-1}) N^t D_r^{-1} N D_c^{-1} u = \lambda (N D_c^{-1}) u,$$

ainsi

$$(N D_c^{-1} N^t D_r^{-1}) (N D_c^{-1} u) = \lambda (N D_c^{-1} u).$$

## Lien entre les deux ACP: formule de transition

Donc  $kND_c^{-1}u$  est un vecteur propre de  $ND_c^{-1}N^tD_r^{-1} = A_c$  associé à la même valeur propre  $\lambda$ .

On cherche ensuite la valeur de  $k$  telle que le vecteur propre soit, de norme 1 pour la métrique  $M_c$ .

# Lien entre les deux ACP: formule de transition

Autrement dit on cherche  $k$  telle que

$$\begin{aligned}(kND_c^{-1}u)^t M_c (kND_c^{-1}u) &= 1 \iff k^2 (u^t D_c^{-1} N^t) M_c N D_c^{-1} u = 1 \\ &\iff k^2 u^t D_c^{-1} (N^t D_r^{-1} N D_c^{-1}) u = 1 \\ &\iff k^2 u^t D_c^{-1} (X_r^t X_c^t) u = 1 \\ &\iff k^2 u^t D_c^{-1} (A_r u) = 1 \\ &\iff k^2 u^t D_c^{-1} \lambda u = 1\end{aligned}$$

$$\Longleftrightarrow \lambda k^2 (u^t M_r u) = 1$$

$$\Longleftrightarrow \lambda k^2 \|u\|_{M_r}^2 = 1$$

$$\Longleftrightarrow k = \frac{1}{\sqrt{\lambda}}.$$

En résumé, si  $u$  est un vecteur propre  $A_r = N^t D_r^{-1} N D_c^{-1}$ , pour la valeur propre  $\lambda$ , alors

$$\tilde{u} := \frac{1}{\sqrt{\lambda}} N D_c^{-1} u$$

est un vecteur propre pour  $A_c = N D_c^{-1} N^t D_r^{-1}$ , pour la même valeur propre  $\lambda$ .

Il existe un lien entre les vecteurs et valeurs propres des deux matrices

$$A_r := X_r^t X_c^t \text{ et } A_c := X_c^t X_r^t.$$

Nous avons

$$A_r \in M(q \times q) \text{ et } A_c \in M(p \times p).$$

Les deux matrices  $A_r$  et  $A_c$  ont les mêmes valeurs propres non nulles.

**Corollaire.** Les deux matrices  $A_r$  et  $A_c$  ont les mêmes valeurs propres non nulles. Par conséquent

$$\operatorname{rg}(A_r) = \operatorname{rg}(A_c),$$

ainsi

$$\tau := \operatorname{rg}(V_r M_r) = \operatorname{rg}(V_c M_c).$$

De plus

$$0 < \tau \leq \min(p-1, q-1).$$

# Lien entre les deux ACP

**Exemple:** Supposons que  $V_r M_r \in \mathcal{M}(4 \times 4)$  et  $V_c M_c \in \mathcal{M}(5 \times 5)$ , et que les valeurs propres de  $V_r M_r$  sont

$$\lambda = \{0.5, 0.3, 0.2, 0\} \rightarrow \text{rg}(V_r M_r) = 3$$

Alors les valeurs propres de  $A_r := X_r^t X_c^t$  sont

$$\lambda = \{1, 0.5, 0.3, 0.2\} \rightarrow \text{rg}(A_r) = 4$$

Ainsi les valeurs propres de  $A_c := X_c^t X_r^t$  sont

$$\lambda = \{1, 0.5, 0.3, 0.2, 0\} \rightarrow \text{rg}(A_c) = 4$$

Alors les valeurs propres de  $V_c M_c$  sont

$$\lambda = \{0.5, 0.3, 0.2, 0, 0\} \rightarrow \text{rg}(V_c M_c) = 3.$$

# Facteurs principaux et composantes principales

**Rappel:**  $u_i$  (sont  $M_r$ —orthonormés) et  $\tilde{u}_i$  (sont  $M_c$ —orthonormés) sont les vecteurs propres de  $A_r$  et  $A_c$  associés aux meme valeurs propres non-nulles,

**Définition.** On appel **facteurs principaux** des profils-lignes (resp. profils-colonnes), les vecteurs

$$w_i := M_r u_i,$$

resp.

$$\tilde{w}_i := M_c \tilde{u}_i.$$



# Facteurs principaux et composantes principales

**Proposition.** Les facteurs principaux des profils-lignes (resp. profils-colonnes) sont  $M_r^{-1}$ —orthonormés (resp.  $M_c^{-1}$ —orthonormé.

**Preuve.** Soient  $u_i, i = 1, \dots, p$  les axe principaux des profils-lignes. On sait que les  $u_i$  sont  $M_r$ —orthonormés, donc

$$w_i^t M_r^{-1} w_j = u_i^t M_r M_r^{-1} M_r u_j = u_i^t M_r u_j = 1 \text{ si } i = j \text{ et } 0 \text{ si } i \neq j.$$

**Proposition.** Les facteurs principaux des profils-lignes (rep. profils colonnes) sont les vecteurs propres  $M_r^{-1}$ —orthonormés (resp.  $M_c^{-1}$ —orthonormés) de la matrice de données  $M_r V_r$  (resp.  $M_c V_c$ ).

**Preuve..** Soit  $u$  un vecteur propre de  $V_r M_r$  associé à la valeur propre  $\lambda$ , alors  $V_r M_r u = \lambda u$ . En multipliant les deux membres de cette équation par  $M_r$ , on obtient

$$M_r (V_r M_r u) = \lambda M_r u.$$

Donc

$$M_r V_r (M_r u) = \lambda (M_r u),$$

où  $M_r u =: w$  est le facteur principal associé à  $u$ . De plus

$$w^t M_r^{-1} w = u^t M_r M_r^{-1} M_r u = u^t M_r u = 1,$$

car  $u$  est  $M_r$ -normé.

**Définition.** Les composantes principales des profils-lignes sont les  $M_r$ -coordonnées des vecteurs colonnes de la matrice de données  $Y_r := X_r - \mathbf{1}_p g_r^t$ . C'est à dire

$$c_k := Y_r w_k, \quad k = 1, \dots, \tau,$$

où  $w_k = M_r u_k$  sont les facteurs principaux.

Les composantes principales des profils-colonnes sont les  $M_c$ -coordonnées des vecteurs colonnes de la matrice de données  $Y_c := X_c - \mathbf{1}_q g_c^t$ . C'est à dire

$$\tilde{c}_k := Y_c \tilde{w}_k, \quad k = 1, \dots, \tau,$$

où  $\tilde{w}_k = M_c \tilde{u}_k$  sont les facteurs principaux.

## Proposition

**Proposition.** *Les composantes principales des profils-lignes sont les  $M_r$ -coordonnées des vecteurs colonnes de la matrice de données  $X_r$ . C'est à dire*

$$c_k := X_r w_k, \quad k = 1, \dots, \tau.$$

*Les composantes principales des profils-colonnes sont les  $M_c$ -coordonnées des vecteurs colonnes de la matrice de données  $X_c$ . C'est à dire*

$$\tilde{c}_k := X_c \tilde{w}_k, \quad k = 1, \dots, \tau.$$

Nous avons  $Y_r := X_r - \mathbf{1}_p g_r^t$  et

$$c_k = Y_r w_k = X_r w_k - \mathbf{1}_p g_r^t w_k = X_r w_k - \mathbf{1}_p g_r^t M_r u_k$$

Nous avons énoncé que les vecteurs propres de  $V_r M_r$  sont  $M_r$ -orthogonaux, donc  $g_r^t M_r u_k = 0$ .

**Preuve.** En effet, nous avons

$$c_k = Y_r w_k = (X_r - \mathbf{1}_p g_r^t) M_r u_k = X_r M_r u_k - \mathbf{1}_p (g_r^t M_r u_k) .$$

Rappelons que  $g_r$  et  $u_k$  sont deux vecteurs propres  $V_r M_r$  associés aux valeurs propres  $\lambda = 0$  et  $\lambda_k \neq 0$ . Nous avons aussi énoncé aussi que les vecteurs propres de  $V_r M_r$  sont  $M_r$ -orthogonaux. Donc  $g_r^t M_r u_k = 0$ , ainsi

$$c_k = X_r w_k .$$

**Proposition.** Nous avons

$$\frac{1}{p} \sum_{j=1}^p c_k(j) = \frac{1}{q} \sum_{j=1}^q \tilde{c}_k(j) = 0, \quad k = 1, \dots, \tau,$$

$$\frac{1}{p} \sum_{j=1}^p c_k^2(j) = \frac{1}{q} \sum_{j=1}^q \tilde{c}_k^2(j) = \lambda_k, \quad k = 1, \dots, \tau,$$

et

$$\begin{cases} \frac{1}{p} \sum_{i=1}^p c_k(i) c_\ell(i) = 0, \text{ pour } k \neq \ell \\ \frac{1}{q} \sum_{j=1}^q \tilde{c}_k(j) \tilde{c}_\ell(j) = 0, \text{ pour } k \neq \ell. \end{cases}$$

**Proposition.** Les facteurs principaux de l'ACP des profils-colonnes, associés aux valeurs propres non nulles, sont, à une constante près, les composantes principales de l'ACP des profils-lignes, et vice-versa. Plus précisément

$$c = \frac{1}{\sqrt{\lambda}} A_c^t \tilde{w} \text{ et } \tilde{c} = \frac{1}{\sqrt{\lambda}} A_r^t w.$$

Ce qui equivalent à

$$c = \sqrt{\lambda} \tilde{w} \text{ et } \tilde{c} = \sqrt{\lambda} w.$$



# Facteurs principaux et composantes principales

**Preuve.** En effet, soit  $u$  un axe principal, des profils-lignes, associé à la valeur propre  $\lambda \neq 0$  et  $\tilde{u}$  l'axe principal correspondant des profils-colonnes. Observons maintenant que, la composante principale associée à  $u$  est

$$\begin{aligned}c &= X_r w = X_r M_r u = (D_r^{-1} N) D_c^{-1} u = \frac{1}{\sqrt{\lambda}} D_r^{-1} N D_c^{-1} N^t D_r^{-1} \tilde{u} \\&= \frac{1}{\sqrt{\lambda}} X_r X_c \tilde{w} = \frac{1}{\sqrt{\lambda}} A_c^t \tilde{w}.\end{aligned}$$

Rappelons que  $\tilde{u} := N D_c^{-1} u / \sqrt{\lambda}$  est un axe principale des profils-colonnes. Donc

$$c = \sqrt{\lambda} D_r^{-1} \tilde{u} = \sqrt{\lambda} M_c \tilde{u} = \sqrt{\lambda} \tilde{w}.$$

Inversement

$$\tilde{c} = X_c M_c \tilde{u} = D_c^{-1} (N^t D_r^{-1} \tilde{u}) = D_c^{-1} (\sqrt{\lambda} u) = \sqrt{\lambda} w.$$

# Facteurs principaux et composantes principales:

Le résultat précédent conduit aux relations fondamentales de l'AFC reliant les composantes principales entre elles, dites les relations *quasi-barycentriques*:

**Proposition.** Soit  $\lambda_1 > \lambda_2 > \dots > \lambda_\tau \neq 0$ . Alors, pour tout  $k \leq \tau$ , on a

$$c_k := \frac{1}{\sqrt{\lambda_k}} X_r \tilde{c}_k \text{ et } \tilde{c}_k := \frac{1}{\sqrt{\lambda_k}} X_c c_k.$$

Par conséquent

$$c_k := \frac{1}{\sqrt{\lambda_k}} X_r X_c c_k$$

**Preuve.** On a

$$\begin{aligned}c_k &= X_r w_k = X_r M_r u_k = X_r D_c^{-1} \left( \frac{1}{\sqrt{\lambda}} N^t D_r^{-1} \tilde{u} \right) \\&= \frac{1}{\sqrt{\lambda}} X_r (D_c^{-1} N^t) (D_r^{-1} \tilde{u}) = \frac{1}{\sqrt{\lambda}} X_r (X_c) (M_c \tilde{u}) \\&= \frac{1}{\sqrt{\lambda}} X_r (X_c \tilde{w}) = \frac{1}{\sqrt{\lambda}} X_r \tilde{c}.\end{aligned}$$

# Facteurs principaux et composantes principales

**Proposition.** On a

$$0 \leq \lambda_k \leq 1, \quad k = 1, \dots, \tau.$$

**Preuve.** Soient  $\mathbf{A} = (a_{ij}) \in \mathcal{M}(p \times q)$ ,  $\mathbf{x} = (x_1, \dots, x_q)^t \in \mathbb{R}^p$  et  $\mathbf{y} = (y_1, \dots, y_p)^t \in \mathbb{R}^q$ . On munit  $\mathbb{R}^q$  et  $\mathbb{R}^p$  des normes  $L_1$  définis par

$$\|\mathbf{x}\| = \sum_{i=1}^q |x_i|, \quad \|\mathbf{y}\| = \sum_{i=1}^p |y_i|.$$

On munit aussi l'espace des matrices  $\mathcal{M}(p \times q)$  de la norme

$$\|\mathbf{A}\| = \sup_{\|\mathbf{x}\| \neq 0} \frac{\|\mathbf{Ax}\|}{\|\mathbf{x}\|} = \sup_{\|\mathbf{x}\| \leq 1} \|\mathbf{Ax}\| = \sup_{\|\mathbf{x}\|=1} \|\mathbf{Ax}\|.$$

# Facteurs principaux et composantes principales

On montre que

$$\|\mathbf{A}\| = \max_{j=1,\dots,q} \sum_{i=1}^p |a_{ij}| ,$$

(voir mon site moodle). Il est clair que

$$\frac{\|\mathbf{Ax}\|}{\|\mathbf{x}\|} \leq \|\mathbf{A}\| \implies \|\mathbf{Ax}\| \leq \|\mathbf{A}\| \|\mathbf{x}\| ,$$

ainsi

$$\|\mathbf{Ax}\| \leq \max_{j=1,\dots,q} \sum_{i=1}^p |a_{ij}| \|\mathbf{x}\| .$$

**Application:**  $\mathbf{A} = X_r \in \mathcal{M}(p \times q)$ . Comme  $c_k$  est un vecteur de  $\mathbb{R}^p$  alors  $\mathbf{x} = X_c c_k$  étant un vecteur de  $\mathbb{R}^p$ . Donc on a

$$\begin{aligned} \lambda \|c_k\| &= \|(X_r)(X_c c_k)\| = \|\mathbf{A}\mathbf{x}\| \leq \|\mathbf{A}\| \|\mathbf{x}\| \\ &\leq \max_{j=1,\dots,q} \sum_{i=1}^p |a_{ij}| \|\mathbf{x}\|. \end{aligned}$$

C'est à dire

$$\lambda \|c_k\| \leq \max_{j=1,\dots,q} \sum_{i=1}^p |a_{ij}| \|X_c c_k\|.$$

# Facteurs principaux et composantes principales

En appliquant ces dernières inégalités  $\mathbf{A} = X_c = (b_{ij}) \in \mathcal{M}(p \times q)$ , on écrit

$$\|X_c c_k\| \leq \|X_c\| \|c_k\| = \left\{ \max_{j=1, \dots, p} \sum_{i=1}^q |b_{ij}| \right\} \|c_k\|.$$

Ce qui implique

$$\lambda \|c_k\| = \|X_r X_c c_k\| \leq \left\{ \max_{j=1, \dots, q} \sum_{i=1}^p |a_{ij}| \right\} \left\{ \max_{j=1, \dots, p} \sum_{i=1}^q |b_{ij}| \right\} \|c_k\|.$$

En simplifiant par  $\|c_k\|$  (qui est évidemment non nulle), on obtien

$$0 < \lambda \leq \left\{ \max_{j=1, \dots, q} \sum_{i=1}^p |a_{ij}| \right\} \left\{ \max_{i=1, \dots, p} \sum_{j=1}^q |b_{ij}| \right\}.$$

# Facteurs principaux et composantes principales

$$0 < \lambda \leq \left\{ \max_{i=1,\dots,p} \sum_{j=1}^q \frac{f_{ij}}{f_{i.}} \right\} \left\{ \max_{j=1,\dots,q} \sum_{i=1}^p \frac{f_{ji}}{f_{.j}} \right\}.$$

Notons que

$$\sum_{j=1}^q \frac{f_{ij}}{f_{i.}} = \frac{1}{f_{i.}} \sum_{j=1}^q f_{ij} = \frac{f_{i.}}{f_{i.}} = 1$$

et

$$\sum_{i=1}^p \frac{f_{ij}}{f_{.j}} = \frac{1}{f_{.j}} \sum_{i=1}^p f_{ij} = \frac{f_{.j}}{f_{.j}} = 1,$$

et parconséquent  $0 < \lambda \leq 1$ .



**Théorème.** Nous avons

$$\frac{f_{ij}}{f_{i.}} - f_{.j} = f_{.j} \sum_{k=1}^{\tau} \frac{1}{\sqrt{\lambda_k}} c_k(i) \tilde{c}_k(j) .$$

Les composantes principales et les valeurs propres expliquent en quoi les fréquences observées s'écartent des fréquences théoriques.

# Formule de reconstitution des données

- Cette formule, appelée formule de reconstitution des données, permet de recalculer les valeurs du tableau initial en fonction des marges et des facteurs.
- Lorsque l'on dépouille les résultats d'une AFC, on limite généralement l'interprétation aux premiers axes principaux. Cela revient à considérer non pas le tableau des données mais son approximation obtenue à l'aide des premiers termes de la somme ci-dessus:

$$f_{ij} - f_{i.}f_{.j} \simeq f_{i.}f_{.j} \sum_{k=1}^2 \frac{1}{\sqrt{\lambda_k}} c_k(i) \tilde{c}_k(j).$$

En d'autres termes

$$f_{ij} \simeq f_{i.}f_{.j} \left( 1 + \sum_{k=1}^2 \frac{1}{\sqrt{\lambda_k}} c_k(i) \tilde{c}_k(j) \right).$$

**Preuve.** Nous avons

$$Y_r := X_r - \mathbf{1}_p g_r^t,$$

dont ces éléments sont

$$y_{i,j} = \frac{f_{ij}}{f_{i.}} - f_{.j}, \quad i = 1, \dots, p \text{ et } j = 1, \dots, q.$$

Rappelons que  $\{u_1, \dots, u_p\}$  est une base  $M_r$ -orthonormée de  $\mathbb{R}^p$  et les  $c_k$ ,  $k = 1, \dots, q$  sont coordonnées de la  $i$ -ième ligne de  $Y_r$  dans la base  $\{u_1, \dots, u_p\}$ . En d'autres termes

$$Y_r = \sum_{k=1}^p c_k u_k = \sum_{k=1}^{\tau} c_k u_k.$$

Rappelons aussi que

$$u_k = \frac{1}{\sqrt{\lambda_k}} N^t D_r^{-1} \tilde{u}_k = \frac{1}{\sqrt{\lambda_k}} X_r^t \tilde{u}_k.$$

Alors

$$\begin{aligned} Y_r &= \sum_{k=1}^{\tau} \frac{1}{\sqrt{\lambda_k}} c_k N^t D_r^{-1} \tilde{u}_k = \sum_{k=1}^{\tau} \frac{1}{\sqrt{\lambda_k}} c_k D_c (D_c^{-1} N^t) (D_r^{-1} \tilde{u}_k) \\ &= \sum_{k=1}^{\tau} \frac{1}{\sqrt{\lambda_k}} c_k D_c X_c \tilde{w}_k = \sum_{k=1}^{\tau} \frac{1}{\sqrt{\lambda_k}} c_k D_c \tilde{c}_k. \end{aligned}$$

En écrivant cette relation coordonnées par coordonnées, on obtient

$$\frac{f_{ij}}{f_{i.}} - f_{.j} = \frac{1}{\sqrt{\lambda_k}} \sum_{k=1}^{\tau} c_k(i) f_{.j} \tilde{c}_k(j) .$$

En multipliant les deux membres de cette équation par  $f_{i.}$  on obtient

$$f_{ij} - f_{.j} f_{i.} = f_{i.} \sum_{k=1}^{\tau} \frac{1}{\sqrt{\lambda_k}} c_k(i) \tilde{c}_k(j) .$$

- Contribution relative d'inertie
- Qualité de la représentation
- Cercle de corrélations

- Contribution relative du profil-ligne  $i$  de de la matrice  $Y_r$  au  $k$ -ième axe  $u_k$  :

$$Ctr(i, k) := \frac{f_{i.} c_k^2(i)}{\sum_{i=1}^p f_{i.} c_k^2(i)} = \frac{f_{i.} c_k^2(i)}{\lambda_k}, \quad i = 1, \dots, p; \quad k = 1, \dots, \tau.$$

- Contribution relative du profil-colonne  $j$  de de la matrice  $Y_c$  au  $k$ -ième axe  $\tilde{u}_k$  :

$$\widetilde{Ctr}(j, k) = \frac{f_{.j} \tilde{c}_k^2(j)}{\sum_{j=1}^q f_{.j} \tilde{c}_k^2(j)} = \frac{f_{.j} \tilde{c}_k^2(j)}{\lambda_k}, \quad j = 1, \dots, q; \quad k = 1, \dots, \tau.$$

# Qualité de représentation sur un axe

- Qualité de la représentation du profil-lignes  $i$  de la matrice  $Y_r$  au  $k$ —ième axe  $u_k$  :

$$\begin{aligned}\cos^2(\widehat{i, u_k}) &= \left( \frac{\langle i, u_k \rangle_{M_r}}{\|i\|_{M_r} \|u_k\|_{M_r}} \right)^2 = \frac{(i^t M_r u_k)^2}{\|i\|_{M_r}^2}, \quad (\|u_k\|_{M_r} = 1) \\ &= \frac{c_k^2(i)}{\sum_{k=1}^{\tau} c_k^2(i)}, \quad i = 1, \dots, p.\end{aligned}$$

- Qualité de la représentation du profil-colonnes  $j$  de la matrice  $Y_c$  au  $k$ —ième axe  $\tilde{u}$  :

$$\begin{aligned}\cos^2(\widehat{j, \tilde{u}_k}) &= \left( \frac{\langle j, \tilde{u}_k \rangle_{M_c}}{\|j\|_{M_c} \|\tilde{u}_k\|_{M_c}} \right)^2 = \frac{\tilde{c}_k^2(j)}{\|j\|_{M_c}^2}, \quad (\|\tilde{u}_k\|_{M_c} = 1) \\ &= \frac{\tilde{c}_k^2(j)}{\sum_{k=1}^{\tau} \tilde{c}_k^2(j)}, \quad j = 1, \dots, q.\end{aligned}$$



**Proposition.** Pour chaque  $j \in \{1, \dots, q\}$  et  $k \in \{1, \dots, \tau\}$ , on a

$$\mathbf{Cor}(c_k, Y_r^{(j)}) = \frac{\sqrt{\lambda_k}}{s_j} u_k(j),$$

où  $s_j = \sqrt{\mathbf{Var}(Y_r^{(j)})}$  et  $u_k(j)$  désigne la  $j$ -ème composante du vecteur propre  $u_k$ .

**Preuve.** On a

$$\mathbf{Cor}(c_k, Y_r^{(j)}) = \frac{\mathbf{Cov}(c_k, Y_r^{(j)})}{\sqrt{\mathbf{Var}(c_k) \mathbf{Var}(Y_r^{(j)})}} = \frac{1}{\sqrt{\lambda_k}} \frac{c_k^t D_r Y_r^{(j)}}{s_j}.$$

Remarquons que  $Y_r^{(j)} = Y_r \delta$ , où  $\delta := (0, \dots, 1, \dots, 0)^t$  est un vecteur à  $q$  composantes. Donc

$$\begin{aligned} \mathbf{Cor}(j, k) &= \frac{1}{\sqrt{\lambda_k}} \frac{(Y_r w_k)^t D_r Y_r \delta}{s_j} = \frac{1}{\sqrt{\lambda_k}} \frac{(Y_r M_r u_k)^t D_r Y_r \delta}{s_j} \\ &= \frac{1}{\sqrt{\lambda_k}} \frac{u_k^t M_r Y_r^t D_r Y_r \delta}{s_j} = \frac{1}{\sqrt{\lambda_k}} \frac{u_k^t M_r V_r \delta}{s_j} \\ &= \frac{1}{\sqrt{\lambda_k}} \frac{(V_r M_r u_k)^t \delta}{s_j} = \frac{1}{\sqrt{\lambda_k}} \frac{\lambda_k u_k^t \delta}{s_j} \\ &= \frac{\sqrt{\lambda_k}}{s_j} u_k^t \delta = \frac{\sqrt{\lambda_k}}{s_j} u_k(j). \end{aligned}$$