

Analyse de données

Chap I Objectif de l'ADD.

Domaine qui regroupe diverses méthodes mathématiques. On souhaite obtenir une organisation de ces données

- Analyse en composantes principales (ACP)
- Analyse factorielle des correspondances (AFC)
- Analyse des composantes multiples (ACM)

Analyse factorielle

Problématique:

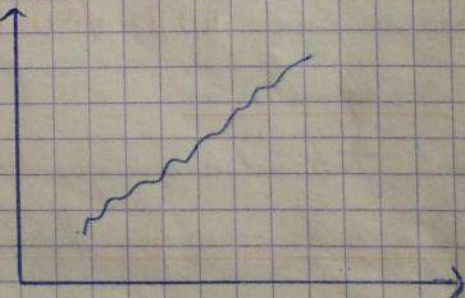
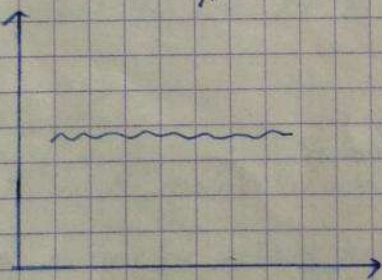
Comment visualiser et interpréter un nuage de points dans un espace de dimension très grande ?

Idee: Définir la "meilleure" projection possible de nuage de points.

Demarche:

Traitement de données initiales afin d'éliminer les biais statistiques et les effets d'échelle.

Effet d'échelle: différentes allures pour une même courbe



- Quantifier la notion d'"information"
 - notion de variance, covariance, corrélation.
 - introduire 7 notions de distance.
- Réaliser des changements d'axes: axes hiérarchisés selon la part d'information contenue
- Choisir les axes de projection en contrôlant la perte d'information
- Réaliser la projection et effectuer l'analyse en la validant par les 7 paramètres que l'on aura obtenus lors de la méthode.

II Type de données

1) Tableau de données

- Les lignes correspondent aux individus statistiques (sauf AFC)
- Les colonnes: * variables statistiques
* modalités des variables statistiques.

2) Variables

- * En ACP: variables **quantitatives** de même importance (\approx) \rightarrow discrète
- * En AFC: 2 variables sur une population (**qualitative** ou **quantitative**)
- * En ACM: des variables sur une population (**qualitative** ou ")

III Calcul matriciel et ADD

1. Matrices symétriques

Propriété: Soit M une matrice quelconque alors les matrices ${}^t M \times M$ et $M \times {}^t M$ sont symétriques

2. Distance

Si u et v , deux vecteurs ^{colonne} dans une base orthonormée

Produit scalaire de U et V : ${}^t U \times V$

Distance euclidienne: $\|U\| = \sqrt{{}^t U \times U}$

3. Diagonalisation

Propriété: Toute matrice symétrique est diagonalisable dans \mathbb{R}

4. Approche géométrique ellipsoïde d'inertie

cf annexes.

Chapitre II

Paramètres d'un nuage de points. Approche statistique.

I Paramètres statistiques

Soit n l'effectif total, X variable, α modalité de X .

moyenne: $\bar{X} = E(X) = \frac{1}{n} \sum \alpha$

variance: $V(X) = E((X-E(X))^2) = \frac{1}{n} \sum (\alpha - \bar{X})^2$

Ecart-type: $\sigma_X = \sqrt{V(X)}$

Covariance: $\text{cov}(X, Y) = E(X-E(X))(Y-E(Y)) = \frac{1}{n} \sum (\alpha - \bar{X})(\beta - \bar{Y})$

Corrélation $\rho(X, Y) = \frac{\text{cov}(X, Y)}{\sigma_X \sigma_Y}$ avec $-1 \leq \rho(X, Y) \leq 1$

variable centrée: $X^c = X - \bar{X}$ $\bar{X}^c = 0$

variable centrée et réduite $X^d = \frac{X - \bar{X}}{\sigma_X}$ $\bar{X}^d = 0, \sigma_{X^d} = 1$

II Matrices fondamentales

1. Matrice de variance-covariance.

Definition: X^c désignant la matrice des données centrées. On appelle matrice de variance-covariance la matrice $\Sigma = \frac{1}{n} {}^t X^c \cdot X^c$

Remarque: Σ matrice symétrique.

Matrice de variance-covariance: $\Sigma = \begin{pmatrix} \text{var}(X_1^c) & \text{cov}(X_1^c, X_2^c) & \text{cov}(X_1^c, X_3^c) \\ \text{cov}(X_2^c, X_1^c) & \text{var}(X_2^c) & \text{cov}(X_2^c, X_3^c) \\ \text{cov}(X_3^c, X_1^c) & \text{cov}(X_3^c, X_2^c) & \text{var}(X_3^c) \end{pmatrix}$

$$M = \begin{bmatrix} a_1 & b_1 & c_1 \\ a_2 & b_2 & c_2 \end{bmatrix}$$

$$M^c = \begin{bmatrix} a_1 - \bar{A} & b_1 - \bar{B} & c_1 - \bar{C} \\ a_2 - \bar{A} & b_2 - \bar{B} & c_2 - \bar{C} \end{bmatrix}$$

$${}^t M^c = \begin{bmatrix} a_1 - \bar{A} & a_2 - \bar{A} \\ b_1 - \bar{B} & b_2 - \bar{B} \\ c_1 - \bar{C} & c_2 - \bar{C} \end{bmatrix}$$

$$\Sigma = \begin{bmatrix} \text{var}(A) & \text{cov}(A, B) & \text{cov}(A, C) \\ \text{cov}(B, A) & \text{var}(B) & \text{cov}(B, C) \\ \text{cov}(C, A) & \text{cov}(C, B) & \text{var}(C) \end{bmatrix}$$

trace $\Sigma = \text{var}(A) + \text{var}(B) + \text{var}(C)$

Matrice symétrique, carré, d'ordre égal au nombre de variable

Une covariance peut être négative

Une variance est toujours positive.

Chapitre II

Paramètres d'un nuage de points. Approche statistique.

I Paramètres statistiques

Soit n l'effectif total, X variable, α modalité de X .

moyenne: $\bar{X} = E(X) = \frac{1}{n} \sum \alpha$

variance: $V(X) = E((X - E(X))^2) = \frac{1}{n} \sum (\alpha - \bar{X})^2$

Ecart-type: $\sigma_X = \sqrt{V(X)}$

Covariance: $\text{cov}(X, Y) = E(X - E(X))(Y - E(Y)) = \frac{1}{n} \sum (\alpha - \bar{X})(\beta - \bar{Y})$

Corrélation $\rho(X, Y) = \frac{\text{cov}(X, Y)}{\sigma_X \sigma_Y}$ avec $-1 \leq \rho(X, Y) \leq 1$

variable centrée: $X^c = X - \bar{X}$ $\bar{X}^c = 0$

variable centrée et réduite $X^d = \frac{X - \bar{X}}{\sigma_X}$ $\bar{X}^d = 0, \sigma_{X^d} = 1$

II Matrices fondamentales

1. Matrice de variance-covariance.

Definition: X^c désignant la matrice des données centrées. On appelle matrice de variance-covariance la matrice $\Sigma = \frac{1}{n} X^c \cdot X^c$

Remarque: Σ matrice symétrique.

Matrice de variance-covariance: $\Sigma = \begin{pmatrix} \text{var}(X_1^c) & \text{cov}(X_1^c, X_2^c) & \text{cov}(X_1^c, X_3^c) \\ \text{cov}(X_2^c, X_1^c) & \text{var}(X_2^c) & \text{cov}(X_2^c, X_3^c) \\ \text{cov}(X_3^c, X_1^c) & \text{cov}(X_3^c, X_2^c) & \text{var}(X_3^c) \end{pmatrix}$

$$M = \begin{bmatrix} a_1 & b_1 & c_1 \\ a_2 & b_2 & c_2 \end{bmatrix}$$

$$M^c = \begin{bmatrix} a_1 - \bar{A} & b_1 - \bar{B} & c_1 - \bar{C} \\ a_2 - \bar{A} & b_2 - \bar{B} & c_2 - \bar{C} \end{bmatrix}$$

$$M^c = \begin{bmatrix} a_1 - \bar{A} & a_2 - \bar{A} \\ b_1 - \bar{B} & b_2 - \bar{B} \\ c_1 - \bar{C} & c_2 - \bar{C} \end{bmatrix}$$

$$\Sigma = \begin{bmatrix} \text{var}(A) & \text{cov}(A, B) & \text{cov}(A, C) \\ \text{cov}(B, A) & \text{var}(B) & \text{cov}(B, C) \\ \text{cov}(C, A) & \text{cov}(C, B) & \text{var}(C) \end{bmatrix}$$

trace $\Sigma = \text{var}(A) + \text{var}(B) + \text{var}(C)$

Matrice symétrique, carré, d'ordre égal au nombre de variable

Une covariance peut être négative

Une variance est toujours positive.

2. Matrice des corrélations R.

Def: X^S (matrice centrée réduite). On appelle la matrice de corrélation R:

$$R = \frac{1}{n} {}^t X^S \cdot X^S$$

avec $R = \begin{bmatrix} 1 & r(x_1^S, x_2^S) & r(x_1^S, x_3^S) \\ r(x_2^S, x_1^S) & 1 & r(x_2^S, x_3^S) \\ r(x_3^S, x_1^S) & r(x_3^S, x_2^S) & 1 \end{bmatrix}$

AFC à récupérer.

Chapitre III

I Aperçu général

p variables quantitatives X_1, X_2, \dots, X_p . Population de n individus.

Si X_1, X_2, \dots, X_p sont corrélés, il existe des redondances.

Ideée de l'ACP

"diminuer les redondances": Transformer les p variables (initiales) X_1, \dots, X_p en p nouvelles variables \rightarrow facteurs f_1, \dots, f_p .

• conserver k facteurs ($k \leq p$) en conservant un maximum d'information.

redondance = corrélation.

\rightarrow Il faut que les facteurs ne soient pas corrélés.

1) La transformation

transformation: $(X_1, \dots, X_p) \rightarrow (f_1, \dots, f_p)$

compression: extraire des p facteurs, k facteurs \rightarrow paramètre: qualité globale d'explication = η^2

transformation: Σ ou $R \rightarrow$ matrice diagonale \rightarrow DIAGONALISATION de la matrice Σ ou R .

Soit U la base de vecteurs propres \rightarrow facteurs

Changement de base sur les données initiales: $A^c \xrightarrow{A^s} F = \begin{cases} A^c U \\ A^s U \end{cases}$ Matrice des composantes principales.
 \rightarrow = coordonnées des individus dans la nouvelle base.

On a donc:

* A^c ou A^s matrice d'origine.

* U matrice de vecteurs propres de la diagonalisation de Σ ou de R .

* Λ matrice diagonale des valeurs propres.

* F matrice des composantes principales.

2) L'interprétation

* "identifier" les facteurs = donner une signification concrète aux facteurs

idée: calculer les corrélations entre les variables et les facteurs

\rightarrow matrice de saturation $S = \frac{1}{n} X^s F \Lambda^{-1/2}$

II ACP d'un exemple.

cf Annexe avec X: prix d'un litre de lait
Y: " " d'eau
Z: " " d'huile.

$A^S \rightarrow$ Matrice Centree reduite \rightarrow moyennes = 0.
Ecart type = 0

matrice correlation \rightarrow termes de la diagonale = 1.

trace = 1+1+1 = 3 = p . = variance totale (permet de mesurer la qte totale d'inf)

3^{em} étape: diagonaliser.

determiner ses valeurs propres puis vecteurs propres.

$\det(R - \lambda I) = 0$ (det valeurs propres). On obtient p valeurs propres.

On classe de la + grande à la + petite: $[\begin{smallmatrix} \lambda_1 & & \\ & \lambda_2 & \\ & & \lambda_3 \end{smallmatrix}]$

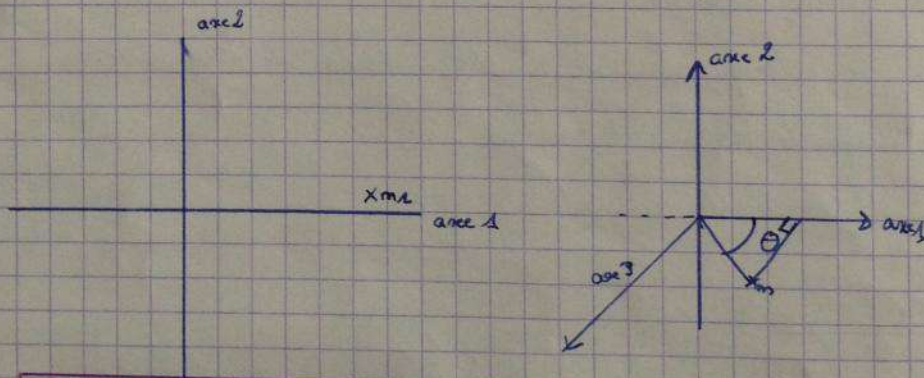
* Vecteurs propres: Si λ est une valeur propre d'une matrice M , un vecteur propre associé à λ est un vecteur u tel que $MU = \lambda U$

Règle: On choisit des vecteurs propres unitaires ou normés, c-à-d de norme = 1.

Propriété les vecteurs propres sont orthogonaux 2 à 2.

4^{em} étape: calcul des composantes principales.

Qualité de représentation d'un individu sur un axe. (qte ou \cos^2)



$$\cos^2(mi; \text{axe } k) = \frac{\text{coord}(mi, \text{axe } k)^2}{\sum \text{coord}(mi, \text{axe } j)^2}$$

Si $qte \approx 1$ alors mi proche de l'axe k .
Si $qte \approx 0$ " " éloigné " k .

Matrice de saturation: - toujours carré

- indique les corrélations entre les variables et les facteurs

- IDC (indice de corrélation) toujours compris entre -1 et 1.

Corrélations: facteur 1 facteur 2 facteur 3

X^s

Y^s

Z^s

Centre de corrélation.

$$S = [d_{ij}]_{\substack{1 \leq i \leq p \\ 1 \leq j \leq p}} \quad \sum_{j=1}^p d_{ij}^2 = 1 \quad (\text{ligne})$$

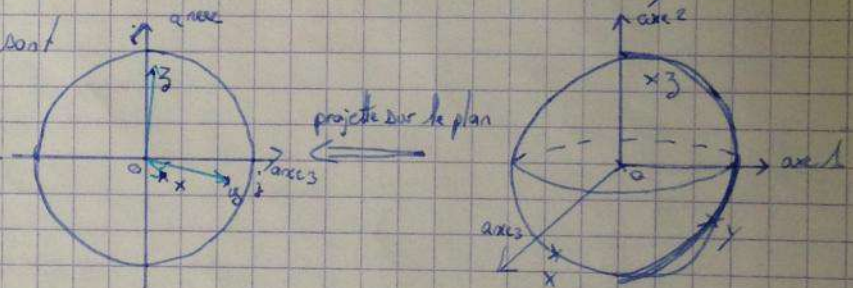
$$\sum_{i=1}^p d_{ij}^2 = \lambda_j \quad (\text{colonne})$$

$$p=3 \quad d_{11}^2 + d_{12}^2 + d_{13}^2 = 1$$

$$\Rightarrow x^2 + y^2 + z^2 = 1$$

• Les points variables sont situés sur une hypersphère de centre 0 et de rayon 1.

• Les projections des points variables sont situés à l'intérieur du cercle de centre 0 de rayon 1.



Interpretation:

En ACP essaye d'identifier les axes. Essaye de donner un nom aux différents axes

Avec quelles variables les axes sont corrélés.

A proximité des axes et des points

2 variables corrélés: si l'un augmente, l'autre augmente

Matrice de corrélation - élément de contrôle

ACP: - variables quantitatives

- variables à peu près la même importance.

- nature des données à observer

ACP centrée X^c, Σ

ACP centrée réduite X^D, R

3 cas de figures pour une ACP : - variable longueur (m: entre 3 et 5, cm: 300 et 500) unités différentes
centrée réduite \Rightarrow variance + grande dans le 2^o cas que dans le 1^{er}

Matrice centrée réduite: Σ colonne \leftarrow

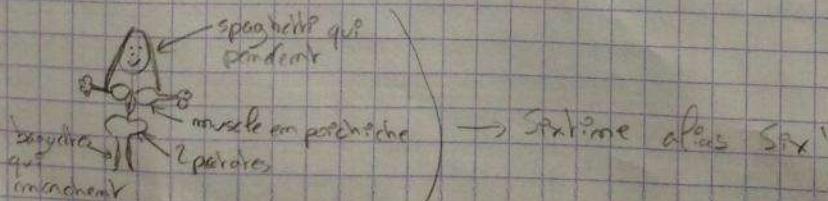
Σ valeur propre = nb de variables = trace de R.

qge: qualité globale d'explication. (mesure la qualité de représentation d'un individu par rapport à l'axe)

Matrice de corrélat: corrélation entre les variables $-1 < 1$ (forte corrélation)

Σ carré (des coeffs) = 1 (des lignes)
colonne = (valeur propre)²

Matrice des saturations:



Analyse factorielle des correspondances (AFC)

Objectif: Etudier les liens éventuels entre 2 variables sur une population.
On est face à 2 variables généralement qualitatives

I Tableau de données

1) Le tableau des effectifs (k)

2 variables X et Y, possédant un certain nombre de modalités

n = nombre de modalités de X → nombre de lignes

p = " " de Y → colonnes.

N = effectif total de la population

k_{ij} = effectif de la modalité $X = X_i$ et $Y = Y_j$
 k_{i0} =
 k_{0j} =

		Y_1	Y_2	...	Y_p	
X	X_1	k_{11}	k_{12}		k_{1p}	k_{10}
	X_2	k_{21}	k_{22}			k_{20}
	...					
	X_n	k_{n1}	k_{n2}		k_{np}	k_{n0}
		k_{01}	k_{02}			k_{0p}

$$\sum_{i=1}^n \sum_{j=1}^p k_{ij} = N$$

$$\sum_{i=1}^n k_{i0} = N$$

$$\sum_{j=1}^p k_{0j} = N$$

Données initiales

			chanson	page	$k_{i.}$
		C	J		
jeune →	J	40	20		60
	AF	40	10		50
	AM	30	20		50
Vieux →	V	10	30		40
	$k_{.j}$	120	80		200

2) Le tableau des fréquences

tableau des fréquences par rapport à l'effectif total

f_{ij} = fréquence de $X = X_i$ et $Y = Y_j$

$$f_{ij} = p(X = X_i \cap Y = Y_j)$$

$$f_{ij} = \frac{k_{ij}}{N}$$

$X \setminus Y$	Y_1		
X_i	f_{ij}		$f_{i.}$
	$f_{.j}$		4

$f_{i.}$ = fréquence de $X = X_i$

$$f_{i.} = \frac{k_{i.}}{N} = \sum_{j=1}^m b_{ij}$$

$f_{.j}$ = fréquence de $Y = Y_j$

$$f_{.j} = \frac{k_{.j}}{N} = \sum_{i=1}^n b_{ij}$$

$$f_{ij} = p(Y = Y_j)$$

(cf. tableau des fréquences)

	C	J	
J	0,18	0,12	$f_{.1} = 0,3$
AF	0,15	0,1	$f_{.2} = 0,25$
AM	0,15	0,1	$f_{.3} = 0,25$
V	0,12	0,08	0,2
	$f_{.1} = 0,6$	$f_{.2} = 0,4$	1

3. Les tableaux des profils- lignes et des profils colonnes.

a) Le tableau des PFL

X \ Y	Y_j	
X_i		
X_n		
$f_{i.}$	$f_{.j}$	$f_{..}$

PFL moyen $\rightarrow p(Y = Y_j)$

$$\frac{k_{ij}}{k_{i.}} = \frac{\frac{k_{ij}}{N}}{\frac{k_{i.}}{N}} = \frac{f_{ij}}{f_{i.}} = \frac{p(X = X_i \cap Y = Y_j)}{p(X = X_i)} = p(Y = Y_j / X = X_i)$$

Indépendance $\Leftrightarrow p(X = X_i \cap Y = Y_j) = p(X = X_i) \cdot p(Y = Y_j)$

$$\Leftrightarrow p(Y = Y_j / X = X_i) = p(Y = Y_j)$$

Y \ X	C	J	
J	4/60	1/10	1
AF	0,8	0,2	1
AM	0,6	0,4	1
V	0,25	0,75	1
PFM	0,6	0,4	1

profil des colonnes $b_{.j}$

b) Le tableau des PFC

X \ Y	Y_j	$f_{i.}$ (PFL moyen)	
X_i		$f_{i.}$	
	1	1	1

$$\frac{k_{ij}}{k_{.j}} = \frac{f_{ij}}{f_{.j}} = \frac{p(X = X_i \cap Y = Y_j)}{p(Y = Y_j)} = p(X = X_i / Y = Y_j)$$

X \ Y	C	J	
J			0,3
AF			0,25
AM			0,25
V			0,2
	1	1	1

PFL: 1 point dans \mathbb{R}^p
 PFC: 1 point dans \mathbb{R}^p

4. Les degrés de liberté

PFL: Somme des composantes d'un PFL = 1 \Rightarrow Chaque PFL vérifie

$$x_1 + x_2 + \dots + x_p = 1 \text{ hyperplan de dim} = p-1$$

\Rightarrow $p-1$ degrés de liberté pour les PFL.

II AFC

1. La distance du X^2 (KHI - Deux)

c) Annexe

2. L'ajustement du nuage des n PFL dans l'espace des p colonnes : matrice à diagonaliser

$$\|u\|^2 = \sum_{i=1}^p u_i^2$$

$$d(x^2, e) = \sum_{j=1}^p \frac{1}{f_{j\cdot}} \left(\frac{f_{ij}}{f_{i\cdot}} - \frac{f_{j\cdot}}{f_{e\cdot}} \right)^2 = \sum_{j=1}^p \left(\frac{b_{ij}}{f_{i\cdot} \sqrt{f_{j\cdot}}} - \frac{f_{e_j}}{f_{e\cdot} \sqrt{f_{j\cdot}}} \right)^2$$

$$A = \begin{bmatrix} b_{ij} \\ \frac{b_{ij}}{f_{i\cdot} \sqrt{f_{j\cdot}}} \end{bmatrix}_{\substack{1 \leq i \leq m \\ 1 \leq j \leq p}} \quad \text{sur PFC} \quad B = \frac{b_{ij}}{f_{i\cdot} \sqrt{f_{j\cdot}}}$$

En pratique $X^* = \begin{bmatrix} b_{ij} \\ \frac{b_{ij}}{\sqrt{f_{i\cdot} f_{j\cdot}}} \end{bmatrix}_{\substack{1 \leq i \leq m \\ 1 \leq j \leq p}}$

La matrice à diagonaliser : $T^* = {}^t X^* X^*$

$X^* \rightarrow$ dimension " $m \times p$ " \rightarrow m lignes
 p colonnes.

${}^t X^* \rightarrow$ " $p \times m$ "

$T^* =$ matrice carrée, symétrique, de dimension " p ".

3. Propriété des valeurs propres.

Prop: la matrice à diagonaliser possède toujours la valeur propre 1, appelée valeur propre triviale.

Ex: $X^* = \begin{bmatrix} b_{ij} \\ \frac{b_{ij}}{\sqrt{f_{i\cdot} f_{j\cdot}}} \end{bmatrix}_{\substack{1 \leq i \leq m \\ 1 \leq j \leq p}} \quad {}^t X^* = \begin{bmatrix} \frac{b_{ki}}{\sqrt{f_{k\cdot} f_{i\cdot}}} \end{bmatrix}_{\substack{1 \leq k \leq p \\ 1 \leq i \leq m}} \quad \begin{bmatrix} a_{ik} \end{bmatrix}_{\substack{1 \leq i \leq p \\ 1 \leq k \leq p}} \times \begin{bmatrix} b_{kj} \end{bmatrix}_{\substack{1 \leq k \leq p \\ 1 \leq j \leq p}} = \begin{bmatrix} \sum a_{ik} b_{kj} \end{bmatrix}_{\substack{1 \leq i \leq p \\ 1 \leq j \leq p}}$

$$T^* = {}^t X^* X^* = \begin{bmatrix} \frac{b_{ki}}{\sqrt{f_{k\cdot} f_{i\cdot}}} \end{bmatrix}_{\substack{1 \leq k \leq p \\ 1 \leq i \leq p}} \begin{bmatrix} \frac{b_{ij}}{\sqrt{f_{i\cdot} f_{j\cdot}}} \end{bmatrix}_{\substack{1 \leq i \leq p \\ 1 \leq j \leq p}} =$$

$$= \begin{bmatrix} \sum_{i=1}^m \frac{b_{ki} \times b_{ij}}{f_{k\cdot} \sqrt{f_{i\cdot} f_{j\cdot}}} \end{bmatrix}_{\substack{1 \leq i \leq p \\ 1 \leq j \leq p}}$$

$$= \begin{bmatrix} 1 & \sum_{i=1}^m \frac{b_{ki} \times b_{ij}}{f_{k\cdot}} \\ \sqrt{f_{i\cdot} f_{j\cdot}} & f_{k\cdot} \end{bmatrix}_{\substack{1 \leq i \leq p \\ 1 \leq j \leq p}}$$

$$T^* G = \begin{bmatrix} t_{ik} \end{bmatrix}_{\substack{1 \leq i \leq p \\ 1 \leq k \leq p}} \times \begin{bmatrix} \sqrt{f_{ok}} \end{bmatrix}_{\substack{1 \leq k \leq p}} = \begin{bmatrix} \sum_{k=1}^p t_{ik} \times \sqrt{f_{ok}} \end{bmatrix}_{\substack{1 \leq i \leq p}}$$

$$\stackrel{?}{=} G = \begin{bmatrix} \sqrt{f_{oi}} \end{bmatrix}_{\substack{1 \leq i \leq p}}$$

$$= \sum_{k=1}^p t_{ik} \sqrt{f_{ok}}$$

ACP centra réduite X^D, R

2 cm de lignes pour une ACP: - variable longueur (m: entre 3 et 5, cm: 300 et 500) unités différentes
- les unités \Rightarrow variance + grande dans le 2^e cas que dans le 1^{er}

$$\sum_{j=1}^p T_{ij} \times \sqrt{f_{oj}} = \sum_{j=1}^p \frac{1}{\sqrt{f_{oi} \times f_{oj}}} \sum_{k=1}^m \frac{f_{ki} \times f_{kj}}{f_{ko}} \times \sqrt{f_{oj}} = \frac{1}{\sqrt{f_{oi}}} \sum_{k=1}^m \frac{f_{ki}}{f_{ko}} \sum_{j=1}^p f_{kj}$$
$$= \frac{1}{\sqrt{f_{oi}}} \sum_{k=1}^m \frac{f_{ki}}{f_{ko}} \times f_{ko} = \frac{1}{\sqrt{f_{oj}}} \sum_{k=1}^m f_{ki} = \frac{f_{oi}}{\sqrt{f_{oi}}} = \sqrt{f_{oi}} \text{ donc } T^+ G = G$$

Inertie totale: L'inertie totale matrice R est égale à la somme des valeurs propres non triviales $R = \sum \lambda_i - 1$

h. Coordonnées des points-lignes

Coordonnée du point ligne i sur l'axe α (engendré par le vecteur propre noté u_α)

$$\psi_{\alpha i} = \sum_{j=1}^p \frac{f_{ij}}{f_{oi} \sqrt{f_{oj}}} \times u_{\alpha j}$$

projection du point i
sur l'axe engendré par
le vc. pr. u_α = coordonnées de
 i sur α

$$\Psi = AU$$

5. Etude des PFC / Ajustement du nuage des p colonnes dans l'espace des m lignes

Matrice à diagonaliser: $W^* = X^{*t} X^* \Rightarrow$ matrice symétrique

$$\text{Matrice des données: } B = \left[\frac{f_{ji}}{f_{oi} \sqrt{f_{oj}}} \right]$$

6. Coordonnées des points-colonnes

$$\phi_{\alpha j} = \sum_{i=1}^m \frac{f_{ij}}{f_{oj} \sqrt{f_{oi}}} \times v_{\alpha i}$$

$$\Phi = BV$$

Analyse des Correspondances Multiples

⑤

Méthode permettant d'étudier les liaisons entre plusieurs variables qualitatives ou quantitatives, c'est donc une généralisation de l'AC.

Q_i = Variable i .

Q_{ij} = valeur de Q_i pour l'individu j .

n = nombre d'individus

q = nombre de variables.

m_k = nombre de modalités de Q_k .

M = nombre total de modalités.

1. Tableau descriptif complet.

On associe à chaque modalité une variable appelée **variable indicatrice**, ne prenant les valeurs 0 et 1 selon que la modalité est vérifiée ou non.

X_j = variable indicatrice de la modalité j

$X_j(k)$ = valeur (0 ou 1) de la variable indicatrice X_j pour l'individu k .

→ Éléments de chaque ligne = q .

Éléments de la colonne associée à X_j ← m_j

→ La fréquence associée, appelée **fréquence marginale** ← p_j

$$p_j = \frac{m_j}{nq}$$

ACM = $n \times \frac{m_j}{m_k}$ d'individus très grand
nb de modalités

Les individus doivent être répartis assez régulièrement.

2. Descriptif de la méthode.

Basé sur l'étude des distances entre PFL et PFC et sur l'analyse des PFL et PFC.

1. On ne s'intéresse pas à l'indépendance.

• Notion de distance (du X^i)

$$d^2(i; l) = \frac{\sum_{j=1}^m \left(\frac{X_j(i)}{q} - \frac{X_j(l)}{q} \right)^2}{p_j} = \frac{1}{q^2} \sum_{j=1}^m \frac{(X_j(i) - X_j(l))^2}{p_j}$$

Ordre de grandeur : fréquence marginale colonne très faible (très peu stable) $\rightarrow p_j$ très petite,
donc quand on divise donne un nombre très grande \rightarrow poids très fort.

\rightarrow D'où regroupement de modalités.

Modalités

$$d^2(j; l) = \frac{\sum_{i=1}^n \left(\frac{X_j(i)}{n_j} - \frac{X_l(i)}{n_l} \right)^2}{\frac{1}{n}} = n \sum_{i=1}^n \left(\frac{X_j(i)}{n_j} - \frac{X_l(i)}{n_l} \right)^2$$

Pense à la xaxcti' à éviter (en regroupant les effectifs)

+ modalité faible, + individus éloignés.

+ eff. modalité faible, + distance entre mod. est grande.

Analyse mathématiques

Même méthode qu'en AFC.

L'inertie totale est égale à $\frac{m}{q} - 1$