

Variables qualitatives : analyse des correspondances

Jean-Marc Lasgouttes

<http://ana-donnees.lasgouttes.net/>

L'analyse factorielle des correspondances

But On cherche à décrire la liaison entre deux variables qualitatives.

Exemple on peut regarder la répartition de la couleur des yeux en fonction de la couleur des cheveux.

Différence avec l'ACP l'ACP se fait dans un cadre différent ; les variables sont quantitatives et donc

- il est possible de faire des opérations mathématiques sur les valeurs des variables ;
- par contre, il n'est en général pas possible de compter les individus qui ont une caractéristique donnée (taille=1, 83m)

Pourquoi deux variables ? le cas de plus de deux variables est l'analyse de correspondance multiples, traité plus tard dans le cours.

Partie I. Les données qualitatives

Variables qualitatives

Soit \mathcal{X} une variable qualitative. On dispose d'un échantillon de n individus sur lesquels la variable est mesurée.

Modalités (ou catégories) les valeurs que peut prendre une variable qualitative ; si la variable a m modalités (valeurs possibles), on note x_i , $1 \leq i \leq m$, ces modalités, ou plus simplement i .

Effectif le nombre d'occurrence de la modalité i dans l'échantillon ; on le note n_i , et on a $\sum_{i=1}^m n_i = n$.

Profil c'est l'ensemble des valeurs n_i/n ; la somme du profil sur les modalités est 1.

Tableau de contingence

Soient \mathcal{X}_1 et \mathcal{X}_2 deux variables qualitatives à m_1 et m_2 modalités respectivement décrivant un ensemble de n individus.

Définition le tableau de contingence est une matrice à m_1 lignes et m_2 colonnes renfermant les effectifs n_{ij} d'individus tels que $\mathcal{X}_1 = i$ et $\mathcal{X}_2 = j$.

$$\mathbf{N} = \begin{pmatrix} n_{11} & n_{12} & \cdots & n_{1m_2} \\ n_{21} & n_{22} & \ddots & \vdots \\ \vdots & & n_{ij} & \vdots \\ n_{m_1 1} & & & n_{m_1 m_2} \end{pmatrix}$$

La constitution de ce tableau est aussi appelé un « tri croisé ».

Marges et profils

Marge en ligne c'est la somme $n_{i.} = \sum_{j=1}^{m_2} n_{ij}$, c'est-à-dire l'effectif total de la modalité i de \mathcal{X}_1 .

On définit aussi le profil marginal des lignes $n_{i.}/n$.

Marge en colonne c'est la somme $n_{.j} = \sum_{i=1}^{m_1} n_{ij}$, c'est-à-dire l'effectif total de la modalité j de \mathcal{X}_2 .

On définit aussi le profil marginal des colonnes $n_{.j}/n$.

Deux lectures possibles selon la variable que l'on privilégie, on peut définir

- le tableau des *profils-lignes* $n_{ij}/n_{i.}$, qui représente la fréquence de la modalité j conditionnellement à $\mathcal{X}_1 = i$; la somme de chaque ligne est ramenée à 100%.
- le tableau des *profils-colonnes* $n_{ij}/n_{.j}$, qui représente la fréquence de la modalité i conditionnellement à $\mathcal{X} = j$; la somme de chaque colonne est ramenée à 100%.

Propriétés des profils

Moyenne la moyenne des profils-lignes (avec poids correspondant aux profils marginaux des lignes) est le profil marginal des colonnes :

$$\sum_{i=1}^{m_1} \frac{n_{i.}}{n} \times \frac{n_{ij}}{n_{i.}} = \frac{n_{.j}}{n},$$

et de même pour les colonnes $\sum_{j=1}^{m_2} \frac{n_{.j}}{n} \times \frac{n_{ij}}{n_{.j}} = \frac{n_{i.}}{n}$.

Indépendance empirique lorsque tous les profils lignes sont identiques, il y a indépendance entre \mathcal{X}_1 et \mathcal{X}_2 , puisque la connaissance de \mathcal{X}_1 ne change pas la répartition de \mathcal{X}_2 . On a pour tout j

$$\frac{n_{1j}}{n_{1.}} = \frac{n_{2j}}{n_{2.}} = \cdots = \frac{n_{m_1 j}}{n_{m_1.}} = \frac{n_{1j} + \cdots + n_{m_1 j}}{n_{1.} + \cdots + n_{m_1.}} = \frac{n_{.j}}{n}$$

et donc $n_{ij} = \frac{n_{i.} \cdot n_{.j}}{n}$.

Partie II. Géométrie de nuages de profils

Analyse des correspondances de deux variables : les données

Effectifs on a un tableau de contingence N à m_1 lignes et m_2 colonnes résultant du croisement de deux variables qualitatives \mathcal{X}_1 et \mathcal{X}_2 à m_1 et m_2 modalités respectivement. On note D_1 et D_2 les matrices diagonales des effectifs marginaux

$$D_1 = \begin{pmatrix} n_{1.} & & & 0 \\ & n_{2.} & & \\ & & \ddots & \\ 0 & & & n_{m_1.} \end{pmatrix}$$

$$D_2 = \begin{pmatrix} n_{.1} & & & 0 \\ & n_{.2} & & \\ & & \ddots & \\ 0 & & & n_{.m_2} \end{pmatrix}$$

Profils le tableau des profils des lignes $n_{ij}/n_{i.}$ est donné par $D_1^{-1}N$ et celui des profils des colonnes $n_{ij}/n_{.j}$ par ND_2^{-1} .

Représentation géométrique des profils

Nuage de points les profils-lignes forment un nuage de m_1 points de \mathbb{R}^{m_2} . Chaque point est affecté d'un poids égal à sa fréquence marginale $n_{i.}/n$, et la matrice des poids est donc $\frac{1}{n}D_1$.

Centre de gravité c'est le profil marginal des colonnes car

$$g_\ell = \frac{1}{n}(D_1^{-1}N)'D_1\mathbf{1}_{m_1} = \left(\frac{n_{.1}}{n}, \dots, \frac{n_{.m_2}}{n}\right)'$$

Profils-colonnes les lignes du tableau $D_2^{-1}N'$ forment un nuage de m_2 points de \mathbb{R}^{m_1} , avec matrice de poids $\frac{1}{n}D_2$ et centre de gravité

$$g_c = \left(\frac{n_{1.}}{n}, \dots, \frac{n_{m_1.}}{n}\right)'$$

Comment étudier ces données

Cas indépendant en cas d'indépendance empirique, on a

$$\frac{n_{ij}}{n_{i.}} = \frac{n_{.j}}{n} \text{ et } \frac{n_{ij}}{n_{.j}} = \frac{n_{i.}}{n}.$$

Les deux nuages sont donc réduits à leurs centres de gravité.

Dimension des nuages comme les profils somment à 1, les m_1 profils-lignes sont situés dans le sous-espace W_1 de dimension $m_2 - 1$ défini par $\sum_{j=1}^{m_2} x_j = 1$ et $x_j \geq 0$.

ACP l'étude de la forme des nuages au moyen de l'analyse en composantes principales permettra de rendre compte de la structure des écarts à l'indépendance.

Partie III. L'AFC : une ACP sur un nuage de profils

La métrique du χ^2

Profils-lignes la norme du χ^2 de la différence de profils-lignes $e_i - e_{i'}$ est définie par

$$\|e_i - e_{i'}\|_{\chi^2} = \sum_{j=1}^{m_2} \frac{n}{n_{.j}} \left(\frac{n_{ij}}{n_{i.}} - \frac{n_{i'j}}{n_{i'.}}\right)^2,$$

ce qui revient à utiliser la métrique diagonale nD_2^{-1} .

Inertie l'inertie totale du nuage des profils-lignes par rapport à g_ℓ est

$$I_{g_\ell} = \sum_{i=1}^{m_1} \frac{n_{i.}}{n} \|e_i - g_\ell\|_{\chi^2}^2 = \sum_{i=1}^{m_1} \sum_{j=1}^{m_2} \frac{n_{i.}}{n_{.j}} \left(\frac{n_{ij}}{n_{i.}} - \frac{n_{.j}}{n}\right)^2$$

$$= \sum_{i=1}^{m_1} \sum_{j=1}^{m_2} \frac{1}{n_{i.}n_{.j}} \left(n_{ij} - \frac{n_{i.}n_{.j}}{n}\right)^2$$

Cette inertie mesure l'écart à l'indépendance.

Pourquoi la métrique du χ^2 ?

Pondération la pondération $n/n_{.j}$ permet de donner des importances comparables aux différentes « variables ».

Équivalence distributionnelle si deux colonnes j et j' de N ont le même profil, il est logique de les regrouper en une seule d'effectif $n_{ij} + n_{ij'}$; on a alors quand $n_{ij}/n_{i.} = n_{ij'}/n_{i.}$

$$\frac{n}{n_{.j}} \left[\frac{n_{ij}}{n_{i.}} - \frac{n_{.j}}{n}\right]^2 + \frac{n}{n_{.j'}} \left[\frac{n_{ij'}}{n_{i.}} - \frac{n_{.j'}}{n}\right]^2$$

$$= \frac{n}{n_{.j} + n_{.j'}} \left[\frac{n_{ij} + n_{ij'}}{n_{i.}} - \frac{n_{.j} + n_{.j'}}{n}\right]^2$$

La distance entre les profils-ligne est donc inchangée.

Propriétés diverses

Profils-colonnes on définit la distance entre deux profils-colonnes e_j et $e_{j'}$ comme

$$\|e_j - e_{j'}\|_{\chi^2} = \sum_{i=1}^{m_1} \frac{n}{n_{i.}} \left(\frac{n_{ij}}{n_{i.}} - \frac{n_{ij'}}{n_{i.}}\right)^2,$$

soit une métrique de matrice nD_1^{-1} .

Dualité On remarque que $I_{g_\ell} = I_{g_c}$, notée donc I_g .

Norme du vecteur centre de gravité g_ℓ

$$\|g_\ell\|_{\chi^2}^2 = \sum_{j=1}^{m_2} \frac{n}{n_{.j}} \left(\frac{n_{.j}}{n}\right)^2 = 1.$$

Forme alternative de l'inertie totale : comme $I_{\mathbf{v}} = I_{\mathbf{g}} + \|\mathbf{v} - \mathbf{g}\|_{\mathbf{M}}^2$, on a

$$I_{\mathbf{g}} = I_0 - \|\mathbf{g}\|_{\chi^2}^2 = \sum_{i=1}^{m_1} \sum_{j=1}^{m_2} \frac{n_{ij}^2}{n_i \cdot n_j} - 1.$$

ACP des deux nuages de profils

Comment ? deux possibilités en dualité exacte

	données	métrique	poids
Profils-lignes	$\mathbf{X} = \mathbf{D}_1^{-1}\mathbf{N}$	$\mathbf{M} = n\mathbf{D}_2^{-1}$	$\mathbf{D} = \frac{\mathbf{D}_1}{n}$
Profils-colonnes	$\mathbf{X} = \mathbf{D}_2^{-1}\mathbf{N}'$	$\mathbf{M} = n\mathbf{D}_1^{-1}$	$\mathbf{D} = \frac{\mathbf{D}_2}{n}$

Autres données

- Centre de gravité $\mathbf{g} = \mathbf{X}'\mathbf{D}\mathbf{1}$ (comme ACP), avec $\mathbf{M}\mathbf{g} = \mathbf{1}$ et $\mathbf{g}'\mathbf{M}\mathbf{g} = 1$
- Matrice de variance-covariance

$$\mathbf{V} = \mathbf{X}'\mathbf{D}\mathbf{X} - \mathbf{g}\mathbf{g}' = (\mathbf{X} - \mathbf{1}\mathbf{g}')'\mathbf{D}(\mathbf{X} - \mathbf{1}\mathbf{g}')$$

Non-nécessité du centrage

Propriété $\mathbf{V}\mathbf{M}$ et $\mathbf{X}'\mathbf{D}\mathbf{X}\mathbf{M}$ ont les mêmes vecteurs propres :

- d'une part \mathbf{g} (associé aux valeurs respectives 0 et 1)
- d'autre part des \mathbf{u} tels que $\mathbf{X}'\mathbf{D}\mathbf{X}\mathbf{M}\mathbf{u} = \mathbf{V}\mathbf{M}\mathbf{u} = \lambda\mathbf{u}$.

Preuve \mathbf{g} satisfait $\mathbf{V}\mathbf{M}\mathbf{g} = \mathbf{0}$ car $\mathbf{M}\mathbf{g} = \mathbf{1}$:

$$\mathbf{V}\mathbf{M}\mathbf{g} = (\mathbf{X} - \mathbf{1}\mathbf{g}')'\mathbf{D}(\mathbf{X} - \mathbf{1}\mathbf{g}')\mathbf{1} = \mathbf{0}.$$

De même

$$\mathbf{X}'\mathbf{D}\mathbf{X}\mathbf{M}\mathbf{g} = \mathbf{V}\mathbf{M}\mathbf{g} + \mathbf{g}\mathbf{g}'\mathbf{M}\mathbf{g} = \mathbf{0} + \mathbf{g} = \mathbf{g}.$$

Les autres vecteurs propres de $\mathbf{V}\mathbf{M}$ sont orthogonaux à \mathbf{g} ($\mathbf{g}'\mathbf{M}\mathbf{u} = 0$) et

$$\mathbf{X}'\mathbf{D}\mathbf{X}\mathbf{M}\mathbf{u} = \mathbf{V}\mathbf{M}\mathbf{u} + \mathbf{g}\mathbf{g}'\mathbf{M}\mathbf{u} = \mathbf{V}\mathbf{M}\mathbf{u} = \lambda\mathbf{u}.$$

Approche On effectue donc une ACP non centrée et on élimine la valeur propre 1 associée à l'axe principal \mathbf{g}

Calcul de l'ACP (profils-lignes)

Facteurs principaux ils sont vecteurs propres de

$$\mathbf{M}\mathbf{X}'\mathbf{D}\mathbf{X} = (n\mathbf{D}_2^{-1})(\mathbf{D}_1^{-1}\mathbf{N})'\frac{\mathbf{D}_1}{n}(\mathbf{D}_1^{-1}\mathbf{N}) = \mathbf{D}_2^{-1}\mathbf{N}'\mathbf{D}_1^{-1}\mathbf{N}.$$

et on a donc pour chaque axe principal k

$$\mathbf{D}_2^{-1}\mathbf{N}'\mathbf{D}_1^{-1}\mathbf{N}\mathbf{u}_k = \lambda_k\mathbf{u}_k$$

Composantes principales la composante principale associée au facteur \mathbf{u}_k est $\mathbf{a}_k = \mathbf{X}\mathbf{u}_k = \mathbf{D}_1^{-1}\mathbf{N}\mathbf{u}_k$; elle est vecteur propre de la matrice $\mathbf{D}_1^{-1}\mathbf{N}\mathbf{D}_2^{-1}\mathbf{N}'$ car

$$\begin{aligned} \mathbf{D}_1^{-1}\mathbf{N}\mathbf{D}_2^{-1}\mathbf{N}'\mathbf{a}_k &= \mathbf{D}_1^{-1}\mathbf{N}\mathbf{D}_2^{-1}\mathbf{N}'\mathbf{D}_1^{-1}\mathbf{N}\mathbf{u}_k \\ &= \lambda_k\mathbf{D}_1^{-1}\mathbf{N}\mathbf{u}_k = \lambda_k\mathbf{a}_k \end{aligned}$$

Analyse des profils-colonnes on échange les indices 1 et 2 et on transpose \mathbf{N} .

Comparaison lignes-colonnes

	ACP profils-lignes	ACP profils-colonnes
Facteurs principaux	Vecteurs propres de $\mathbf{D}_2^{-1}\mathbf{N}'\mathbf{D}_1^{-1}\mathbf{N}$	Vecteurs propres de $\mathbf{D}_1^{-1}\mathbf{N}\mathbf{D}_2^{-1}\mathbf{N}'$
Composantes principales	Vecteurs propres de $\mathbf{D}_1^{-1}\mathbf{N}\mathbf{D}_2^{-1}\mathbf{N}'$ normalisés par $\text{var } \mathbf{a}_k = \mathbf{a}_k'\frac{\mathbf{D}_1}{n}\mathbf{a}_k = \lambda_k$	Vecteurs propres de $\mathbf{D}_2^{-1}\mathbf{N}'\mathbf{D}_1^{-1}\mathbf{N}$ normalisés par $\text{var } \mathbf{b}_k = \mathbf{b}_k'\frac{\mathbf{D}_2}{n}\mathbf{b}_k = \lambda_k$

Comparaison les deux analyses conduisent aux mêmes valeurs propres et les facteurs principaux de l'une sont les composantes principales de l'autre (à un facteur près).

Partie IV. Aspects pratiques

Interprétation des résultats

Coordonnées des points Les coordonnées des points-lignes et points-colonnes s'obtiennent en cherchant les vecteurs propres des produits des deux tableaux de profils. Ce sont les grandeurs principales à obtenir.

Projection des nuages il est possible de projeter les deux nuages de points sur la même représentation. On justifiera plus tard le sens de cette représentation et son interprétation.

Cercle des corrélations il n'a aucun intérêt ici, puisque les véritables variables sont qualitatives.

(non) effet de taille comme les composantes variables sont centrées ($\sum_{i=1}^{m_1} n_i \cdot a_{ik} = \sum_{j=1}^{m_2} n_j \cdot b_{jk} = 0$), on sait que les coordonnées des \mathbf{a}_k et \mathbf{b}_k ne peuvent être toutes de même signe ; il n'y a donc jamais d'effet de « taille ».

Contributions à l'inertie

Contribution des profils-lignes On sait que $\lambda_k = \sum_{i=1}^{m_1} \frac{n_i}{n} (a_{ik})^2$, où a_{ik} est la coordonnée du profil-ligne i sur la k -ième composante principale de l'ACP sur les profils-lignes. On définit donc la contribution de la modalité i à l'axe principal k comme

$$\frac{n_i}{n} \cdot \frac{(a_{ik})^2}{\lambda_k}.$$

On considérera les modalités ayant l'influence la plus importante (typiquement $> \alpha n_i/n$, $\alpha = 2$ ou 3) comme constitutives des axes ; on regardera aussi le signe de la coordonnée.

Il n'y a pas ici de modalités sur-représentées, puisqu'on ne peut pas les retirer.

Contribution des profils-colonnes pour les mêmes raisons, la contribution de la modalité j de \mathcal{X}_2 à l'axe k est

$$\frac{n_j}{n} \cdot \frac{(b_{jk})^2}{\lambda_k}.$$

Qualité de la représentation

Profils-lignes l'AFC est une ACP, et on peut donc mesurer la qualité de la représentation de la modalité i (son profil-ligne) par un sous-espace factoriel. La qualité (le \cos^2 de l'angle entre le point et sa projection) s'écrit encore, pour le plan formé des k^* premiers axes :

$$\frac{\sum_{k=1}^{k^*} (a_{ik})^2}{\sum_{k=1}^{m_2} (a_{ik})^2}.$$

Comme pour l'ACP, > 0.8 signifie « bien représenté » et < 0.5 veut dire « mal représenté ». Les valeurs sont souvent données en pourcents.

Profils-colonne Le principe est le même, mais la formule devient, pour la modalité j :

$$\frac{\sum_{k=1}^{k^*} (b_{jk})^2}{\sum_{k=1}^{m_1} (b_{jk})^2}.$$

Formules de transition

But on cherche une relation entre les vecteurs \mathbf{a}_k et \mathbf{b}_k pour éviter de faire deux diagonalisation de matrice. Par exemple, si $m_1 < m_2$, on diagonalisera la matrice $\mathbf{D}_1^{-1} \mathbf{N} \mathbf{D}_2^{-1} \mathbf{N}'$.

Formules un calcul simple donne les formules suivantes

$$\mathbf{b}_k = \frac{1}{\sqrt{\lambda_k}} \mathbf{D}_2^{-1} \mathbf{N}' \mathbf{a}_k, \text{ soit } b_{jk} = \frac{1}{\sqrt{\lambda_k}} \sum_{i=1}^{m_1} \frac{n_{ij}}{n_{.j}} a_{ik},$$

$$\mathbf{a}_k = \frac{1}{\sqrt{\lambda_k}} \mathbf{D}_1^{-1} \mathbf{N} \mathbf{b}_k, \text{ soit } a_{ik} = \frac{1}{\sqrt{\lambda_k}} \sum_{j=1}^{m_2} \frac{n_{ij}}{n_{.i}} b_{jk}.$$

Méthode comme \mathbf{a}_k est (à une normalisation près) le facteur principal associé à \mathbf{b}_k , on sait que $\mathbf{b}_k = \alpha \mathbf{D}_2^{-1} \mathbf{N}' \mathbf{a}_k$. Pour déterminer α , il suffit d'écrire que $\mathbf{b}_k' \frac{\mathbf{D}_2}{n} \mathbf{b}_k = \lambda_k$.

Le χ^2 d'écart à l'indépendance

Utilité Il permet d'évaluer la dépendance entre les variables.

Définition c'est la grandeur suivante (parfois aussi notée χ^2 ou X^2)

$$d^2 = \sum_{i=1}^{m_1} \sum_{j=1}^{m_2} \frac{(n_{ij} - \frac{n_{i.} n_{.j}}{n})^2}{\frac{n_{i.} n_{.j}}{n}} = n \left[\sum_{i=1}^{m_1} \sum_{j=1}^{m_2} \frac{n_{ij}^2}{n_{i.} n_{.j}} - 1 \right].$$

$d^2 = 0 \iff$ les variables sont indépendantes.

Contribution au χ^2 c'est le terme

$$\frac{(n_{ij} - \frac{n_{i.} n_{.j}}{n})^2}{\frac{n_{i.} n_{.j}}{n}}$$

qui permet de mettre en évidence les associations significatives entre modalités de deux variables.

Borne supérieure comme $n_{ij} \leq n_{i.}$, on a

$$\sum_{i=1}^{m_1} \sum_{j=1}^{m_2} \frac{n_{ij}^2}{n_{i.} n_{.j}} \leq \sum_{i=1}^{m_1} \sum_{j=1}^{m_2} \frac{n_{ij}}{n_{.j}} = \sum_{j=1}^{m_2} \frac{\sum_{i=1}^{m_1} n_{ij}}{n_{.j}} = \sum_{j=1}^{m_2} \frac{n_{.j}}{n_{.j}} = m_2,$$

et donc $d^2 \leq n(m_2 - 1)$. On fait de même pour m_1 et

$$\varphi^2 = \frac{d^2}{n} \leq \min(m_1 - 1, m_2 - 1).$$

Dépendance fonctionnelle si $\varphi^2 = m_2 - 1$, alors pour chaque i soit $n_{ij} = n_{i.}$, soit $n_{ij} = 0$: il existe une unique case non nulle par ligne. \mathcal{X}_2 est donc fonctionnellement liée à \mathcal{X}_1 .

Dépendance inverse cette relation ne signifie pas que \mathcal{X}_1 est fonctionnellement liée à \mathcal{X}_2 , sauf si $m_1 = m_2$. On peut alors représenter le tableau comme une matrice diagonale.

Caractère significatif du χ^2

Problème à partir de quelle valeur de d^2 doit-on considérer que les variables \mathcal{X}_1 et \mathcal{X}_2 sont dépendantes ?

Méthode on suppose que \mathcal{X}_1 et \mathcal{X}_2 sont issus de tirages de deux variables aléatoires indépendantes. On peut alors montrer que d^2 est une réalisation d'une variable aléatoire D^2 qui suit une loi $\chi_{(m_1-1)(m_2-1)}^2$.

Définition Loi du khi-deux à ℓ degrés de libertés χ_ℓ^2 est la loi de la variable $\sum_{i=1}^{\ell} U_i^2$, où les U_i sont des variables gaussiennes réduites indépendantes.

Le test du χ^2 Ingrédients :

- on se fixe un risque d'erreur α (0.01 ou 0.05 en général)
- on calcule la valeur d_c^2 telle que $P(\chi_{(m_1-1)(m_2-1)}^2 > d_c^2) = \alpha$.
- Si $d^2 > d_c^2$ on considère que l'événement est trop improbable et que donc que l'hypothèse originale d'indépendance doit être rejetée.

$d^2 > d_c^2 \implies$ variables liées $d^2 < d_c^2 \implies$ pas de conclusion

Mode de calcul du χ^2

Calcul par table du χ^2 Traditionnellement, on trouvait ces valeurs dans une table précalculée pour $\ell \leq 30$.

- la ligne indique le nombre de degrés de liberté ℓ ;
- la colonne indique la probabilité cumulative $P(\chi_\ell^2 > d_c^2)$;
- la case de la table donne la valeur de d_c^2 .

Quand $\ell > 30$, on considère que $\sqrt{2}\chi_\ell^2 - \sqrt{2\ell - 1}$ est distribué comme une variable gaussienne centrée réduite $N(0, 1)$.

Utilisation de la p -value L'utilisation d'un logiciel de statistique permet en général de calculer directement la p -value $P(\chi_{(m_1-1)(m_2-1)}^2 > d^2)$ et rend inutile de se fixer un seuil d'erreur préalable.

Décomposition de l'inertie

Nombre de valeurs propres Comme $\mathbf{D}_1^{-1}\mathbf{N}\mathbf{D}_2^{-1}\mathbf{N}'$ est carrée de dimension m_1 et que $\mathbf{D}_2^{-1}\mathbf{N}'\mathbf{D}_1^{-1}\mathbf{N}$ est de dimension m_2 , le nombre de valeurs propres non nulles est $\min(m_1, m_2)$. Mais comme une des valeurs propres est 1 (associée à \mathbf{g}) et n'est pas intéressante :

Il y a au plus $\min(m_1 - 1, m_2 - 1)$ valeurs propres non nulles

φ^2 et valeurs propres l'inertie totale (et donc la somme des valeurs propres) est égale à φ^2 . Donc si $m_1 < m_2$, on obtient $\varphi^2 = \sum_{k=1}^{m_1-1} \lambda_k$.

Choix du nombre de valeurs propres On se contente souvent de regarder le premier plan principal car

- la règle de Kaiser $\lambda_k > \varphi^2/(m_1 - 1)$ s'applique mal ;
- la règle du coude reste valide, mais est subjective ;
- il existe un test sur de la part d'inertie non expliquée, mais il est un peu compliqué.

Partie V. Analyse des correspondances multiples

Analyse des correspondances multiples

But on veut étendre l'AFC au cas de $p \geq 2$ variables $\mathcal{X}_1, \mathcal{X}_2, \dots, \mathcal{X}_p$ à m_1, m_2, \dots, m_p modalités. Ceci est particulièrement utile pour l'exploration d'enquêtes où les questions sont à réponses multiples.

Problème l'analyse des correspondances utilise une table de contingence qui nécessite $p = 2$.

Méthode on cherche un moyen différent d'analyser $p > 2$ variables et on vérifie que les résultats sont comparables à l'AFC pour $p = 2$.

Les données

Données brutes chaque individu est décrit par les numéros des modalités qu'il possède pour chacune des p variables. Il n'est pas possible de faire des calculs sur ce tableau, où les valeurs sont arbitraires.

Tableau disjonctif on remplace la v -ième colonne par m_v colonnes d'indicatrices : on met un zéro dans chaque colonne, sauf celle correspondant à la modalité de l'individu i qui reçoit 1.

Exemple On interroge 6 personnes sur la couleur de leurs cheveux (CB, CC et CR pour blond, châtain et roux), la couleur de leurs yeux (YB, YV et YM pour bleu, vert et marron) et leur sexe (H/F). On a donc trois variables (avec respectivement 3, 3 et 2 modalités) mesurées sur 6 individus. Les tableaux brut (ci-dessous à gauche) sont équivalents aux tableaux disjonctifs (à droite).

$$\begin{pmatrix} \text{CB} \\ \text{CB} \\ \text{CC} \\ \text{CC} \\ \text{CR} \\ \text{CB} \end{pmatrix} \begin{pmatrix} \text{YB} \\ \text{YV} \\ \text{YB} \\ \text{YM} \\ \text{YV} \\ \text{YB} \end{pmatrix} \begin{pmatrix} \text{H} \\ \text{H} \\ \text{F} \\ \text{H} \\ \text{F} \\ \text{F} \end{pmatrix} \begin{pmatrix} 1 & 0 & 0 \\ 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \\ 1 & 0 & 0 \end{pmatrix} \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 1 & 0 & 0 \\ 0 & 0 & 1 \\ 0 & 1 & 0 \\ 1 & 0 & 0 \end{pmatrix} \begin{pmatrix} 1 & 0 \\ 1 & 0 \\ 0 & 1 \\ 1 & 0 \\ 0 & 1 \\ 0 & 1 \end{pmatrix}$$

Tableau disjonctif et tableau de contingence

Tableau disjonctif à la variable \mathcal{X}_v on associe le tableau disjonctif \mathbf{X}_v à n lignes et m_v colonnes.

Tableau de contingence on vérifie facilement que le tableau de contingence des variables \mathcal{X}_v et \mathcal{X}_w est donné par

$$\mathbf{N}_{vw} = \mathbf{X}'_v \mathbf{X}_w.$$

Effectifs marginaux la matrice diagonale des effectifs marginaux de la variable \mathcal{X}_v est donnée par

$$\mathbf{D}_v = \mathbf{X}'_v \mathbf{X}_v.$$

Exemple (suite) Table de contingence Cheveux/Yeux et matrice d'effectif marginaux de la couleur de cheveux

$$\mathbf{N}_{12} = \begin{pmatrix} 2 & 1 & 0 \\ 1 & 0 & 1 \\ 0 & 1 & 0 \end{pmatrix} \quad \mathbf{D}_1 = \begin{pmatrix} 3 & 0 & 0 \\ 0 & 2 & 0 \\ 0 & 0 & 1 \end{pmatrix}$$

Tableau disjonctif joint

Définition c'est la matrice $\mathbf{X} = (\mathbf{X}_1 | \mathbf{X}_2 | \dots | \mathbf{X}_p)$, qui possède n lignes et $m_1 + \dots + m_p$ colonnes. Chaque colonne représente une *catégorie*, c'est-à-dire une modalité d'une variable.

Exemple pour l'exemple de variables précédentes, on a le tableau disjonctif joint suivant

$$\mathbf{X} = \left(\begin{array}{ccc|ccc|ccc} 1 & 0 & 0 & 1 & 0 & 0 & 1 & 0 & 0 \\ 1 & 0 & 0 & 0 & 1 & 0 & 1 & 0 & 0 \\ 0 & 1 & 0 & 1 & 0 & 0 & 0 & 0 & 1 \\ 0 & 1 & 0 & 0 & 0 & 1 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 & 1 & 0 & 0 & 0 & 1 \\ 1 & 0 & 0 & 1 & 0 & 0 & 0 & 0 & 1 \end{array} \right)$$

Chaque ligne somme à 3. Les sommes de colonnes sont

$$(3 \quad 2 \quad 1 \mid 3 \quad 2 \quad 1 \mid 3 \quad 3)$$

Le tableau de Burt

Définition c'est un super-tableau de contingence des variables $\mathcal{X}_1, \dots, \mathcal{X}_p$, formé de tableaux de contingence et de matrices d'effectifs marginaux. :

$$\mathbf{B} = \mathbf{X}'\mathbf{X} = \begin{bmatrix} \mathbf{X}'_1\mathbf{X}_1 & \mathbf{X}'_1\mathbf{X}_2 & \dots & \mathbf{X}'_1\mathbf{X}_p \\ \mathbf{X}'_2\mathbf{X}_1 & \mathbf{X}'_2\mathbf{X}_2 & & \\ \vdots & & \ddots & \vdots \\ \mathbf{X}'_p\mathbf{X}_1 & \dots & & \mathbf{X}'_p\mathbf{X}_p \end{bmatrix} = \begin{bmatrix} \mathbf{D}_1 & \mathbf{N}_{12} & \dots & \mathbf{N}_{1p} \\ \mathbf{N}_{21} & \mathbf{D}_2 & & \\ \vdots & & \ddots & \vdots \\ \mathbf{N}_{p1} & \dots & & \mathbf{D}_p \end{bmatrix}$$

Exemple Toujours pour les mêmes variables

$$\mathbf{B} = \left(\begin{array}{ccc|ccc|cc} 3 & 0 & 0 & 2 & 1 & 0 & 2 & 1 \\ 0 & 2 & 0 & 1 & 0 & 1 & 1 & 1 \\ 0 & 0 & 1 & 0 & 1 & 0 & 0 & 1 \\ \hline 2 & 1 & 0 & 3 & 0 & 0 & 1 & 2 \\ 1 & 0 & 1 & 0 & 2 & 0 & 1 & 1 \\ 0 & 1 & 0 & 0 & 0 & 1 & 1 & 0 \\ \hline 2 & 1 & 0 & 1 & 1 & 1 & 3 & 0 \\ 1 & 1 & 1 & 2 & 1 & 0 & 0 & 3 \end{array} \right)$$

Partie VI. L'ACM : une AFC sur tableau disjonctif

Comment utiliser l'AFC pour analyser p variables

But on cherche à faire une représentation des $m_1 + \dots + m_p$ catégories comme points d'un espace de faible dimension.

Méthode on fait une AFC sur le tableau disjonctif joint $\mathbf{X} = (\mathbf{X}_1 | \mathbf{X}_2 | \dots | \mathbf{X}_p)$.

Les lignes la somme des éléments de chaque ligne de \mathbf{X} est égale à p . Le tableau des profils-lignes est donc $\frac{1}{p}\mathbf{X}$.

Les colonnes la somme des éléments de chaque colonne de \mathbf{X} est égale à l'effectif marginal de la catégorie correspondante. Le tableau des profils colonnes est donc $\mathbf{X}\mathbf{D}^{-1}$, où \mathbf{D} est la matrice diagonale par blocs

$$\mathbf{D} = \begin{pmatrix} \mathbf{D}_1 & & \mathbf{0} \\ & \ddots & \\ \mathbf{0} & & \mathbf{D}_p \end{pmatrix}$$

Les coordonnées factorielles des catégories

Notation On note $\mathbf{a}_k = (\mathbf{a}_{1k}, \dots, \mathbf{a}_{pk})'$ le vecteur à $m_1 + \dots + m_p$ composantes des coordonnées factorielles des catégories sur l'axe k .

Calcul de l'AFC sur \mathbf{X} comme la matrice des profils lignes est $\frac{1}{p}\mathbf{X}$ et celle des profils colonnes $\mathbf{X}\mathbf{D}^{-1}$, \mathbf{a}_k est vecteur propre de

$$(\mathbf{X}\mathbf{D}^{-1})' \frac{1}{p}\mathbf{X} = \frac{1}{p}\mathbf{D}^{-1}\mathbf{X}'\mathbf{X} = \frac{1}{p}\mathbf{D}^{-1}\mathbf{B}$$

et donc l'équation des coordonnées des catégories est

$$\frac{1}{p}\mathbf{D}^{-1}\mathbf{B}\mathbf{a}_k = \mu_k \mathbf{a}_k$$

avec la convention de normalisation $\frac{1}{np}\mathbf{a}_k' \mathbf{D}\mathbf{a}_k = \mu_k$.

Propriétés des valeurs propres

Valeur propres triviales La valeur propre 1 est associée (comme en AFC) à la composante $\mathbf{z}^0 = (1, \dots, 1)$ dans l'espace des individus. Les autres vecteurs propres lui sont orthogonaux, et donc de moyenne nulle.

Autres valeurs propres Si $n > \sum_{v=1}^p m_v$, le rang de \mathbf{X} est $\sum_{v=1}^p m_v - p + 1$ et le nombre de valeurs propres non trivialement égales à 0 ou 1 est

$$q = \sum_{v=1}^p m_v - p.$$

Somme La somme des valeurs propres non triviales est

$$\sum_{k=1}^q \mu_k = \text{Tr} \left(\frac{1}{p}\mathbf{D}^{-1}\mathbf{B} \right) - 1 = \frac{1}{p} \sum_{v=1}^p m_v - 1 = \frac{q}{p}$$

et leur moyenne vaut donc $1/p$.

Résolution dans le cas $p = 2$

On note \mathbf{a}_k (resp. \mathbf{b}_k) les m_1 premières (resp. m_2 dernières) coordonnées de la composante principale k et μ_k la valeur propre correspondante :

$$\frac{1}{2}\mathbf{D}^{-1}\mathbf{B} \begin{bmatrix} \mathbf{a}_k \\ \mathbf{b}_k \end{bmatrix} = \frac{1}{2} \begin{bmatrix} \mathbf{I}_{m_1} & \mathbf{D}_1^{-1}\mathbf{N} \\ \mathbf{D}_2^{-1}\mathbf{N}' & \mathbf{I}_{m_2} \end{bmatrix} \begin{bmatrix} \mathbf{a}_k \\ \mathbf{b}_k \end{bmatrix} = \mu_k \begin{bmatrix} \mathbf{a}_k \\ \mathbf{b}_k \end{bmatrix}.$$

On obtient les équations

$$\begin{cases} \mathbf{D}_1^{-1}\mathbf{N}\mathbf{b}_k = (2\mu_k - 1)\mathbf{a}_k \\ \mathbf{D}_2^{-1}\mathbf{N}'\mathbf{a}_k = (2\mu_k - 1)\mathbf{b}_k \end{cases},$$

et donc on retrouve les coordonnées des modalités de lignes et de colonnes dans l'AFC classique (avec $\lambda_k = (2\mu_k - 1)^2$) :

$$\begin{cases} \mathbf{D}_2^{-1}\mathbf{N}'\mathbf{D}_1^{-1}\mathbf{N}\mathbf{b}_k = (2\mu_k - 1)^2\mathbf{b}_k \\ \mathbf{D}_1^{-1}\mathbf{N}\mathbf{D}_2^{-1}\mathbf{N}'\mathbf{a}_k = (2\mu_k - 1)^2\mathbf{a}_k \end{cases}.$$

Différences ACM/AFC pour $p = 2$

Nombre de valeurs propres on a *a priori* $m_1 + m_2 - 2$ valeurs propres non nulles, ce qui est plus important que dans le cas classique. En particulier pour chaque λ_k , on a deux μ_k possibles

$$\begin{cases} \mu_k = \frac{1+\sqrt{\lambda_k}}{2} & \text{associée à } \begin{bmatrix} \mathbf{a}_k \\ \mathbf{b}_k \end{bmatrix} \\ \mu'_k = \frac{1-\sqrt{\lambda_k}}{2} & \text{associée à } \begin{bmatrix} \mathbf{a}_k \\ -\mathbf{b}_k \end{bmatrix} \end{cases}$$

On ne garde donc que les valeurs $\mu_k > \frac{1}{2}$. On peut montrer qu'il y en a $\min(m_1 - 1, m_2 - 1)$.

Inertie L'interprétation de la part d'inertie expliquée par les valeurs propres est maintenant très différente. En particulier les valeurs propres qui étaient très séparées dans l'AFC de \mathbf{N} le sont beaucoup moins dans celle de \mathbf{X} .

Partie VII. Aspects pratiques

Formules barycentriques

Les coordonnées des individus Soit \mathbf{c}_k le vecteur à n composantes des coordonnées des n individus sur l'axe factoriel associé à la valeur propre μ_k . D'après les résultats sur l'AFC, on a

$$\mathbf{c}_k = \frac{1}{\sqrt{\mu_k}} \frac{1}{p} \mathbf{X} \mathbf{a}_k \quad \text{et donc} \quad c_{ik} = \frac{1}{\sqrt{\mu_k}} \frac{1}{p} \sum_{j \text{ catégorie de } i} a_{jk}$$

Les seuls termes non nuls dans le calcul de $\mathbf{X} \mathbf{a}_k$ sont les coordonnées de la catégorie de chaque variable possédée par l'individu. Comme on est dans le cadre de l'AFC, la variance de \mathbf{c}_k est toujours $\text{var } \mathbf{c}_k = \frac{1}{n} \mathbf{c}'_k \mathbf{c}_k = \mu_k$.

Barycentre des catégories À $1/\sqrt{\mu_k}$ près, la coordonnée d'un individu est égale à la moyenne arithmétique simple des coordonnées des catégories auxquelles il appartient.

Les coordonnées des catégories On a de même la seconde formule

$$\mathbf{a}_k = \frac{1}{\sqrt{\mu_k}} \mathbf{D}^{-1} \mathbf{X}' \mathbf{c}_k \quad \text{c-à-d} \quad a_{jk} = \frac{1}{\sqrt{\mu_k}} \frac{1}{n_j} \sum_{i \text{ de catégorie } j} c_{ik}$$

Les seuls termes non nuls de $\mathbf{X}' \mathbf{c}_k$ sont les coordonnées des individus ayant une catégorie donnée. Là encore, $\text{var } \mathbf{a}_k = \mu_k$.

Barycentre des individus À $1/\sqrt{\mu_k}$ près, la coordonnée d'une catégorie est égale à la moyenne arithmétique des coordonnées des n_j individus de cette catégorie.

Barycentres et représentation

Représentation commune Les points représentatifs des catégories sont barycentres des groupes d'individus. On peut donc représenter individus et catégories dans un même plan factoriel.

Moyennes Comme \mathbf{c}_k est une variable de moyenne nulle, la formule de barycentre indique que pour chaque variable \mathcal{X}_i , les coordonnées de ses catégories (pondérés par les effectifs) sont de moyenne nulle. Aucun centrage n'est donc nécessaire

Échelle pour que les catégories se trouvent visuellement au barycentre des individus qui les représentent on peut remplacer \mathbf{a}_k par

$$\boldsymbol{\alpha}_k = \mathbf{D}^{-1} \mathbf{X}' \mathbf{c}_k = \sqrt{\mu_k} \mathbf{a}_k$$

Sélection de variables et axes

Sélection des variables on décide souvent de ne garder qu'un nombre réduit de variables actives et de garder les autres comme variables supplémentaires.

Sélection des axes

- règle courante : garder les axes tels que $\mu_k > 1/p$ (la moyenne des valeurs propres est $1/p$).
- les axes intéressants sont ceux que l'on peut interpréter, en regardant les contributions des variables actives et les valeurs-tests associées aux variables supplémentaires (définies plus tard).
- En pratique on se contente souvent d'interpréter le premier plan principal.

Inertie expliquée elle est moins intéressante qu'en ACP.

Catégories et axes factoriels

Catégorie Comme $\text{var } \mathbf{a}_k = \sum_j \frac{n_j}{np} (a_{jk})^2 = \mu_k$, la contribution de la catégorie j à l'axe k est

$$\frac{n_j}{np} \frac{(a_{jk})^2}{\mu_k},$$

intéressante si elle est supérieure au poids n_j/np (à un facteur près comme en ACP et AFC).

Variable la contribution totale de la variable \mathcal{X}_v à l'axe factoriel est

$$\frac{1}{\mu_k} \frac{1}{np} \sum_{j \text{ modalité de } \mathcal{X}_v} n_j (a_{jk})^2$$

Qualité de la représentation pour le sous-espace formé par les k^* premiers axes, la qualité de la représentation de la catégorie j est le cosinus carré habituel

$$\frac{\sum_{k=1}^{k^*} (a_{jk})^2}{\sum_{k=1}^q (a_{jk})^2}.$$

Individus et axes factoriels

La normalisation de \mathbf{c}_k est $\sum_{i=1}^n (c_{ik})^2 = n\mu_k$, où c_{ik} est la coordonnée de l'individu i sur l'axe factoriel k associé à la valeur propre μ_k .

Contribution d'un individu pour l'individu i , c'est

$$\frac{1}{n} \frac{(c_{ik})^2}{\mu_k}$$

Cette contribution est jugée en la comparant au poids $1/n$ comme en ACP et AFC.

Qualité de la représentation pour le sous-espace formé par les k^* premiers axes, la qualité de la représentation de l'individu i est

$$\frac{\sum_{k=1}^{k^*} (c_{ik})^2}{\sum_{k=1}^q (c_{ik})^2}.$$

Contribution à l'inertie totale

Soit $\mathbf{x}^j = (x_i^j)$ le vecteur colonne de \mathbf{X} correspondant à une catégorie j . On rappelle que l'inertie totale vaut

$$I_g = \sum_{j \in \text{catégories}} \frac{n_j}{np} d^2(\mathbf{z}^j, \mathbf{g}) = \frac{1}{p} \sum_{v=1}^p m_v - 1,$$

où la distance du profil-colonne j au centre de gravité des profils-colonnes $\mathbf{g} = \mathbf{1}/n$ est

$$\begin{aligned} d^2(\mathbf{z}^j, \mathbf{g}) &= \sum_{i=1}^n \frac{np}{p} \left(\frac{x_i^j}{n_j} - \frac{1}{n} \right)^2 = n \sum_{i=1}^n \left(\frac{x_i^j}{n_j^2} + \frac{1}{n^2} - \frac{2x_i^j}{nn_j} \right) \\ &= n \left(\frac{n_j}{n_j^2} + \frac{n}{n^2} - \frac{2n_j}{nn_j} \right) = \frac{n}{n_j} - 1 \end{aligned}$$

Contribution d'une catégorie La contribution absolue de la catégorie j à l'inertie est

$$\frac{n_j}{np} d^2(\mathbf{z}^j, \mathbf{g}) = \frac{1}{p} \left(1 - \frac{n_j}{n} \right),$$

qui est une fonction décroissante de l'effectif. Il faut donc éviter les catégories d'effectif trop faible, qui d'ailleurs se retrouveront dans les premiers axes

Contribution d'une variable La contribution de la variable \mathcal{X}_v est

$$\sum_{j \text{ modalité de } \mathcal{X}_v} \frac{1}{p} \left(1 - \frac{n_j}{n} \right) = \frac{m_v - 1}{p}$$

Elle est d'autant plus grande que le nombre de modalités de \mathcal{X}_i est élevé. Il faut donc éviter les disparités trop grandes entre le nombre de modalités (quand on a le choix du découpage...)

Correspondances multiples et ACP non linéaire

Problème l'ACP vise à trouver une combinaison *linéaire* $u_1 \mathbf{x}^1 + \dots + u_p \mathbf{x}^p$ des variables qui soit de variance maximale. Si les relations entre variables ne sont pas linéaires, l'ACP échoue à extraire des données intéressantes.

Extension non-linéaire on cherche des fonctions ϕ^1, \dots, ϕ^p qui maximisent la variance

$$\text{var}(\phi^1(\mathbf{x}^1) + \dots + \phi^p(\mathbf{x}^p))$$

Fonctions en escalier On peut prendre des fonctions en escalier : on découpe l'intervalle de variation de \mathbf{x}^j en m_j classes et on se donne un vecteur $\mathbf{a}_j = (a_{j1}, \dots, a_{jm_j})$ de poids. Alors $\phi^j(x) = a_{j\ell}$ si x est dans la ℓ -ième classe.

Discretisation des variables on définit le tableau disjonctif \mathbf{X}_j indiquant quelle modalité (classe) de la variable j est prise par chaque individu. Alors

$$\phi^j(\mathbf{x}^j) = \mathbf{X}_j \mathbf{a}_j.$$

Reformulation de l'ACP non-linéaire on cherche le vecteur $\mathbf{a} = (\mathbf{a}_1, \dots, \mathbf{a}_p)$ qui maximise la variance

$$\text{var}(\mathbf{X}_1 \mathbf{a}_1 + \dots + \mathbf{X}_p \mathbf{a}_p) = \text{var}(\mathbf{X} \mathbf{a})$$

La solution est la première composante de l'ACM du tableau disjonctif joint \mathbf{X} .

Conclusion le découpage en classes des variables numériques permet d'obtenir une analyse non linéaire des données. Elle n'est possible que si on a suffisamment d'observations par classe.

Partie VIII. Interprétation externe

Les variables supplémentaires

Leur usage est très courant en analyse des correspondances multiples.

Variables quantitatives on calcule « à la main » leur corrélation avec les axes factoriels et on les place sur un cercle de corrélations. Si $\hat{\mathbf{z}}$ est une version centrée-réduite de la variable, alors

$$\text{cor}(\hat{\mathbf{z}}, \mathbf{c}_k) = \frac{1}{\sqrt{\mu_k}} \frac{1}{n} \sum_{i=1}^n \hat{z}_i c_{ik}$$

On peut aussi les découper en classes et les traiter comme des variables qualitatives.

Variables qualitatives on calcule directement les coordonnées de leurs modalités en utilisant la formule de barycentre des individus : la coordonnée la catégorie supplémentaire \hat{j} sur l'axe principal k est

$$a_{jk} = \frac{1}{\sqrt{\mu_k}} \frac{1}{n_j} \sum_{i \text{ de catégorie } \hat{j}} c_{ik}$$

Valeurs-test pour les variables supplémentaires qualitatives

But on cherche à savoir si une catégorie \hat{j} d'effectif n_j et de coordonnée a_{jk} sur cet axe est liée à cet axe.

Idee du calcul si les n_j individus d'une catégorie étaient pris au hasard, la moyenne de leurs coordonnées serait une variable aléatoire centrée (les \mathbf{c} sont de moyenne nulle) et de variance $\frac{\mu_k}{n_j} \frac{n - n_j}{n - 1}$. De plus, la moyenne des coordonnées est égale à $\sqrt{\mu_k} a_{jk}$.

Valeur-test c'est la version centrée et réduite de la moyenne des coordonnées

$$a_{jk} \sqrt{n_j} \sqrt{\frac{n - 1}{n - n_j}}.$$

Quand n_j et $n - n_j$ sont assez grand (en général > 30), elle est significative si elle est supérieure à 2 ou 3 en valeur absolue. On ne doit pas l'utiliser sur les variables actives.

Partie IX. Récapitulatif

Notations AFC

Notation	taille	description
m_1 et m_2	entiers	nombre de modalités des variables 1 et 2
$\mathbf{N} = (n_{ij})$	$m_1 \times m_2$	table de contingence
$\mathbf{D}_1 = \text{diag}(n_{i.})$	$m_1 \times m_1$	effectifs (marges) de lignes
$\mathbf{D}_2 = \text{diag}(n_{.j})$	$m_2 \times m_2$	effectifs (marges) de colonnes
$n_{ij}/n_{i.}$ et $n_{ij}/n_{.j}$ $d^2 = n\varphi^2$	$m_1 \times m_2$ réel > 0	profils lignes et colonnes χ^2 d'écart à l'indépendance
\mathbf{a}_k et \mathbf{b}_k	m_1 et m_2	coordonnées des lignes et colonnes sur l'axe k
λ_k	réel > 0	Valeur propre associée à l'axe k

Notations ACM

Notation	taille	description
m_1, \dots, m_p	entiers	nombres de modalités des variables
$\mathbf{X} = (x_i^j)$	$n \times (\text{nb. cat.})$	tableau disjonctif
$\mathbf{N}_{k\ell}$	$m_k \times m_\ell$	table de contingence des variables k et ℓ
$\mathbf{D} = \text{diag}(n_i)$	nb. cat.	effectifs (marges) des catégories
\mathbf{B}	$(\text{nb. cat.})^2$	matrice de Burt
μ_k	réel > 0	Valeur propre associée à l'axe k
\mathbf{a}_k	nb. cat.	coordonnées des catégories sur l'axe k
\mathbf{c}_k	n	coordonnées des individus sur l'axe k

nb. cat. = $m_1 + \dots + m_p$.

Points communs entre AFC et ACM

But	décrire les liaisons entre plusieurs variables qualitatives
Cas $p = 2$	les coordonnées des modalités sont les mêmes pour les deux analyses
Représentation	toutes les modalités peuvent être représentées sur le même diagramme
Contribution d'une modalité à un axe	$\text{poids} \times \frac{(\text{coordonnée})^2}{\text{valeur propre}}$
Qualité de la représentation d'une modalité par un sous espace	$\cos^2 \theta = \frac{\sum_{\text{axes du sous esp.}} (\text{coord sur l'axe})^2}{\sum_{\text{tous les axes}} (\text{coord sur l'axe})^2}$

Différences entre AFC et ACM

	AFC	ACM
Individus	non	oui
Données	tableau de contingence profils lignes/colonnes	tableau disjonctif tableau de Burt
Poids d'une modalité	$\frac{n_{i.}}{n}$ (profil-ligne) $\frac{n_{.j}}{n}$ (profil-colonne)	$\frac{n_j}{np}$
Nb de val. propres	$\min(m_1 - 1, m_2 - 1)$	$\sum_{v=1}^p m_v - p$
Axes à conserver	pas de règle Kaiser ; peut-être part d'inertie.	$\mu > \frac{1}{p}$
Variables supplémentaires	pas vraiment de sens	qualitatives et quantitatives