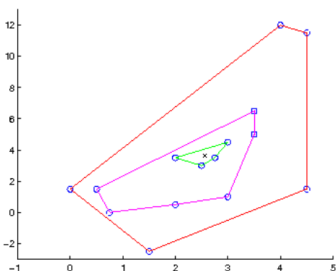


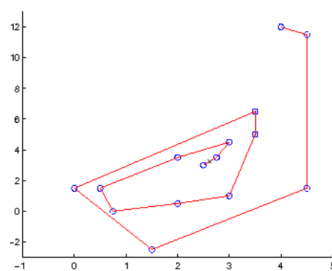
Description statistique des données

M8 - Chapitre 2

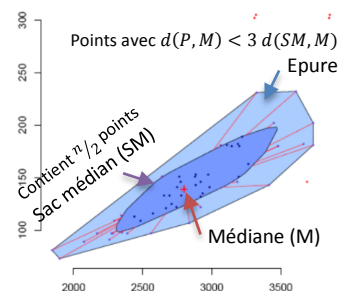
• Variables :	$X = \{a_1, \dots, a_l\}$ $Y = \{b_1, \dots, b_k\}$	Tableau de contingence : <table border="1"> <thead> <tr> <th>x \ y</th> <th>b_1</th> <th>...</th> <th>b_j</th> <th>...</th> <th>b_k</th> <th>marge</th> </tr> </thead> <tbody> <tr> <td>a_1</td> <td>n_{11}</td> <td>...</td> <td>n_{1j}</td> <td>...</td> <td>n_{1k}</td> <td>$n_{1\bullet}$</td> </tr> <tr> <td>\vdots</td> <td>\vdots</td> <td></td> <td>\vdots</td> <td></td> <td>\vdots</td> <td>\vdots</td> </tr> <tr> <td>a_i</td> <td>n_{i1}</td> <td>...</td> <td>n_{ij}</td> <td>...</td> <td>n_{ik}</td> <td>$n_{i\bullet}$</td> </tr> <tr> <td>\vdots</td> <td>\vdots</td> <td></td> <td>\vdots</td> <td></td> <td>\vdots</td> <td>\vdots</td> </tr> <tr> <td>a_l</td> <td>n_{l1}</td> <td>...</td> <td>n_{lj}</td> <td>...</td> <td>n_{lk}</td> <td>$n_{l\bullet}$</td> </tr> <tr> <td>marge</td> <td>$n_{\bullet 1}$</td> <td>...</td> <td>$n_{\bullet j}$</td> <td>...</td> <td>$n_{\bullet k}$</td> <td>$n_{\bullet\bullet} = n$</td> </tr> </tbody> </table>	x \ y	b_1	...	b_j	...	b_k	marge	a_1	n_{11}	...	n_{1j}	...	n_{1k}	$n_{1\bullet}$	\vdots	\vdots		\vdots		\vdots	\vdots	a_i	n_{i1}	...	n_{ij}	...	n_{ik}	$n_{i\bullet}$	\vdots	\vdots		\vdots		\vdots	\vdots	a_l	n_{l1}	...	n_{lj}	...	n_{lk}	$n_{l\bullet}$	marge	$n_{\bullet 1}$...	$n_{\bullet j}$...	$n_{\bullet k}$	$n_{\bullet\bullet} = n$
x \ y	b_1		...	b_j	...	b_k	marge																																												
a_1	n_{11}		...	n_{1j}	...	n_{1k}	$n_{1\bullet}$																																												
\vdots	\vdots			\vdots		\vdots	\vdots																																												
a_i	n_{i1}		...	n_{ij}	...	n_{ik}	$n_{i\bullet}$																																												
\vdots	\vdots		\vdots		\vdots	\vdots																																													
a_l	n_{l1}	...	n_{lj}	...	n_{lk}	$n_{l\bullet}$																																													
marge	$n_{\bullet 1}$...	$n_{\bullet j}$...	$n_{\bullet k}$	$n_{\bullet\bullet} = n$																																													
• Échantillon :	$S_n = \{(x_1, y_1), \dots, (x_n, y_n)\}$ $(x_i, y_i) \in X \times Y$																																																		
• Lignes :	$n_{i\bullet} = \sum_{j=1}^k n_{ij}$																																																		
• Colonnes :	$n_{\bullet j} = \sum_{i=1}^l n_{ij}$																																																		
• Total :	$n_{\bullet\bullet} = \sum_{i=1}^l \sum_{j=1}^k n_{ij}$																																																		
• Fréquence :	$f_{ij} = \frac{n_{ij}}{n} = \hat{P}_{ij} = \hat{P}(X = a_i \cap Y = b_j)$																																																		
• Profil ligne :	$l_{ij} = \frac{n_{ij}}{n_{i\bullet}} = \hat{P}(Y = b_j X = a_i)$																																																		
• Profil colonne :	$c_{ij} = \frac{n_{ij}}{n_{\bullet j}} = \hat{P}(X = a_i Y = b_j)$																																																		
• Distribution marginale :	$\frac{n_{i\bullet}}{n} = \hat{P}_i \quad \frac{n_{\bullet j}}{n} = \hat{P}_j$																																																		
• Moyenne :	$\min_{(\bar{x}, \bar{y}) \in \mathbb{R}^2} \sum_{i=1}^n (x_i - \bar{x})^2 + (y_i - \bar{y})^2$																																																		
• Médiane :	$\min_{(M_x, M_y) \in \mathbb{R}^2} \sum_{i=1}^n x_i - M_x + y_i - M_y $																																																		
• Médiane euclidienne :	$\min_{(M_x, M_y) \in \mathbb{R}^2} \sum_{i=1}^n \sqrt{ x_i - M_x ^2 + y_i - M_y ^2}$																																																		
• Information mutuelle : (X, Y discrètes)	$I(X, Y) = \sum_{i \in \Omega_X} \sum_{j \in \Omega_Y} \hat{P}_{ij} \log \left(\frac{\hat{P}_{ij}}{\hat{P}_i \hat{P}_j} \right)$ $X \text{ indep } Y \Rightarrow I(X, Y) = 0$																																																		
• Dépendance : (X, Y quantitatives)	<div style="display: flex; justify-content: space-around;"> <div style="border: 1px solid black; padding: 5px;"> $s_{XY} = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})$ $= \frac{1}{n} \sum_{i=1}^n (x_i y_i - \bar{x} \bar{y})$ <p style="text-align: center; font-size: small;">Covariance</p> </div> <div style="border: 1px solid black; padding: 5px;"> $s_X^2 = s_{XX} = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2$ $s_Y^2 = s_{YY} = \frac{1}{n} \sum_{i=1}^n (y_i - \bar{y})^2$ <p style="text-align: center; font-size: small;">Variance</p> </div> <div style="border: 1px solid black; padding: 5px; font-size: small;"> $r = \frac{s_{XY}}{s_X s_Y}$ <p>Corrélation $\in [-1; 1]$ $r \approx 0 \Rightarrow x \text{ indep } y$ $r \approx 1 \Rightarrow x \nearrow y \nearrow$ $r \approx -1 \Rightarrow x \nearrow y \searrow$</p> </div> </div>																																																		
• Dépendance : (X, Y qualitatives) Distance du χ^2	$D_{\chi^2} = \sum_{i \in \Omega_X} \sum_{j \in \Omega_Y} \frac{\left(\frac{n_{ij} - \frac{n_{i\bullet} n_{\bullet j}}{n}}{n} \right)^2}{\frac{n_{i\bullet} n_{\bullet j}}{n}} = n \sum_{i \in \Omega_X} \sum_{j \in \Omega_Y} \frac{\left(\frac{\hat{P}_{ij}}{\hat{P}_i \hat{P}_j} - 1 \right)^2}{\hat{P}_i \hat{P}_j}$ <small>eff obs / eff th</small>																																																		
• Dépendance : (X qual, Y quant) Coef. de détermination	$S_{Y/X}^2 = \frac{\sum_{j \in \Omega_Y} n_j (\bar{y}_j - \bar{y})^2}{\sum_{i=1}^n (y_i - \bar{y})^2}$																																																		



Médiane de Tukey



Médiane de Jarvis



Sac médian

Description statistique des données

M8 – Chapitre 2

I. Matrice de covariance et matrice de corrélation

$$X_{C_{i,j}} = X_{i,j} - \bar{X}_j$$

Matrice centrée

$$X_{R_{i,j}} = \frac{X_{i,j} - \bar{X}_j}{S_j}$$

Matrice centrée réduite

$$\Sigma = \frac{X_C^T X_C}{n}$$

Matrice de covariance

$$\rho = \frac{X_R^T X_R}{n}$$

Matrice de corrélation

II. L'ACP

1. Le problème

On veut reconstruire au mieux X par le produit uv^T , c'est-à-dire minimiser :

$$\min_{u,v} J(u,v) = \min_{u,v} \|X - uv^T\|_F^2 = \min_{u,v} -2(Xv)^T u + \|u\|^2 + \|v\|^2 = \min_{u,v} \sum_i^n \sum_j^p (x_{ij} - u_i v_j)^2$$

2. La solution

$$\begin{cases} \nabla_u J(u) = -2Xv + 2\|v\|^2 u = 0 & \Leftrightarrow Xv = \|v\|^2 u \\ \nabla_v J(v) = -2X^T u + 2\|u\|^2 v = 0 & \Leftrightarrow -X^T u = \|u\|^2 v \end{cases} \Leftrightarrow X^T X v = \frac{\|u\|^2 \|v\|^2}{\lambda} v$$

La solution est un vecteur propre v de la matrice $X^T X$. C'est celle qui a la plus grande valeur propre car à l'optimum on a $J(u,v) = -\lambda$.

3. Remarques sur les valeurs et vecteurs propres

- Les valeurs propres sont ordonnées : $\lambda_1 \geq \lambda_2$
- Les vecteurs propres sont normés : $v_i^T v_i = 1$
- Les vecteurs propres sont orthogonaux : $v_i^T v_j = 0$

4. L'ACP

$$[v, \lambda] = \text{eig}(X) \quad v : \text{vecteurs propres} \quad \text{diag}(\lambda) : \text{valeurs propres}$$

Les composantes principales : $u = Xv$

Corrélations entre variables observées normalisées X_n (centrées réduites) et les composantes principales u :

$$\text{cor}(X_R, u_i) = \frac{X_R^T u_i}{\sqrt{n} \|u_i\|} = \frac{\|u_i\|}{\sqrt{n}} v = \frac{\sqrt{\lambda_i}}{\sqrt{n}} v$$

Rôle des variables :

$$V_n = \sqrt{\frac{\lambda}{n}} V$$