

Partie 1 : TD

Question 1

- 1) La droite de projection orthogonale a pour but de minimiser la somme des carrés des distances M_{imi} , $\frac{1}{n} \sum_{i=1}^m M_{imi}^2$.

Cherchons à retrouver pourquoi cette droite D doit passer par le point moyen du nuage.

En débutant avec le calcul de l'inertie d'une variable statistique par rapport à un point A .
L'inertie est définie comme la moyenne du carré des distances de la variable X au point A .

$$\text{On a : } I_T = \sum_{i=1}^m \frac{1}{n} d^2(A, X_i) = \sum_{i=1}^m \frac{1}{n} \left[\sum_{j=1}^p (a_j - x_{ij})^2 \right] = \frac{1}{n} \sum_{i=1}^m \sum_{j=1}^p (a_j - x_{ij})^2$$

Or l'inertie I_T est dérivable puisqu'elle s'écrit comme une somme d'éléments eux-mêmes dérivables.

$$\text{Ainsi : } I_T' = \frac{1}{n} \sum_{i=1}^m 2 \sum_{j=1}^p (a_j - x_{ij}) = \frac{1}{n} \sum_{i=1}^m 2 \sum_{j=1}^p (a_j - x_{ij}) = \frac{1}{n} \sum_{i=1}^m 2 a_j \sum_{j=1}^p (a_j - x_{ij})$$

On cherche à résoudre $I_T' = 0 \Leftrightarrow a_j = \frac{1}{n} \sum_{i=1}^m x_{ij} = g$ avec g , le centre de gravité du nuage de point.

Ainsi I_T' s'annule en g , qui est donc un minimum global de I_T .

On peut donc en conclure que D doit passer par le point moyen du nuage

Question 2 : Etude de la projection associée à l'ACP en 2D

2.1 Posons : $\mathcal{Y} : \mathbb{R}^2 \mapsto \mathbb{R}^{2c}$

$(u_1, u_2) \mapsto [(A - \lambda I_2)u_1; (A - \lambda I_2)u_2]$ ainsi \mathcal{Y} définit le noyau de l'endomorphisme

On à retrouver les conditions pour lesquelles l'application linéaire du noyau de l'endomorphisme, soit \mathcal{Y} ne se réduise pas à 0.

Comme $\varphi: \mathbb{R}^2 \mapsto \mathbb{R}^{2e}$, nous sommes en dimension finie. Alors il suffit de montrer que si $\varphi(u_1, u_2) = (0; 0)$ n'implique pas que (u_1, u_2) soit égal au vecteur nul, c'est à dire que $(u_1, u_2) = (0, 0)$

$$\text{Ainsi: } \varphi(u_1, u_2) = (0; 0) \Rightarrow \begin{cases} (A - \lambda I_2) u_1 = 0 \\ (A - \lambda I_2) u_2 = 0 \end{cases} \Rightarrow \begin{cases} A u_1 - \lambda u_1 = 0 \\ A u_2 - \lambda u_2 = 0 \end{cases} \Rightarrow \begin{cases} \lambda_1 = \lambda_1 u_1 \\ \lambda_2 = \lambda_2 u_2 \end{cases}$$

1) après l'écriture précédente on a $(u_1, u_2) = (1, 1)$ ainsi nous avons un contre exemple.

2) l'application n'est donc pas injective, donc bijective dans certaines conditions.

la condition: $\text{Ker}(A - \lambda I_2) \neq 0 \Leftrightarrow \det(A - \lambda I_2) = 0$

2.2 L'équation du 2^{me} degré dont les valeurs λ_1 et λ_2 sont solutions est: $\lambda^2 - \text{tr}(A)\lambda + \det(A) = 0$
où A est la matrice de variance/covariance. $A = \begin{pmatrix} \text{var}(X) & \text{cov}(X, Y) \\ \text{cov}(Y, X) & \text{var}(Y) \end{pmatrix}$

donc $\text{tr}(A) = V(X) + V(Y) = \underline{S^2(X) + S^2(Y)}$

$\det(A) = V(X) \cdot V(Y) - \text{Cov}^2(X, Y) = \underline{S^2(X)S^2(Y) - \text{cov}^2(X, Y)}$

2.3 Calculons le discriminant du polynôme.

$$\begin{aligned} \text{ainsi: } \Delta &= (V(X) + V(Y))^2 - 4[V(X)V(Y) - \text{cov}^2(X, Y)] \\ &= V^2(X) + 2V(X)V(Y) + V^2(Y) - 4V(X)V(Y) + 4\text{cov}^2(X, Y) \\ &= V^2(X) + V^2(Y) - 2V(X)V(Y) + 4\text{cov}^2(X, Y) \\ &= [V(X) - V(Y)]^2 + 4\text{cov}^2(X, Y) \end{aligned}$$

ainsi $[V(X) - V(Y)]^2 > 0$ et $4\text{cov}^2(X, Y) > 0$

donc $\Delta > 0$, ainsi il y a deux racines $(\lambda_1, \lambda_2) \in \mathbb{R}^2$

ainsi $\lambda_1 = \frac{1}{2} [v(x) + v(y) + \sqrt{(v(x) - v(y))^2 + 4 \text{Cov}^2(x, y)}]$

$$\lambda_2 = \frac{1}{2} [v(x) + v(y) - \sqrt{(v(x) - v(y))^2 + 4 \text{Cov}^2(x, y)}]$$

2.4 (u_1, u_2) une base orthonormée; $Au_1 = \lambda_1$ et $Au_2 = \lambda_2$.

et le vecteur propre U associé à la valeur propre λ , avec $U = \begin{pmatrix} S^2(Y) - \lambda \\ -\text{Cov}(X, Y) \end{pmatrix}$

et U vecteur propre associé à λ , de la matrice à xxi $AU = \lambda U$

Ainsi: $AU = \begin{pmatrix} v(x) & \text{cov}(x, y) \\ \text{cov}(x, y) & v(y) \end{pmatrix} \begin{pmatrix} S^2(Y) - \lambda \\ -\text{Cov}(X, Y) \end{pmatrix}$

$$\Rightarrow \begin{pmatrix} v(x)(S^2(Y) - \lambda) - \text{cov}^2(x, y) \\ \text{cov}(y, x)(S^2(Y) - \lambda) - v(y)\text{cov}(x, y) \end{pmatrix}$$

$$\Rightarrow \begin{pmatrix} S^2(x)S^2(y) - \lambda S^2(x) - \text{cov}^2(x, y) \\ \text{cov}(x, y)(S^2(y) - \lambda) - S^2(y)\text{cov}(x, y) \end{pmatrix}$$

d'après
q.2.2

$$\Rightarrow \begin{pmatrix} \det(A) - \lambda S^2(x) \\ -\lambda \text{Cov}(x, y) \end{pmatrix} \quad \text{et} \quad \lambda^2 - (S^2(x) + S^2(y))\lambda + \det(A) = 0$$

$$\Rightarrow \begin{pmatrix} -\lambda^2 + (S^2(x) + S^2(y))\lambda - \lambda S^2(x) \\ -\lambda \text{Cov}(x, y) \end{pmatrix}$$

$$\Rightarrow \begin{pmatrix} -\lambda^2 + \lambda S^2(x) \\ -\lambda \text{Cov}(x, y) \end{pmatrix}$$

$$\Rightarrow \lambda \begin{pmatrix} -\lambda + S^2(X) \\ -\text{Cov}(X, Y) \end{pmatrix}$$

$$\Rightarrow \lambda U$$

Ainsi, nous avons bien $AU = \lambda U$, avec U un vecteur propre relatif à la valeur propre λ

2.5 $(u_1; u_2)$ étant une BON, $\|u_1\|^2 = 1$, $\|u_2\|^2 = 1$

$$\text{et par ailleurs } \|U\| = \sqrt{\|U\|^2} = \sqrt{{}^t U U}$$

$$\text{avec } {}^t U U = (S^2(Y) - \lambda \quad -\text{Cov}(X, Y)) \begin{pmatrix} S^2(Y) - \lambda \\ -\text{Cov}(X, Y) \end{pmatrix}$$

$$= \begin{pmatrix} (S^2(Y) - \lambda)^2 + \text{Cov}^2(X, Y) \end{pmatrix}$$

$$\text{ainsi } \|U\| = \sqrt{(S^2(Y) - \lambda)^2 + \text{Cov}^2(X, Y)}$$

2.6 le vecteur propre relatif à la valeur propre λ :

$$U_{\text{normé}} = \begin{pmatrix} \frac{S^2(Y) - \lambda}{\sqrt{(S^2(Y) - \lambda)^2 + \text{Cov}^2(X, Y)}} \\ \frac{-\text{Cov}(X, Y)}{\sqrt{(S^2(Y) - \lambda)^2 + \text{Cov}^2(X, Y)}} \end{pmatrix}$$

$$2.7 \text{ Ainsi, } u_1 = \frac{1}{\sqrt{(S^2(Y) - \lambda_1)^2 + \text{Cov}^2(X, Y)}} \begin{pmatrix} S^2(Y) - \lambda_1 \\ -\text{Cov}(X, Y) \end{pmatrix}$$

$$u_2 = \frac{1}{\sqrt{(S^2(Y) - \lambda_2)^2 + \text{Cov}^2(X, Y)}} \begin{pmatrix} S^2(Y) - \lambda_2 \\ -\text{Cov}(X, Y) \end{pmatrix}$$

TP n°1 _ACP

Question n°3.1.

On définit la matrice par $X = \begin{pmatrix} 0 & 1 & 0 & 1 & 1 & 0 \\ 0 & 1 & 1 & 0 & 1 & 0 \end{pmatrix}^T$ ainsi que la matrice des poids $D = \frac{1}{6} I_n$, puisqu'il y a 6 individus dans notre jeu de données.

Une rapide boucle sur les colonnes de X permet de nous donner la moyenne de chacune des variables, soit le vecteur $M = (0.5 \quad 0.5)$.

La matrice Z centrée est alors définie par $Z = \begin{pmatrix} -0.5 & 0.5 & -0.5 & 0.5 & 0.5 & -0.5 \\ -0.5 & 0.5 & 0.5 & -0.5 & 0.5 & -0.5 \end{pmatrix}^T$.

Finalement, en "re-biaisant" la fonction `cov()` sur R , on obtient la matrice de covariance $A = \frac{n-1}{n} \text{cov}(Z) = \begin{pmatrix} 0.25 & 0.8333 \\ 0.8333 & 0.25 \end{pmatrix}$ il vient donc directement $S^2(X) = S^2(Y) = 0.25$ et $\text{cov}(X, Y) = 0.83$.

Question n°3.2.

Pour la suite du TP, nous avons développé une fonction générique qui réalise les calculs nécessaires à l'ACP pour une matrice de taille $n \times p$.

Pseudo-code :

ACP(X , centrer = FALSE) :

 Récupérer les dimensions n et p de X

 Pour chaque colonne de X : $Z \leftarrow \text{colonne} - \text{moyenne}(\text{colonne})$

 Si (centrer == TRUE) :

$Z \leftarrow Z / \text{var}(\text{colonne})$

 Afficher Z

 Calculer la matrice de covariance A

 Afficher A

 Diagonaliser A

$U \leftarrow$ vecteurs propres

$\text{Lambda} \leftarrow$ valeurs propres

 Composantes $\leftarrow X \cdot U$

Retourner Z , A , U , Lambda , Composantes

Après exécution du code avec notre matrice X , on obtient $\lambda_1 = 0.33334$ et $\lambda_2 = 0.16667$.

On a donc bien l'égalité $\text{tr}(A) = \lambda_1 + \lambda_2 = 0.5$ ainsi que $\det(A) = \lambda_1 \lambda_2$.

Question n°3.3.

D'après le TD : $u_i = (S^2(Y) - \lambda_i; -\text{cov}(X, Y))$, d'où les expressions des vecteurs propres $u_1 = (0.7071; 0.7071)$ et $u_2 = (-0.7071; 0.7071)$.

Notre fonction R nous renvoie les mêmes valeurs de vecteurs propres.

Question n°3.4.

On a $I_s(u_1)$ que l'on calcule directement avec une boucle sur R, d'où $I_s(u_1) = 0.3333 = \lambda_1$. De même avec $I_s(u_2) = 0.16667 = \lambda_2$.

On obtient aisément par calcul matriciel sur R les égalités $u^T A u = \lambda u^T u = \lambda \|u\|^2 = \lambda$.

Question n°3.5.

Après calcul matriciel sur R, on obtient une inertie totale de 0.5, ce qui correspond bien à la somme des deux inerties statistiques portées par chacun des axes principaux.

Question n°3.6.

Le graphique ci-dessous représente les 6 points ainsi que le tracé des deux vecteurs portant les axes principaux.

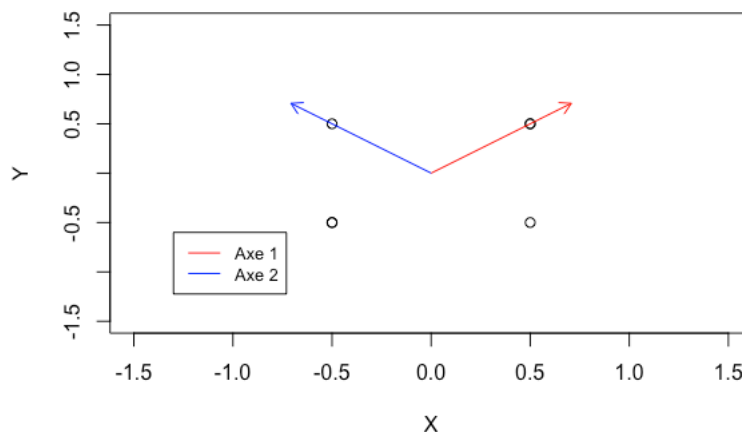


Figure 1 : Tracé des deux axes principaux dans le plan

Question n°3.7.

Le taux d'inertie expliqué par la 1^{ère} composante est d'environ 66%, alors que la 2nd composante explique quant à elle 33%. On remarque ici que la deuxième composante ne compte que pour 1/3 de l'inertie totale, on ne considérant que la proportion d'inertie expliquée par chacun des axes, il pourrait alors être judicieux de ne conserver que la première composante.

Question n°4.

On obtient les coordonnées des individus dans notre repère (O, C1, C2) par un rapide calcul vectoriel. Ces coordonnées correspondent aux lignes de la matrice Composantes retournées par notre fonction ACP.

On possède alors la matrice $C = \begin{pmatrix} -0.707 & 0.707 & 0 & 0 & 0.707 & -0.707 \\ 0 & 0 & 0.707 & -0.707 & 0 & 0 \end{pmatrix}^T$. La figure ci-dessous représente les points projetés dans notre nouveau repère. On constate qu'en 2 dimensions, en conservant les deux composantes, notre ACP revient simplement à faire une rotation de notre nuage de points.

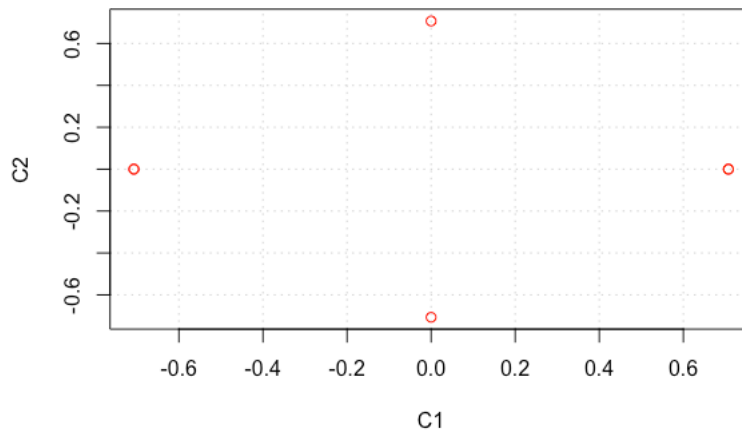


Figure 2 : Points projetés dans le repère (O, C1, C2)

Question n°5.

Sur l'espace E à 2 dimensions, la qualité de projection des individus vaut forcément 1 puisque l'on ne réduit pas la dimension. En revanche, en calculant pour chaque individu la qualité de sa projection sur l'espace E à une dimension engendré par le 1^{er} vecteur propre, on obtient un vecteur de la forme $Q = (1 \ 1 \ 0 \ 0 \ 1 \ 1)$. Pour la qualité des individus si l'on considère le sous-espace engendré par le 2nd vecteur propre, on obtient $Q = (0 \ 0 \ 1 \ 1 \ 0 \ 0)$.

On remarque que les individus sont parfaitement projetés sur un axe ou bien sur l'autre : effectivement les axes C1 et C2 correspondent en fait aux deux diagonales reliant les sommets du rectangle formés par les 6 individus. Les individus 3 et 4 appartiennent donc à l'axe C2, alors que les autres appartiennent à l'axe C1. Le repère (0, C1, C2) étant orthonormé, chaque axe correspond au donc noyau de la projection orthogonale sur l'autre axe.

Question n°6 et n°7.

De la même manière que précédemment, on obtient une matrice de contribution des individus

à une composante notée $C_t = \begin{pmatrix} 0.25 & 0 \\ 0.25 & 0 \\ 0 & 0.5 \\ 0 & 0.5 \\ 0.25 & 0 \\ 0.25 & 0 \end{pmatrix}$. Comme 4 des individus se trouvent sur l'axe C1,

ils n'y contribuent chacun qu'à 25%, alors que les deux individus de l'axe C2 y contribuent à 50% chacun (effectivement chaque individu est affecté d'un même poids).