

ACP: Analyse en Composantes principales

Introduction:

ACP est une approche qui est employée pour résumer l'information contenue dans une matrice (tableau) de données comportant un nombre important de variables numériques.

⚠ Si le table de données ne contient que deux variables, alors c'est facile d'étudier le lien entre ces deux variables en appliquant une analyse bivariée ; donc on calcule la corrélation entre ces variables.

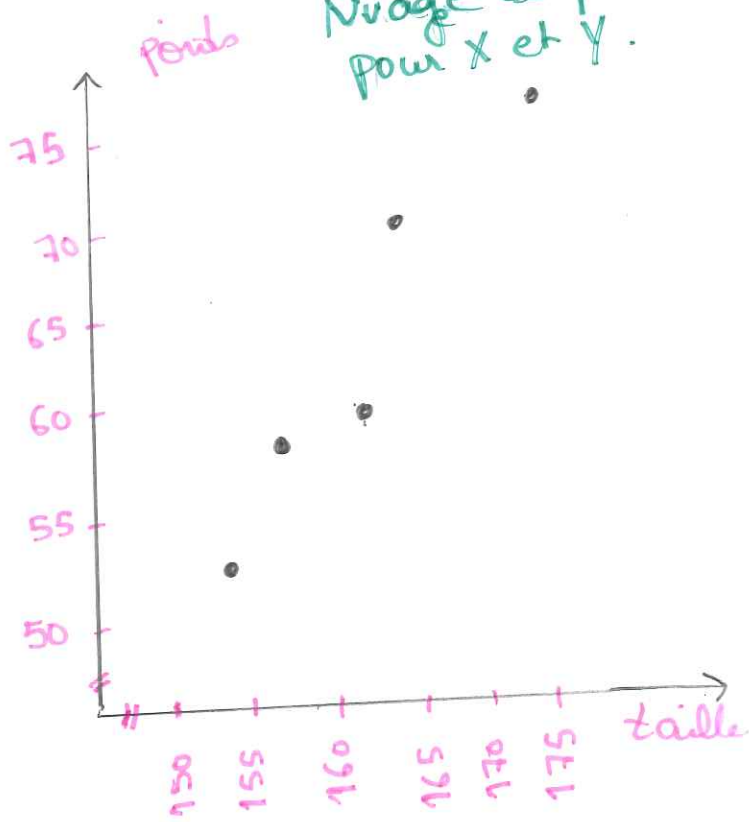
Mais Si le tableau contient beaucoup des variables cela devient impossible à appréhender. Alors, l'ACP permet alors de passer d'un nuage de points qui évolue dans un espace à P dimensions ($P = \text{nb de variables ou nb de colonnes}$) à une représentation en deux dimensions (ou éventuellement plusieurs représentations à 2 dimensions)

Exp:

On donne le tableau suivant :

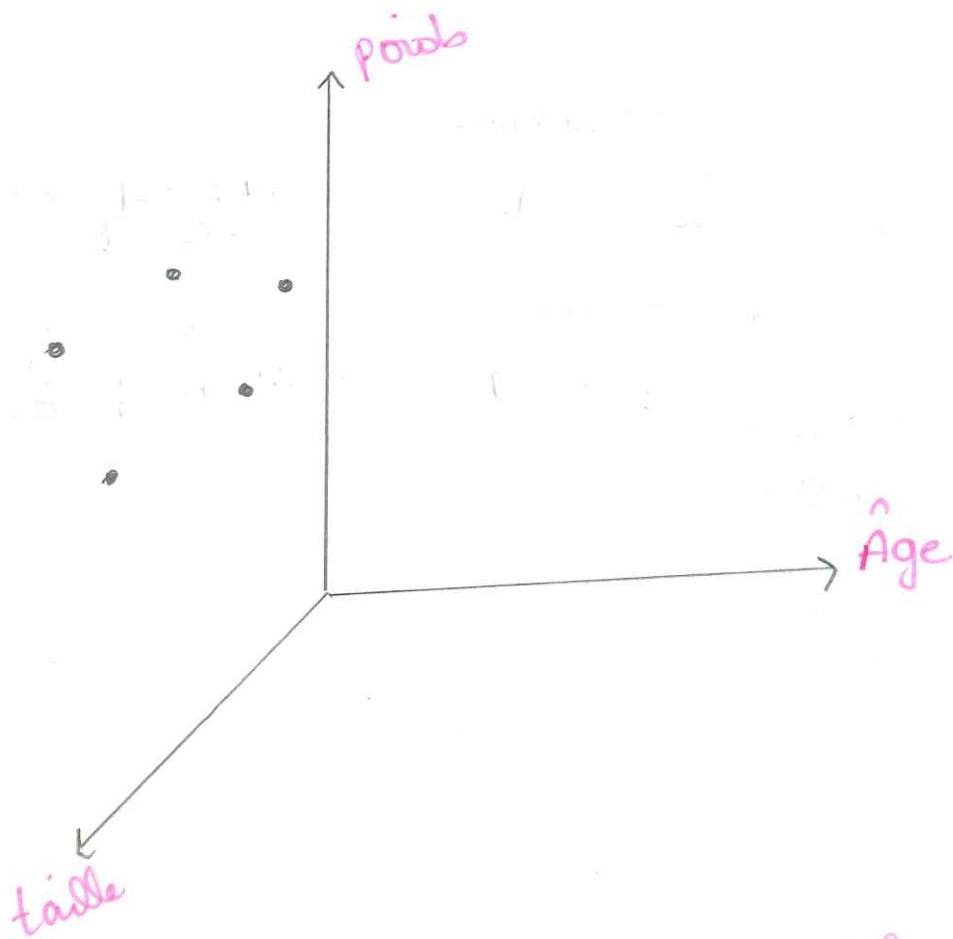
① Etude de deux variables :

Etudiant	taille X	Poids Y
Sarah	158	57
Olivier	169	70
Paul	174	75
Manon	155	52
Serena	163	58



② Si on a 3 variables : On a ajouté une 3^e variable Âge :

Etudiant	taille X	Poids Y	Âge Z
Sarah	158	57	21
Olivier	169	70	19
Paul	174	75	18
Manon	155	52	22
Serena	163	58	23



③ Si on ajoute une 4ème variable exemple : l'année d'obtention de Bac : ça sera difficile de représenter ces données, car on a 4 variables alors 4 axes (x, y, z, t) ; dans ce cas on applique

un ACP :

⚠ BUT :

• Maintenant ACP va rechercher l'axe dans lequel les observations sont le plus dispersées. on appelle cet axe : première composante principale C_1 .

• Ensuite : IL va chercher l'axe perpendiculaire à C_1 et qui va être la deuxième composante principale C_2 . (les axes C_1 et C_2 sont orthogonaux $\rightarrow \cos(\hat{C}_1, \hat{C}_2) = \text{correlat}(C_1, C_2) = 0$)

③

On considère le jeu de données avec (n) individus et (P) variables :

observation

V_1	V_2	...	V_P
X_{11}	X_{12}	...	X_{1P}
\vdots	\vdots		\vdots
X_{k1}	X_{k2}	...	X_{kP}
\vdots	\vdots		\vdots
X_{n1}	...		X_{nP}

Individu e_1

Individu e_n

P variables.

Centre de gravité : est le point dont les coordonnées sont les valeurs moyennes des variables :

$$G = (\bar{X}_1, \dots, \bar{X}_P)$$

Inertie :

$$I = \sum_{k=1}^n \|x_k - G\|^2$$

observation

Norme euclidienne

centre de gravité

Mesure l'information contenue dans un nuage de points. c'est la somme des distances au carré entre les observations et le centre de gravité.

Propriétés :

(4)

- Norme euclidienne pour vecteur $U(x, y)$
 $\|U\| = \sqrt{x^2 + y^2}$
- Norme euclidienne entre 2 pts A et B :
 $\|AB\| = \sqrt{(x_B - x_A)^2 + (y_B - y_A)^2}$

Propriétés de l'inertie:

$$I = \text{tr}(V) = \sum_{k=1}^p S_k^2$$

$$[S_k^2 = \text{Var}(X_k)]$$

V = matrice de var-cov de X

matrice symétrique
définie positive

$$V = \begin{pmatrix} S_1^2 & \dots & C_{X_1 X_p} \\ \vdots & \ddots & \vdots \\ C_{X_1 X_p} & \dots & S_p^2 \end{pmatrix}$$

covariance
entre X_1 et X_p

$$I = \sum_{k=1}^p \lambda_k$$

$$\lambda_1 + \dots + \lambda_p$$

valeurs propres de matrice
de var-cov \underline{V} .

→ Si les variables sont centrées réduites alors:

$$I = P$$

nb variables

Interprétation

* Inertie: Mesure les informations expliquées par l'axe.

* On fait la somme de % de variance donnée pour C_1 et C_2 pour l'avoir

Ici: 96,5% d'inertie expliquée par le plan (C_1, C_2)

• Centrer et réduire les variables:

① Distance entre 2 observations (individus) $\{k\}$ et $\{k'\}$:

$$\|e_k - e_{k'}\|^2 = \sum_{i=1}^p (x_{ki} - x_{k'i})^2$$

② Centrer et réduire les variables:

but: Si les variables ont pas m[^]e unité, certains variables à valeurs faibles « disparaissent » de l'information au profit de celles ayant de fortes valeurs.

• Comment faire? Replacer chaque variable X_i par:

$$X_i \rightarrow \frac{X_i - \bar{X}_i}{S_i} ; i = 1, \dots, p$$

Moyenne variable X_i
écart-type de variable X_i .

★ Produit scalaire:

→ Le produit scalaire entre 2 variables X_i et X_j est:

$$\langle X_i, X_j \rangle = \frac{1}{n} \sum_{k=1}^n x_{ki} x_{kj}$$

d'où la norme: $\|X_i\|^2 = \frac{1}{n} \sum_{k=1}^n (x_{ki})^2$

Si les variables sont centrées alors;

$$\langle X_i, X_j \rangle = \text{cov}(X_i, X_j) = C_{X_i X_j}$$

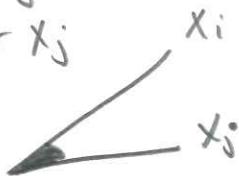
$$\|X_i\|^2 = S_i^2$$

$$\cos(\hat{X_i, X_j}) = \frac{\langle X_i, X_j \rangle}{\|X_i\| \cdot \|X_j\|} = \frac{C_{X_i X_j}}{S_i^2 \cdot S_j^2}$$

$= r_{X_i X_j}$ → corrélation entre X_i et X_j

$$\rightarrow \boxed{\cos(\hat{X_i, X_j}) = r_{X_i X_j}}$$

$r_{X_i X_j} > 0$, alors $\cos(\hat{X_i, X_j}) > 0$, donc il y a un angle aigu (entre 0 et $\frac{\pi}{2}$) entre X_i et X_j



$r_{X_i X_j} < 0$ alors $\cos(\hat{X_i, X_j}) < 0$ donc il y a un angle obtus (entre $\frac{\pi}{2}$ et π) entre X_i et X_j



X_i et X_j corrélation négative

$r_{X_i X_j} = 0$ alors $\cos(\hat{X_i, X_j}) = \frac{\pi}{2}$ donc les variables sont orthogonales
↓
pas corrélation linéaire

Principe ACP:

• Trouver des espaces des petites dimensions, sur lesquels les projections des observations minimisent la déviation de la réalité.

• On cherche un sous-espace F_q de \mathbb{R}^p (de dimension $q=2, 3, \dots$) sur lequel on projette le nuage de points.

→ On la projection de e_k dans le nouveau plan F_q .
Il faut que la distance entre l'individu e_k et sa projection sur F_q soit minimale.

→ $U_k \rightarrow$ Le vecteur propre unitaire de la matrice V associé à la k -ième plus grande valeur propre.

- L'inertie du nuage projeté sur U_k est λ_k .
- L'inertie du nuage projeté sur F_q est $\lambda_1 + \dots + \lambda_q$.
- L'inertie totale : $I = \lambda_1 + \dots + \lambda_q$.

→ Les vecteurs propres U_1, \dots, U_k sont appelés axes principaux.

* Le 1^{er} axe principal U_1 est associé à la plus grande valeur propre λ_1 . ; pareil U_2 pour λ_2 .

→ La projection des individus sur un axe principal est une nouvelle variable appelée composante principale.

→ la 1^{ère} composante C_1 représente les coordonnées des projections des individus sur l'axe U_1 .

Interpretation :

Représentation des observations :

	eigenvalue	% Var
C1	2.663	66.57
C2	1.19	29.94
C3		
C4		

cumulative Var.

$$\lambda_1 = 2.663$$

$$\lambda_2 = 1.19$$

C1 contient 66,5% de l'information (des informations)

C2 contient 29,94% de l'information

4 Variables



on obtient 4 Cmp principal

alors : ces deux axes
contient $66,5 + 29,94$
 $= 96,5\%$ de l'information.

% variance de C_i : $\frac{\lambda_i}{\lambda_1 + \dots + \lambda_p}$

$$\text{Var}(C_i) = S_{C_i}^2 = \lambda_i$$

A. Savoir :

$$C_i = \tilde{X} \cdot U_i$$

Composante principale i

Matrice de données centrée réduite

axe principal

Var C_i

Var X_i

$$\sum \lambda_i = p$$

nb variable

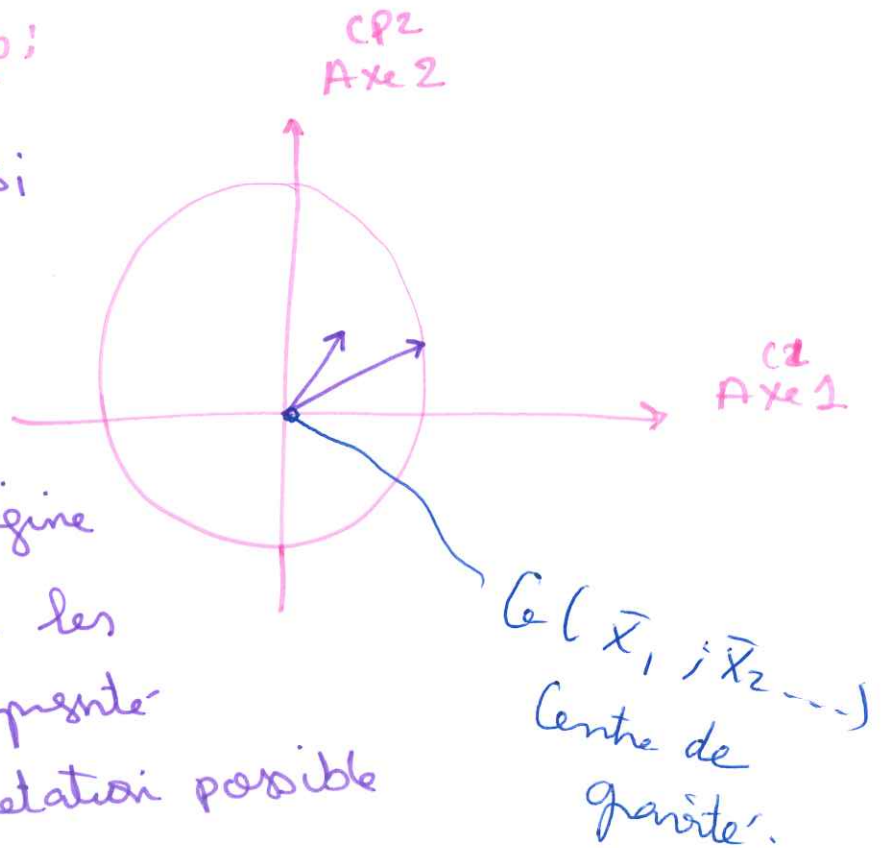
9

done $\left[\sum S_{C_i}^2 = \sum S_{X_i}^2 = \sum \lambda_i = p = I \right]$

Représentation des variables:

Variables bien représentées si elles sont proches du Cercle.

Variables proches de l'origine sont peu corrélées avec les axes; sont mal représentées; alors: pas d'interprétation possible pour ces variables.



Les individus sont bien représentés s'ils ne sont pas trop éloignés de l'axe sur lequel ont les projette.

→ Donc, il faut vérifier cos entre l'individu et l'axe (proche de 1).

→ valable si l'individu loin du centre de gravité.

