

Analyse en Composantes Principales

4. Mesures de liaison

4.1 Rappels

4.2 Mesures de liaison

1) Pour les individus

Soient X_i et X_j deux individus quelconques, donc deux vecteurs de \mathfrak{R}^n

On utilise **la distance Euclidienne**, elle est définie par :

$$d(X_i, X_j) = \left(\sum_{k=1}^n (x_i^k - x_j^k)^2 \right)^{1/2}.$$

2) Pour les variables

Soient X et Y deux variables de \mathfrak{R}^m , alors nous appliquons les deux mesures suivantes :

- **Covariance** : Elle permet de préciser le sens de la liaison entre les variables,

$$\text{cov}(X, Y) = \frac{1}{m} \times \sum_{j=1}^m (x_j - \bar{X}) \times (y_j - \bar{Y}).$$

Remarque La covariance n'est que $\langle \tilde{X}, \tilde{Y} \rangle_M$ pour \tilde{X} et \tilde{Y} les deux Variables centrées de X et Y respectivement et M est la matrice des poids.

- **Coefficient de corrélation** : Ce n'est qu'une normalisation de la covariance Il mesure l'intensité de la liaison linéaire entre les variables,

$$\rho_{X,Y} = \frac{\text{cov}(X,Y)}{\sigma_X \cdot \sigma_Y} = \frac{\text{cov}(X,Y)}{\sqrt{\text{var}(X)} \cdot \sqrt{\text{var}(Y)}}.$$

Remarque

Le coefficient de corrélation **prend ses valeurs entre -1 et 1**.

Le coefficient de corrélation mesure le **cosinus de l'angle inscrit** entre les variables. En effet,

$$\frac{\text{cov}(X,Y)}{\sqrt{\text{var}(X)} \cdot \sqrt{\text{var}(Y)}} = \frac{\langle \tilde{X}, \tilde{Y} \rangle_M}{\|\tilde{X}\|_M \cdot \|\tilde{Y}\|_M} = \cos(\theta_{\tilde{X}, \tilde{Y}}).$$

D'où,

$$\rho_{X,Y} = \cos(X,Y).$$

Définition

Deux variables sont dites :

- **décorrélées** si leur coefficient est nul, et
- **linéairement liées** si leur coefficient de corrélation est de **module égal à 1**.

Exemple Soit la matrice des données

$$X = \begin{pmatrix} 1 & 6 \\ 3 & 2 \end{pmatrix}.$$

Alors : les variables sont $X^1 = {}^t(1,3)$ **et** $X^2 = {}^t(6,2)$.

Comme :

$$\text{cov}(X,Y) = \frac{1}{m} \times \sum_{j=1}^m (x_j - \bar{X}) \times (y_j - \bar{Y}), \quad \rho_{X,Y} = \frac{\text{cov}(X,Y)}{\sigma_X \cdot \sigma_Y} = \frac{\text{cov}(X,Y)}{\sqrt{\text{var}(X)} \cdot \sqrt{\text{var}(Y)}}$$

Alors :

- Le **centre de gravité** du nuage est donné par :

$$g = (\bar{X}^1, \bar{X}^2) = \left(\frac{1+3}{2}, \frac{6+2}{2} \right) = (2,4).$$

La matrice **centrée** est la matrice :

$$Y = \begin{pmatrix} 1-2 & 6-4 \\ 3-2 & 2-4 \end{pmatrix}, \quad Y = \begin{pmatrix} -1 & 2 \\ 1 & -2 \end{pmatrix}.$$

- Les **variances** sont données par :

$$\text{Var}(X^1) = \frac{1}{2} \cdot \sum_{j=1}^2 (x_j^k - \bar{X}^1)^2 = (1/2) \cdot (1+1) = 1$$

$$\text{Var}(X^2) = \frac{1}{2} \cdot \sum_{j=1}^2 (x_j^k - \bar{X}^2)^2 = (1/2) \cdot (4+4) = 4.$$

D'où,

$$\text{cov}(X^1, X^2) = \frac{1}{2} \times (-2 - 2) = -2 < 0$$

(Produit scalaire relativement à la matrice des poids)

Sens de la liaison entre les deux variables

Donc,

Les deux variables sont de sens contraire.

$$\rho_{X^1, X^2} = \frac{\text{cov}(X^1, X^2)}{\sigma_{X^1} \cdot \sigma_{X^2}} = \frac{\text{cov}(X, Y)}{\sqrt{\text{var}(X)} \cdot \sqrt{\text{var}(Y)}} = \frac{-2}{1 \cdot 2} = -1$$

Donc,

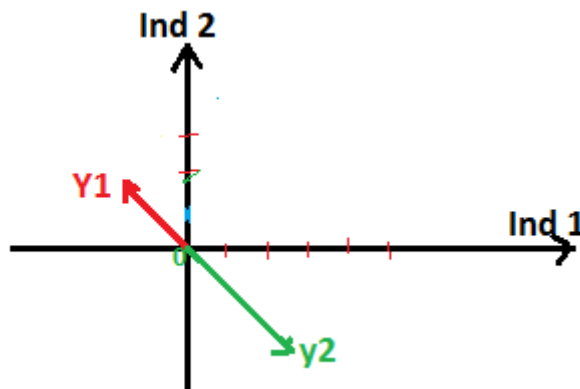
L'angle entre les deux variables est égal à π .

C'est à dire

Les deux variables sont colinéaires mais de sens contraire

Remarque

Ces résultats peuvent être interprétés graphiquement :



Presentation des variables dans l'espace des individus.

5. Analyse des points

Puisque notre but est de *réduire la dimension* de l'espace des données,

Ça ne peut être que par *projections* !

Donc,

L'analyse des données en composantes principales consiste à **étudier les projections**
 des points du nuage sur

- i) un *axe*, dimension = 1
- ii) un *plan*, dimension = 2
- iii) un *hyperplan*, de dimension = 3

Ce qui revient à déterminer le *meilleur ajustement du nuage* par

un sous espace de \mathbb{R}^n pour $n \geq 1$.

i) Ajustement du nuage par une droite (un axe)

Soit la droite (Δu) passant par l'origine, et engendré par le vecteur unitaire \vec{u} i.e $\|\vec{u}\| = 1$

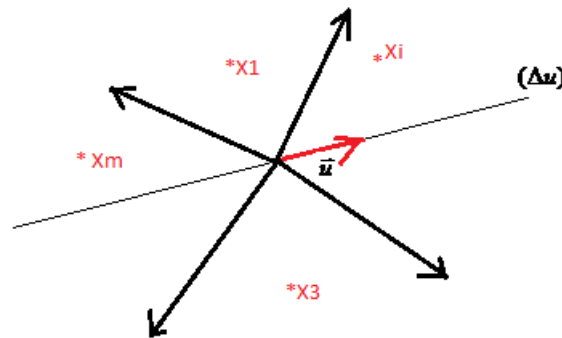


Figure 1. Présentation de l'axe de projection et des individus dans l'espace \mathcal{R}^n .

Alors la projection de X_i sur (Δu) notée x_i est présentée par :

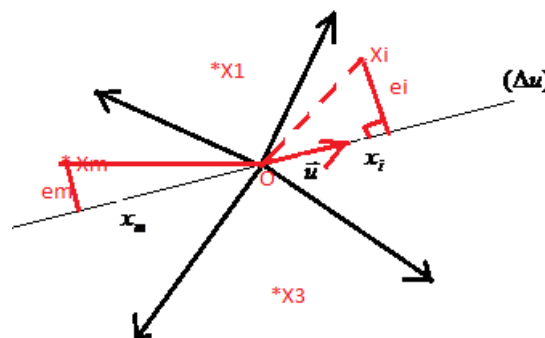


Figure 2. Illustration des projections des individus sur l'axe (Δu) .

Elle est définie par :

$$x_i = X_i \times \vec{u}$$

et ce n'est que le produit scalaire : $\langle X_i, \vec{u} \rangle = \langle \vec{u}, X_i \rangle$.

Partant du **théorème de Pythagore** (voir figure 2), nous avons :

$$x_i^2 + e_i^2 = \|X_i\|^2.$$

Ce qui donne

$$x_i^2 = \|X_i\|^2 - e_i^2.$$

Donc, ça peut être interprété comme suit :

x_i Représente *l'information projetée* sur la l'axe i.e la droite
 e_i n'est que l'erreur ou bien *l'information perdue*.

Pour réduire la dimension de l'espace, il faut donc :

maximiser l'information x_i , il suffit de minimiser la perte d'information e_i .

pour chaque individu X_i , $i = 1, 2, \dots, m$.

Donc pour avoir **l'information totale maximale**, il faut :

- 1) **Projeter** tous les points (*individus*) sur la droite (Δu) .
- 2) **Choisir** \vec{u} tel que la *somme des carrées de ces projections soit maximale* c'est à dire

maximiser : $\sum_{i=1}^m x_i^2$.

Or $x_i^2 = \langle X_i \vec{u}, X_i \vec{u} \rangle = {}^t(X_i \vec{u}) \times (X_i \vec{u})$, alors :

$$\begin{aligned} \sum_{i=1}^m x_i^2 &= \sum_{i=1}^m {}^t(X_i \vec{u}) \times (X_i \vec{u}) = \sum_{i=1}^m {}^t \vec{u} \cdot {}^t X_i \cdot X_i \cdot \vec{u} \\ &= {}^t \vec{u} \times \left(\sum_{i=1}^m {}^t X_i \cdot X_i \right) \times \vec{u} = {}^t \vec{u} \times ({}^t X \cdot X) \times \vec{u} \end{aligned}$$

D'où tout revient à maximiser : ${}^t \vec{u} \times ({}^t X \cdot X) \times \vec{u}$ sous la contrainte : $\|\vec{u}\| = 1$.

Que nous pouvons présenter comme suit :

Maximiser :
$$\begin{cases} {}^t \vec{u} \cdot ({}^t X \cdot X) \cdot \vec{u} \\ \|\vec{u}\|^2 = {}^t \vec{u} \cdot \vec{u} = 1 \end{cases} \quad (P).$$

Il s'agit d'un problème d'optimisation avec contraintes.

Pour que nous puissions résoudre ce problème, nous donnons le rappel suivant :

Rappels

Multiplicateur de Lagrange Le multiplicateur de Lagrange est une méthode d'optimisation permettant de trouver les points stationnaires d'une fonction dérivable sous-contraintes.

Formellement, l'écriture du Lagrangien est donnée par :

$$L(X, \lambda) = f(x) + \lambda \cdot g(x).$$

Avec :

- X : les variables figurant dans la fonction à optimiser.
- $f(x)$: la fonction à optimiser.

- λ : le multiplicateur de Lagrange = inconnu à déterminer.
- $g(x)$: la contrainte à imposer dans le problème à résoudre.

A lors **la solution est obtenue** en résolvant le système des dérivées partielles suivant (notons bien qu' il s'agit d'une condition nécessaire d'existence de solution) :

$$Sol = \begin{cases} \frac{\partial L}{\partial x_i} = \frac{\partial f}{\partial x_i} + \lambda \frac{\partial g}{\partial x_i} = 0, & \forall i \text{ et } X = (x_i)_i \\ \frac{\partial L}{\partial \lambda} = g(x) = 0 \end{cases}.$$