

# Introduction à l'analyse des données

O. Mory

LMI-SFA-UNA

2 février 2018

- Analyse en Composantes Principales (A.C.P.) (K. Pearson, 1901)
- Analyse factorielle des correspondances simples (A.F.C.)(Hirschfeld, 1936)
- Analyse factorielle des correspondances multiples (A.C.M.)(Guttman, 1941)
- Méthodes de classification automatique

# Analyse des données

Analyser des données, c'est extraire d'une masse d'informations brutes, des éléments de réponse aux questions qui résultent des objectifs globaux poursuivis.

# Analyse des données

## Traitement de données en masse

- Grand nombre d'individus
- Grand nombre de variables

# Analyse des données

## Développement parallèle à l'informatique

- Fichiers volumineux  $\Rightarrow$  demande de méthodes
- Capacité de calcul  $\Rightarrow$  méthodes praticables

# Analyse des données

Les outils mathématiques de l'analyse des données :

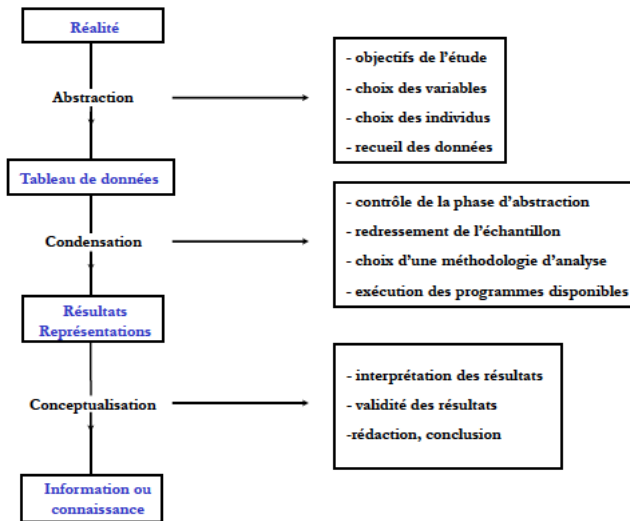
- Algèbre linéaire
- Calcul matriciel

# Analyse des données

Les outils mathématiques de l'analyse des données :

- Analyse en Composantes Principales (A.C.P.)
- Analyse factorielle des correspondances  
analyse de variables qualitatives
  - Correspondances simples (A.F.C.) (Étude d'un tableau de contingence)
  - Correspondances multiples (A.C.M.) (Utile lors du dépouillement d'enquêtes)

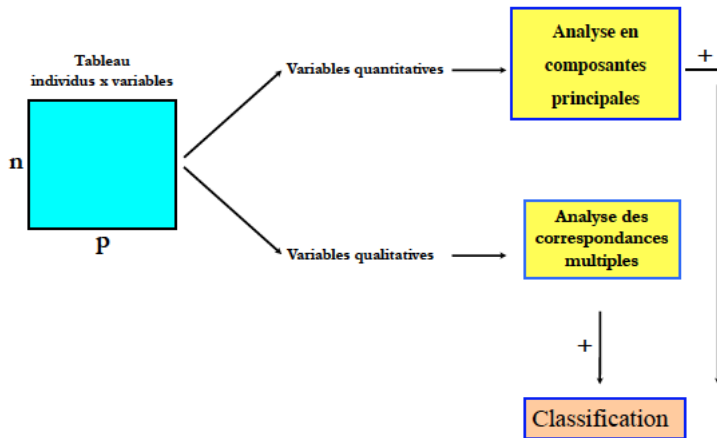
# Introduction





# Introduction

## METHODES DESCRIPTIVES



# ACP

Données :

$n$  individus observés sur  $p$  variables quantitatives. L'A.C.P. permet d'explorer les liaisons entre variables et les ressemblances entre individus

Résultats

- ⇒ Visualisation des individus (Notion de distances entre individus)
- ⇒ Visualisation des variables (en fonction de leurs corrélations)

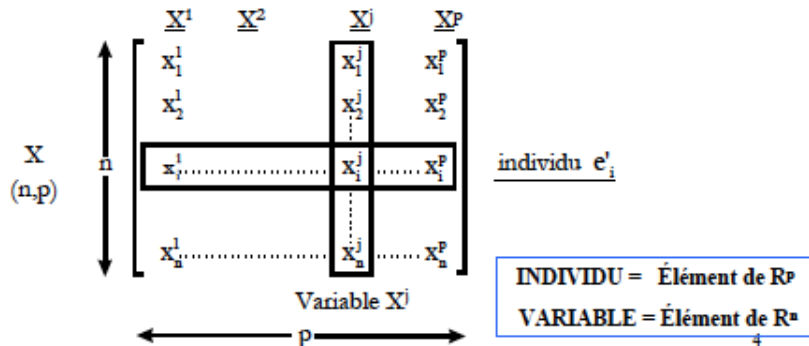
# Interprétation des axes

- Mesurer la qualité des représentations obtenues
  - critère global
  - critères individuels
- Utilisation éventuelle de variables supplémentaires

- ① LES DONNÉES
- ② PRINCIPE DE L'A.C.P.
- ③ LE CHOIX DE LA DISTANCE ENTRE INDIVIDUS
- ④ INERTIE TOTALE

# Les données

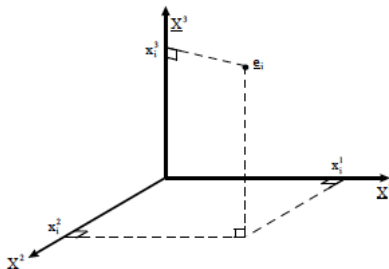
$p$  variables quantitatives observées sur  $n$  individus.



# On cherche à représenter le nuage des individus

A chaque individu noté  $e_i$ , on peut associer un point dans  $R^p$  = espace des individus.

A chaque variable du tableau  $X$  est associé un axe de  $R^p$ .



Comment faire la représentation 4D ?

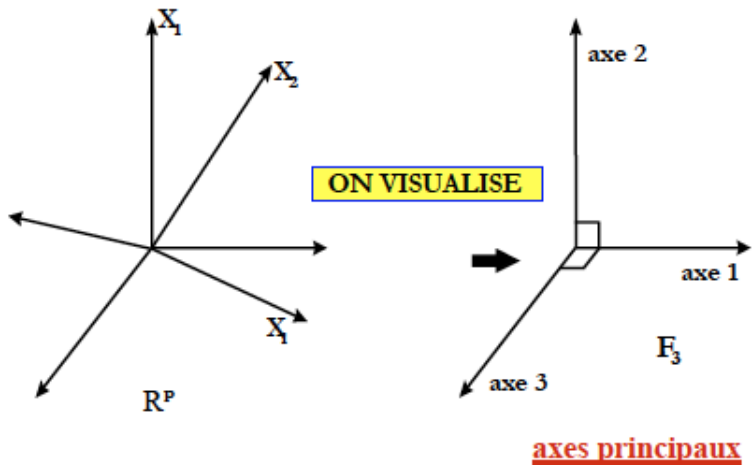
# Principe de l'ACP

## Principe

On cherche une représentation des  $n$  individus, dans un sous-espace  $F_k$  de  $R^p$  de dimension  $k < p$  ( $k$  petit 2, 3 ... ; par exemple un plan)  
Autrement dit, on cherche à définir  $k$  nouvelles variables combinaisons linéaires des  $p$  variables initiales qui feront perdre le moins d'information possible.

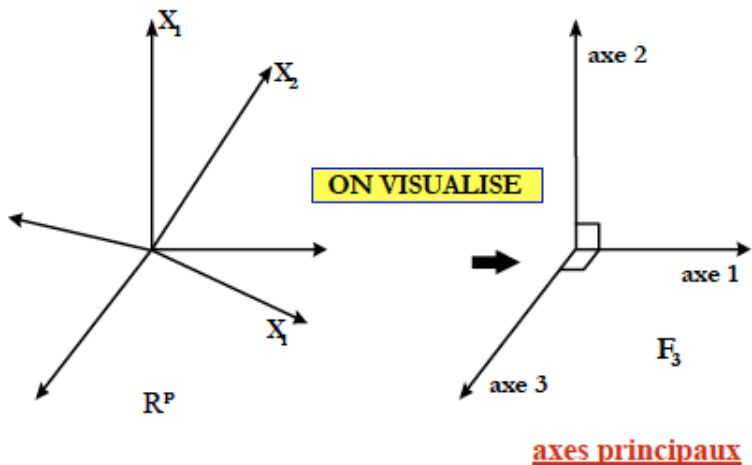
## Vocabulaire

- Ces variables seront appelées « **composantes principales** »,
- les axes qu'elles déterminent : « **axes principaux** »
- les formes linéaires associées : « **facteurs principaux** »



Qu'espere t'on ?





Qu'espre t'on ?  $\Rightarrow$  Perdre le moins d'information possible

①

$F_k$  devra être « ajusté » le mieux possible au nuage des individus: la somme des carrés des distances des individus à  $F_k$  doit être minimale.



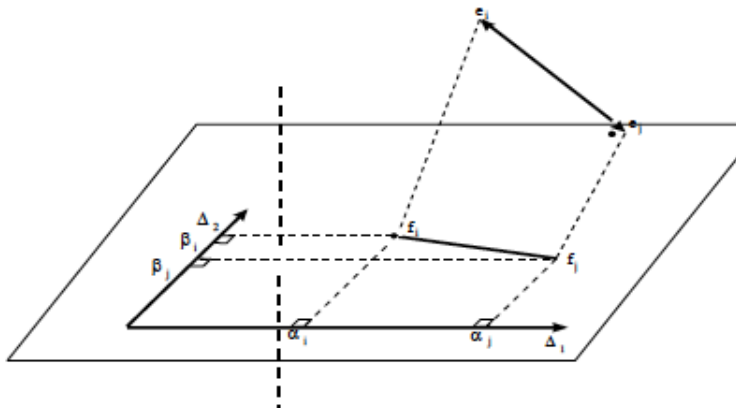
②

$F_k$  est le sous-espace tel que le nuage projeté ait une **inertie** (dispersion) maximale.

① et ② sont basées sur les notions de :

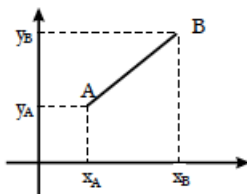
**distance**

**projection orthogonale**



La distance entre  $f_i$  et  $f_j$  est inférieure ou égale à celle entre  $e_i$  et  $e_j$

# Notion de distance euclidienne



Dans le plan:

$$d^2(A, B) = (x_B - x_A)^2 + (y_B - y_A)^2$$

Dans l'espace  $\mathbb{R}^p$  à  $p$  dimensions, on généralise cette notion : la distance euclidienne entre deux individus s'écrit:

$$e_i = (x_i^1 \ x_i^2 \ \dots \ x_i^p) \quad e_j = (x_j^1 \ x_j^2 \ \dots \ x_j^p)$$

$$d^2(e_i, e_j) = (x_i^1 - x_j^1)^2 + (x_i^2 - x_j^2)^2 + \dots + (x_i^p - x_j^p)^2$$

$$d^2(e_i, e_j) = \sum_{k=1}^p (x_i^k - x_j^k)^2$$

Le problème des unités ?

10

## Que faire ?

# Normaliser

Pour résoudre ce problème, on choisit de transformer les données en données centrées-réduites.

L'observation  $x_i^k$  est alors remplacée par :

**UNITÉS D'ÉCART TYPE:**

$$\frac{x_i^k - \bar{X}^k}{s_k}$$

où :  $\bar{X}^k$  = moyenne de la variable  $X^k$

$s_k$  = écart-type de la variable  $X^k$

# Inertie

$$I_{\underline{g}} = \sum_{i=1}^n \frac{1}{n} d^2(e_i, \underline{g})$$

ou de façon plus générale

$$I_{\underline{g}} = \sum_{i=1}^n p_i d^2(e_i, \underline{g})$$

avec  $\sum_{i=1}^n p_i = 1$

L'inertie est la somme pondérée des carrés des distances des individus au centre de gravité  $\underline{g}$

L'inertie mesure la dispersion totale du nuage de points.

# Inertie

L'inertie est donc aussi égale à la somme des variances des variables étudiées.

En notant  $V$  la matrice de variances-covariances :

$$V = \begin{pmatrix} s_1^2 & s_{12} & \dots & s_{1p} \\ \vdots & s_2^2 & & \vdots \\ \vdots & & & \vdots \\ s_{p1} & & & s_p^2 \end{pmatrix}$$

$$I_g = \sum_{i=1}^p s_i^2$$

$$I_g = \text{Tr}(V)$$

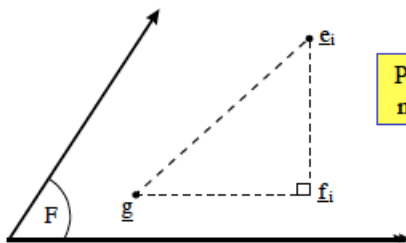
## Remarque

Dans le cas où les variables sont centrées réduites, la variance de chaque variable vaut 1.

L'inertie totale est alors égale à  $p$  (nombre de variables).

# Inertie

## Équivalence des deux critères concernant la perte d'information



Projection orthogonale du nuage sur un sous-espace

Soit  $F$  un sous-ensemble de  $\mathbb{R}^p$

$\underline{f}_i$  la projection orthogonale de  $\underline{e}_i$  sur  $F$

$$\|\underline{e}_i - \underline{g}\|^2 = \|\underline{e}_i - \underline{f}_i\|^2 + \|\underline{f}_i - \underline{g}\|^2 \quad \forall i = 1 \dots n$$

14



# Inertie

On va chercher  $F$  tel que :

$$\textcircled{1} \quad \sum_{i=1}^n p_i \|\mathbf{e}_i - \mathbf{f}_i\|^2 \text{ soit minimal}$$

ce qui revient d'après le théorème de Pythagore à maximiser :

$$\textcircled{2} \quad \sum_{i=1}^n p_i \|\mathbf{f}_i - \mathbf{g}\|^2$$

# Inertie

$$\|\underline{e}_i - \underline{g}\|^2 = \|\underline{e}_i - \underline{f}_i\|^2 + \|\underline{f}_i - \underline{g}\|^2 \quad \forall i = 1 \dots n$$

$$\text{Donc : } \underbrace{\sum_{i=1}^n p_i \|\underline{e}_i - \underline{g}\|^2}_{\text{Inertie totale}} - \underbrace{\sum_{i=1}^n p_i \|\underline{e}_i - \underline{f}_i\|^2}_{\text{minimiser cette quantité (carrés des distances entre points individus et leurs projections)}} = \underbrace{\sum_{i=1}^n p_i \|\underline{f}_i - \underline{g}\|^2}_{\text{maximiser l'inertie du nuage projeté}}$$

$\Leftrightarrow$

# Axes principaux

- On appelle axes principaux d'inertie les axes de direction les vecteurs propres de  $V$  normés à 1.
- Il y en a  $p$ .
- Le premier axe est celui associé à la plus grande valeur propre . On le note  $u^1$
- Le deuxième axe est celui associé à la deuxième valeur propre . On le note  $u^2$

# Composantes principales

- À chaque axe est associée une variable appelée composante principale.
- La composante  $c_1$  est le vecteur renfermant les coordonnées des projections des individus sur l'axe 1.
- La composante  $c_2$  est le vecteur renfermant les coordonnées des projections des individus sur l'axe 2.
- Pour obtenir ces coordonnées, on écrit que chaque composante principale est une combinaison linéaire des variables initiales.

$$c^1 = u_1^1 * X^1 + u_1^2 * X^2 + \dots + u_1^p * X^p$$

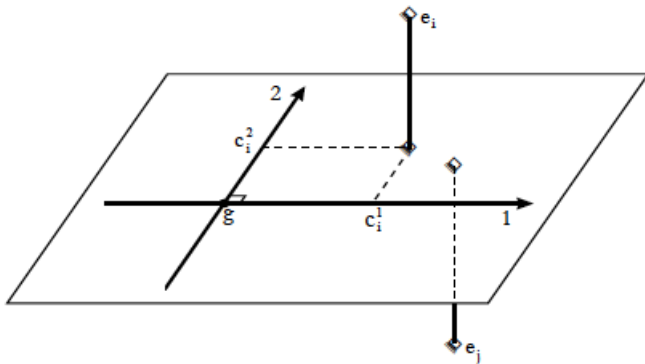
# PROPRIÉTÉS DES COMPOSANTES PRINCIPALES

- La variance d'une composante principale est égale à l'inertie portée par l'axe principal qui lui est associé.
  - 1ère composante  $c^1$  variance :
  - 2ème composante  $c^2$  variance :
  - 3ème composante  $c^3$  variance :
- Les composantes principales sont non corrélées deux à deux. En effet, les axes associés sont orthogonaux.  $\lambda_1$

# REPRÉSENTATION DES INDIVIDUS

La  $j^{\text{ème}}$  composante principale  $\underline{c}^j = \begin{pmatrix} c_1^j \\ c_2^j \\ \vdots \\ c_n^j \end{pmatrix}$  fournit les coordonnées des  $n$  individus sur le  $j^{\text{ème}}$  axe principal.

Si on désire une **représentation plane** des individus, la meilleure sera celle réalisée grâce aux **deux premières composantes principales**.

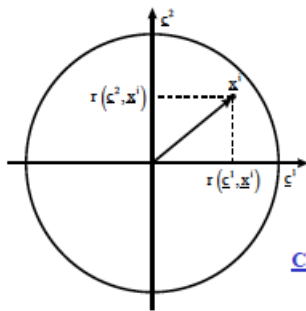


**Attention à la qualité de représentation de chaque individu!**

# REPRÉSENTATION DES VARIABLES

Les « proximités » entre les composantes principales et les variables initiales sont mesurées par les covariances, et surtout **les corrélations**.

$r(\underline{c}^j, \underline{x}^i)$  est le **coefficient de corrélation linéaire** entre  $\underline{c}^j$  et  $\underline{x}^i$



**CERCLE DES CORRÉLATIONS**

23



# INTERPRETATION DES « PROXIMITÉS » ENTRE VARIABLES

On utilise un **produit scalaire** entre variables permettant d'associer aux paramètres courants : écart-type, coefficient de corrélation linéaire des représentations géométriques.

$$\left\langle \underline{x}^i, \underline{x}^j \right\rangle = \frac{1}{n} \sum_{k=1}^n x_k^i x_k^j$$

On suppose les **variables centrées**.

$$\langle \underline{x}^i, \underline{x}^j \rangle = \text{Cov}(\underline{x}^i, \underline{x}^j)$$

$$\|\underline{x}^i\|^2 = \langle \underline{x}^i, \underline{x}^i \rangle = \frac{1}{n} \sum_{k=1}^n (x_k^i)^2$$

$$\|\underline{x}^i\|^2 = s_i^2 \quad \text{Variance de } \underline{x}^i$$

$$\|\underline{x}^i\| = s_i \quad \text{Écart-type de } \underline{x}^i$$

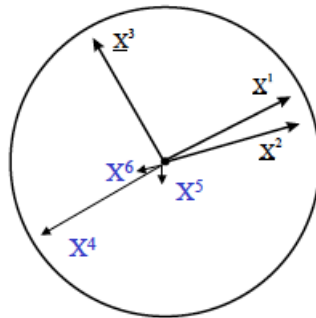
# Coefficient de corrélation linéaire

$$\text{Cos}(\widehat{\underline{X}^i, \underline{X}^j}) = \frac{\langle \underline{X}^i, \underline{X}^j \rangle}{\|\underline{X}^i\| \|\underline{X}^j\|} = \frac{\text{Cov}(\underline{X}^i, \underline{X}^j)}{s_i s_j} = r(\underline{X}^i, \underline{X}^j)$$

Le cosinus de l'angle formé par les variables  $X^i$  et  $X^j$  est le coefficient de corrélation linéaire de ces deux variables

$X^1$  et  $X^2$  ont une  
corrélation proche de 1.

$X^1$  et  $X^3$  ont une  
corrélation proche de 0.



CERCLE DES CORRÉLATIONS

# VALIDITÉ DES REPRÉSENTATIONS

$$\frac{\lambda_i}{\lambda_1 + \lambda_2 + \dots + \lambda_p}$$

mesure la part d'inertie expliquée par l'axe  $i$ .

Exemple :

$$\frac{\lambda_1 + \lambda_2}{\sum_{i=1}^p \lambda_i}$$

est la part d'inertie expliquée par le premier plan principal.

Ce critère (souvent exprimé en pourcentage) mesure le degré de reconstitution des carrés des distances.

La réduction de dimension est d'autant plus forte que les variables de départ sont plus corrélées.

28

# Combien d'axes ?

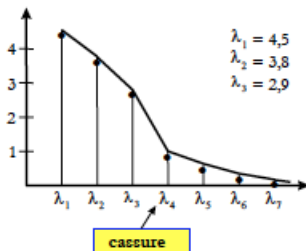
Différentes procédures sont complémentaires:

- ① Pourcentage d'inertie souhaité : a priori
  - ② Diviser l'inertie totale par le nombre de variables initiales
- ⇒ inertie moyenne par variable : I.M.

Conserver tous les axes apportant une inertie supérieure à cette valeur I.M.  
(inertie > 1 si variables centrées réduites).

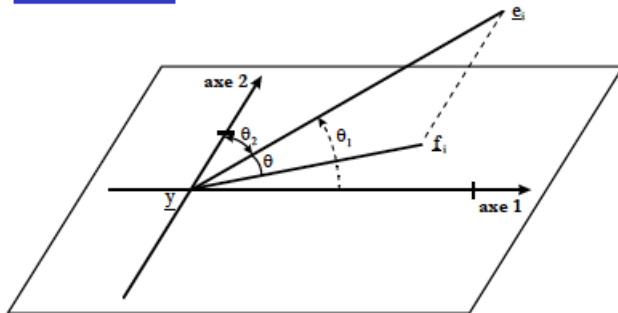
## ③ Histogramme

Conserver les axes associés aux valeurs propres situées avant la cassure.



29

## Cosinus carrés



$$\cos^2 \theta = \cos^2 \theta_1 + \cos^2 \theta_2$$

30

# Contributions

Il est très utile aussi de calculer pour chaque axe la **contribution apportée** par les divers individus à cet axe.

Considérons la  $k^{\text{ième}}$  composante principale  $\underline{c}^k$ , soit  $\underline{c}_i^k$  la valeur de la composante pour le  $i^{\text{ème}}$  individu.

$$\sum_{i=1}^n \frac{1}{n} (\underline{c}_i^k)^2 = \lambda_k$$

La **contribution** de l'individu  $\underline{e}_i$   
à la composante n°  $k$  est définie par

$$\frac{\frac{1}{n} (\underline{c}_i^k)^2}{\lambda_k}$$



# REPRÉSENTATION DES VARIABLES

Le cercle des corrélations est la projection du nuage des variables sur le plan des composantes principales.

