

# Analyse de données multivariées

## *Analyses factorielles*

### **Visualisation et réduction de dimension - statistiques exploratoires**

Mireille Batton-Hubert - Institut H.Fayol  
 UP3: 'Analyse de données et données temporelles'  
*Majeure Sciences des données- octobre 2020*

1

## UP3: Analyse de données et données temporelles

Contenu (30 h) :

- Module (15h) statistiques exploratoires : comprendre et faire bon usage des techniques pour décrire, visualiser et analyser les données
- Module séries temporelles (15 h)
  - Visualisation et réduction de dimension : analyse en composantes principales (3h) ( analyse factorielle) ACP (PCA) # cours + **TP1**
  - ACP et décomposition : (3h) # cours + **TP2** : *TP1 + TP2 : 1 note*
  - Cas des variables qualitatives. Analyse factorielle des correspondances (3h) AFC cours + TP3
  - Analyse discriminante descriptive (3h). AFD : cours + TP4
  - Mise en œuvre sur un cas gros TP (AFD + ACP) (soit 3h en séance de TP + travail personnel) : **TPF** noté

**Pas d'examen mais rendus de TP notés : date limite des rendus 28 octobre**

2

# Définition de l'analyse de données

- On s'intéresse à des données statistiques représentant un échantillon composé de  $n$  individus décrits par un ensemble de variables  $X$  de dimension  $p$  (souvent de grande dimension)
- Ces données sont :
  - quantitatives sur  $\mathbb{R}$
  - qualitatives (ordinales ; catégorielle)
- Les techniques de traitement de ces données constituent l'analyse de données multi-variées (multidimension)
  - Analyse en composantes principales : ACP - PCA (Principal Component Analysis)
  - Analyse des correspondances
  - Analyse factorielle discriminante
- Ces méthodes appartiennent aux **statistiques multidimensionnelles** (regroupe des méthode de description + réduction de dimension et des méthodes de classification)

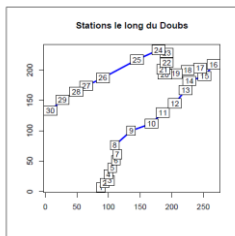
2

# Des exemples de données

- Un échantillon de données avec des individus :
- Du plus simple* (un échantillon de notes , de mesures d'individus ...)

*une étude d'un botaniste : dimensions de 15 fleurs d'iris avec 3variables  
longueur du sépale  $x_1$ , largeur du sépale  $x_2$ , longueur du pétale  $x_3$*

- au plus sophistiqué* : des données de capteurs ,
- un suivi de couts / produits financiers
- Échantillon pouvant être en continu



11 variables physico-chimiques :  
distance à la source, altitude,  
Pente, débit, pH, dureté de l'eau  
phosphate, nitrate, ammoniacque,  
oxygène, demande biologique en  
oxygène

Fleur n°	$x_1$	$x_2$	$x_3$
1	5.1	3.5	1.4
2	4.9	3.0	1.4
3	4.7	3.2	1.3
4	4.6	3.1	1.5
5	5.0	3.6	1.4
6	7.0	3.2	4.7
7	6.4	3.2	4.5
8	6.9	3.1	4.9
9	5.5	2.3	4.0
10	6.5	2.8	4.6
11	6.3	3.3	6.0
12	5.8	2.7	5.1
13	7.1	3.0	5.9
14	6.3	2.9	5.6
15	6.5	3.0	5.8

- Données temporelles ou non ...
- De grande dimension (ACP dite de grande dimension...)

4

## Objectifs et organisation

- ✓ Comprendre ce qu'est l'analyse de données multivariées
- ✓ Comment mettre en œuvre ces analyses et leur signification
- ✓ Comment obtenir et interpréter des résultats
- séances : cours (50%) + TD / TP info (50%) ( apprentissage des outils et programmation) qui serviront à l'étude de cas

### Organisation du cours :

- cours + TP pour chaque technique factorielle : à chaque TP un rendu de l'avancement sur *campus*

Utilisation de 2 TP pour l'étude du cas (3h00) qui sera notée : rendu fait pour le 28 octobre

5

## L'enjeu de l'analyse de données

Ces données statistiques : nombreuses variables

- dont aucune n'est parfaite
- La perception d'un phénomène appréhendé comme la combinaison d'un grand nombre de variables
- **Comment faire pour tenir compte de l'ensemble de l'information ?**

### Moyens ?

Faire de tableaux croisés ( var  $X \times Y$  ) :

Si 10 variables soit 45 tableaux croisés si 100 var 4950 tableaux croisés

Autre méthode : dite des indices exp: indice d'inflation , indice du développement humain, BIP40 avec une somme pondérée mais formule arbitraire avec des pondérations ...

Trouver des méthodes pour synthétiser **les variables sans trop les déformer**

6

## En fonction du type de variables: différentes méthodes de réduction de dimension

### Méthodes factorielles de représentation

- Analyse en composantes principales (ACP) : variables quantitatives
- Analyse factorielle de correspondance (AFC) analyse d'un tableau croisé de deux variables qualitatives (caractère)
- Analyse des correspondances multiples (ACM) plusieurs variables qualitatives
- Analyse factorielle discriminante ou analyse discriminante décisionnelle ( analyse discriminante ): liaisons existants entre un caractère à expliquer qualitatif et un ensemble de var. quantitatives
- Analyse factorielle multiple ( AFM)

7

## Analyse en composantes principales (ACP) Principal component analysis (PCA)

Appelée aussi la transformation de Karhunen–Loève (KLT)<sup>1</sup>

[Karl Pearson](#) (1901)

[Harold Hotelling](#) ( 1930)

8

# 1. Objectifs de l'ACP (1)

In *Dictionary of statistics and methodology* » (Vogt, 1993) :

‘L’analyse en composantes principales : Ensemble de méthodes permettant de procéder à des transformations linéaires d’un grand nombre de variables intercorrélées de manière à obtenir un nombre relativement limité de composantes non corrélées. Cette approche facilite l’analyse en regroupant les données en des ensembles plus petits et en permettant d’éliminer les problèmes de multicolinéarité entre les variables. L’analyse en composantes principales s’apparente à l’analyse factorielle, mais c’est une technique indépendante qui est souvent utilisée comme première étape à une analyse factorielle ‘ (Vogt, 1993, page 177).



**A pour premier objectif de réduire la dimension de données**

**Exemple des Iris**

Iris - échantillon :  $p=3$  variables -  $n$  individus = 15

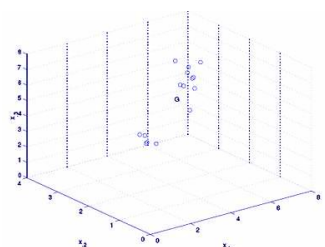
9

# 1. Objectifs de l'ACP (2)

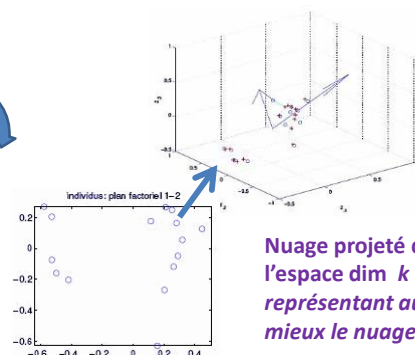
**A pour premier objectif de réduire la dimension de données**

**exemple des Iris - échantillon :  $p=3$  variables -  $n$  individus = 15**

Trouver des axes ( ~indices ) qui respectent la forme du nuage multidimensionnel initial c.a.d la forme des relations entre les variables



Nuage initiale  
dans l'espace  
de dim  $p=3$



Nuage projeté dans  
l'espace dim  $k=2 < p$   
représentant au  
mieux le nuage

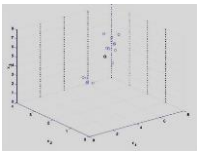
10

# 1. Objectifs de l'ACP(3)

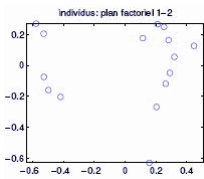
**Données :** n individus observés sur  $p$  variables quantitatives.  
L'A.C.P. permet d'explorer les liaisons entre variables et les ressemblances entre individus.

**Pour construire un espace plus petit  $k < p$  qui conserve ces propriétés**

$p=3$



$k=2$



**Résultats : dans ce nouvel espace de dimension  $k$**   
Visualisation des individus : (Notion de distances entre individus)  
Visualisation des variables : (en fonction de leurs corrélations)  
Interprétation des résultats

11

# 2. Pourquoi : veut-on réduire (1) ?

Les ossements de cranes de 42 canidés, identifiés comme chiens et loups, ont été mesurés. Les variables sont :

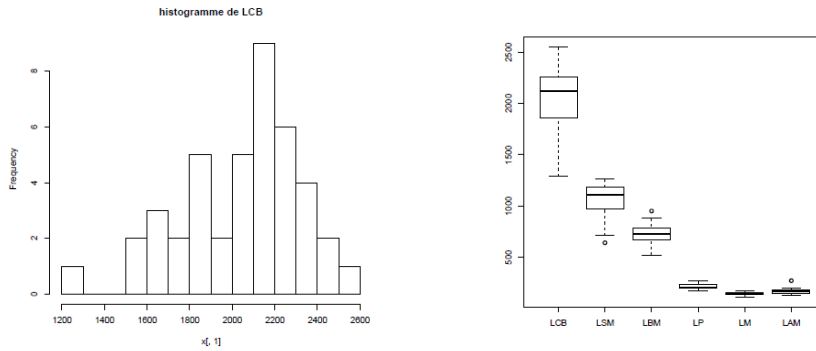
- LCB : longueur condylo-basale
- LSM : longueur de la mâchoire supérieure
- LBM : largeur bi-maxillaire
- LP : longueur de la carnassière supérieure
- LM : longueur 1° molaire supérieure
- LAM : largeur 1° molaire supérieure

Source : JAMBU - Classification

TYPE	LCB	LSM	LBM	LP	LM	LAM
BULL-DOG 1	1290	640	950	175	112	138
BULL-DOG 2	1540	740	760	200	142	165
BOXER	1580	710	710	167	125	133
SAINT-BERNARD	2200	1110	880	225	154	180
BULL-MASSIF	1900	930	780	197	132	140
DOGUE ALLEMAND 1	2410	1190	870	210	147	183
DOGUE ALLEMAND 2	2420	1200	850	199	153	176
DOGUE ALLEMAND 3	2550	1260	860	214	150	180
DOGUE ALLEMAND 4	2340	1130	830	213	148	170
SETTER 1	2010	1050	700	198	143	159
SETTER 2	1960	1060	670	185	126	142
SETTER 3	2200	1170	700	198	143	156
GROENDAL	2050	1050	700	190	124	149
BAS-ROUGE	2160	1120	750	196	140	164
BRIARD 1	2280	1220	780	225	142	178
BRIARD 2	2120	1110	730	205	137	166
...	...	...	...	...	...	...

12

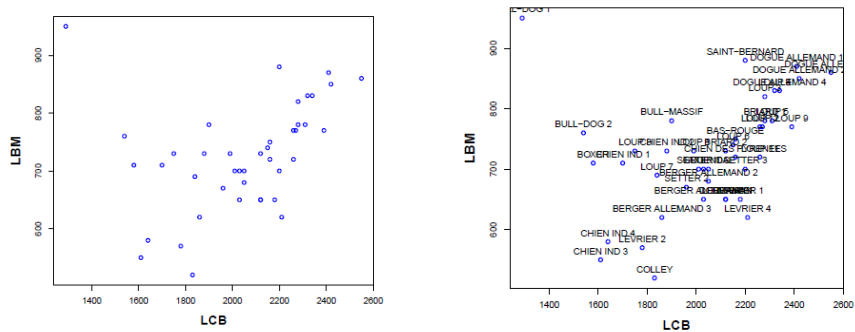
## 2. Pourquoi : veut-on réduire (2) ?



Statistique classique : distribution et fréquences des observations par variables

13

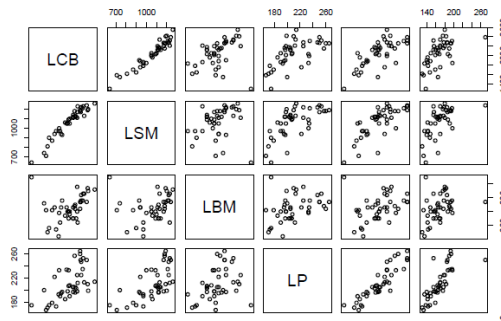
## 2. Pourquoi : veut-on réduire (3) ?



Statistique classique : comparaison de 2 variables : relation entre 2 variables

14

## 2. Pourquoi : veut-on réduire (4) ?



ACP : Décrire au mieux la variabilité des données et doit permettre de :

- Une réduction des données décrites par  $p$  variables à un nombre restreint  $k$  de descripteurs
- Une visualisation des données (si  $k=2$  ou 3)
- Une interprétation des données : liaisons entre variables

Filter l'information (enlever les redondances et le bruit) : compression de données

- ☐ Réduire la dimension, en particulier pour la visualisation des données
- ☐ Déterminer les individus qui se ressemblent ou, au contraire, s'opposent
- ☐ Identifier les variables sur lesquelles sont fondées les ressemblances ou dissemblances (variables discriminantes)
- ☐ Obtenir de nouvelles variables non corrélées (décomposition orthogonale)
- ☐ Faire une analyse de la corrélation entre toutes les variables à la fois

## 3. Analyse de données /statistiques

**Statistique classique** : un nombre restreint de caractères mesurés sur un petit ensemble d'individus

→ développer des méthodes d'estimation et de tests statistiques

**En réalité**

→ Un individu est observé par un grand nombre de caractères

Statistique  
inférentielle

**Analyse de données** : étude globale des individus et des variables avec une représentation graphique

Statistique  
descriptive ou  
explicative

**Différents points de vue de l'analyse :**

- Rechercher des ressemblances ou différences entre individus sur plusieurs caractères (profil) : proximité entre individus : regroupement en catégorie homogène → *Classification*
- Description de relations entre caractères : deux caractères sont liés s'ils varient de la même façon pour les différents individus : si tous les caractères ont un rôle identique : mettre en évidence les groupes de caractères qui sont corrélés ou indépendants : individus et var. dans espace géométrique réduit en perdant le moins d'information possible → **A.D**
- **Différences entre stats classiques et A.D**

16



## 4. Problème de réduction de $dim$ (1)

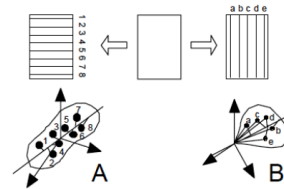
On cherche :

A projeter un ensemble de points  $x_{i,j}$  représenté dans un espace de dimension  $p$  dans un espace  $dim\ k < p$  qui respecte l'agencement des points du nuage :  
ce que l'on cherche → **une projection (à caractériser)**

Pour représenter (non visible au delà de  $p=3$ )

- les individus dans l'espace des variables  
situation A :  $n$ -points-lignes dans  $\mathbb{R}^p$

- les variables dans l'espace des individus  
situation B :  $p$ -points-lignes dans  $\mathbb{R}^n$



→ **choix d'un espace vectoriel sur  $\mathbb{R}^p$  qui sera doté d'un produit scalaire**

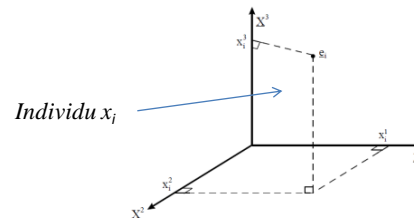
Qui permet de représenter les points individus et les variables dans un même espace

17

### 4.1 Choix : un espace vectoriel euclidien

Propriété de cet espace vectoriel réel  $\mathbb{R}^n$  : il est doté de l'addition, de la multiplication par un scalaire qui sont des opérateurs sur les matrices ;  
une base canonique  $\mathbb{R}^n$  de peut être définie ( $e_1, e_2, \dots, e_n$ ) tout élément  $X$  de  $\mathbb{R}^n$  est une combinaison linéaire des vecteurs de la base

$$\text{Soit } \sum_{i=1}^n x_i e_i$$



On peut définir un produit scalaire  $\Phi$  :

Une application de  $\mathbb{R}^n \times \mathbb{R}^n$  dans  $\mathbb{R}$

Noté  $\langle X|Y \rangle_\Phi$  (noté aussi  $\langle X|\Phi|Y \rangle$ ) le nombre réel  $\Phi(X, Y)$

On peut définir ainsi

un produit scalaire canonique :

un produit scalaire défini par une matrice diagonale à éléments positifs



18

## 4.2 Formalisation du problème

Construire un critère permettant de construire puis d'évaluer la projection :

- Conserver au mieux les caractéristiques du nuage :  
soit la forme du nuage  
qui est liée à la distance entre les individus
- Se doter d'une distance entre individus

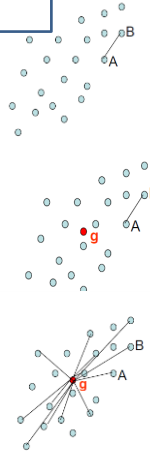
→ Distance euclidienne dans un espace vectoriel sur  $\mathbb{R}^p$

$$d^2(X_A, X_B) = \sum_{j=1}^p (X_A^j - X_B^j)^2$$

→ Caractériser le nuage de points (dans espace des individus)

centre de gravité du nuage :  $g = \left(\frac{1}{n}\right) \sum_{i=1}^n X_i$

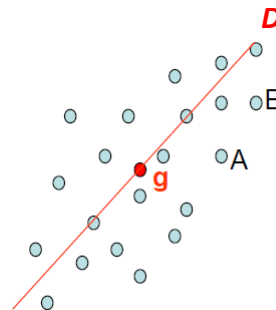
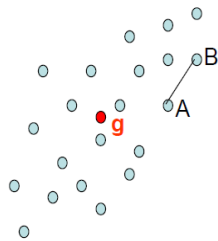
Inertie totale du nuage :  $I_T = \sum_{i=1}^n \frac{1}{n} d^2(g, X_i)$   
 $\rightarrow I_T = \sum_{i=1}^n \left(\frac{1}{n} \sum_{j=1}^p (x_{ij} - g_j)^2\right)$



## 4.3 Formalisation du problème

Définir une projection tel que :

Ces sous-espaces représentant au mieux ce nuage de point en respectant 2 critères : le **critère de proximité** et la **fidélité des distances** (Inertie du nuage et variance).



**Projection sur une droite : projeté orthogonal**

Les données sont fournies par une matrice  $X$   
De  $n$  individus décrits par  $p$  variables

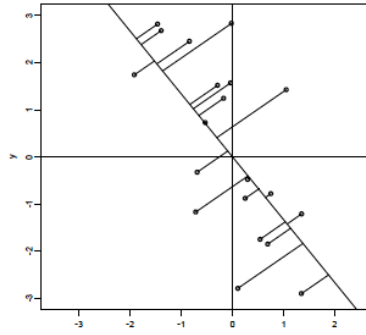
$$\begin{pmatrix} X_{11} & \cdots & X_{1p} \\ \vdots & \ddots & \vdots \\ X_{n1} & \cdots & X_{np} \end{pmatrix}$$

21

## 4.4 Enoncé (1)

Soit deux variables  $X_1$  et  $X_2$  mesurées sur  $n$  individus.

$$\begin{pmatrix} X_{11} & X_{12} \\ \vdots & \vdots \\ X_{n1} & X_{n2} \end{pmatrix}$$



↓  
Deux variables  $X_1$  et  $X_2$  notées :  $X$  et  $Y$  mesurées sur  $n$  individus

L'étude de la liaison entre  $x$  et  $y$  est la recherche d'une droite optimum.  
Quel est le critère d'optimisation ?

22

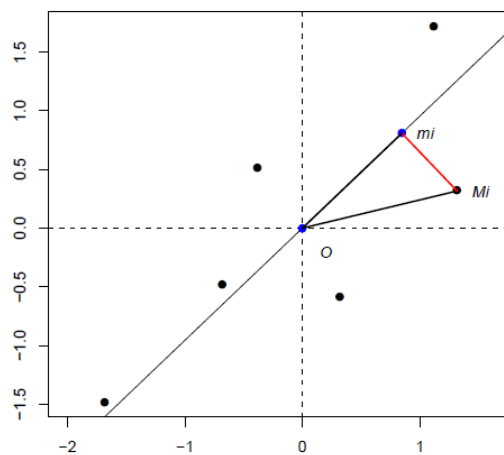
## 4.4 Enoncé (2)

Cette droite minimise

$$\frac{1}{n} \sum_{i=1}^n M_i m_i^2$$

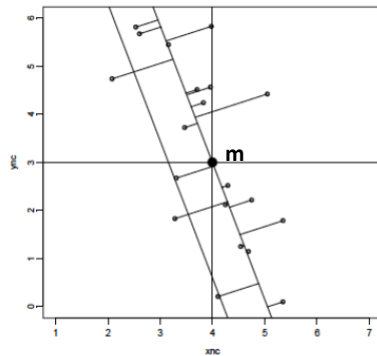
où

- $M_i$  est le point  $i$  du plan,
- $m_i$  est la projection orthogonale de  $M_i$  sur la droite.



**Minimisation de la distorsion entre le point d'origine  $M_i$  et son projeté  $m_i$**   
**Minimisation de la moyenne des carrés des distances des points à cette droite D**

## 4.5 Pourquoi cette droite doit passer par le centre de gravité - point moyen du nuage ?



$(m(x)m(y))$  le point moyen

$D$  et  $D'$  : deux droites parallèles :  
 $D$  passe par le point moyen  $m$  et  
 $D'$  passe par un point quelconque  $m'$

Que vaut :

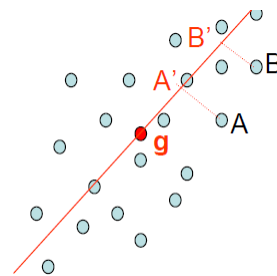
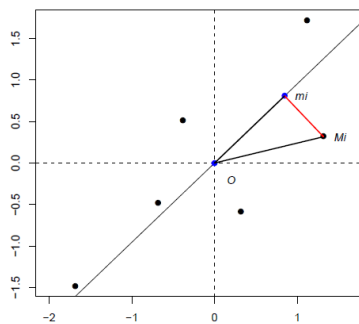
$$\|M_i - m_i'\|^2$$

**A retrouver ...**

24

## 4.6 Soit un problème de minimisation

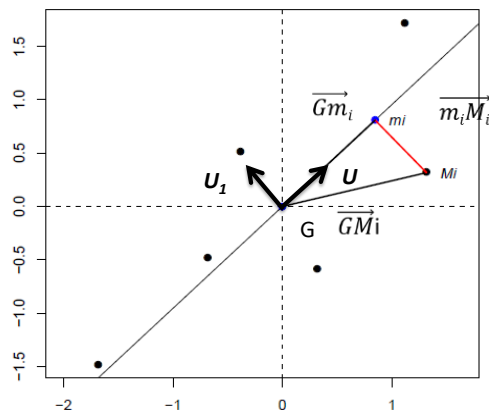
Recherche un sous-espace qui passe par G centre de gravité du nuage (point moyen)



**Minimisation de la distorsion entre le point d'origine  $M_i$  et son projeté  $m_i$  (inertie)**  
**Minimisation de la moyenne des carrés des distances des points à cette droite D (distance)**

25

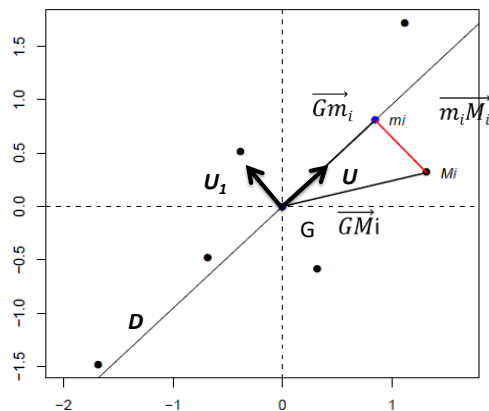
### 4.7 Formulation de la projection sur une droite D



[ACP\\_projection.pdf](#)

27

### 4.7 Formulation de la projection sur une droite D



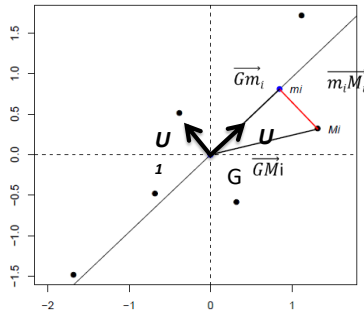
$$u_1 = \begin{pmatrix} \alpha \\ \beta \end{pmatrix}, \text{ avec } \alpha^2 + \beta^2 = 1.$$

$$u = \begin{pmatrix} \beta \\ -\alpha \end{pmatrix}.$$

Un vecteur :  $\alpha X + \beta Y$   
qui est :  
- combinaison linéaire de 2 variables centrées  
- minimise :  $\|\alpha X + \beta Y\|_{D^\perp}^2$

28

## 4.8 Décomposition de l'Inertie du nuage (a)



Avec : Matrice des variables centrées

$$Z = \begin{pmatrix} x_{10} & y_{10} \\ \vdots & \vdots \\ x_{n0} & y_{n0} \end{pmatrix}$$

**Inertie Totale du nuage  $I_T$**

$$I_T = \frac{1}{n} \sum_{i=1}^{i=n} \|\overrightarrow{GM_i}\|^2 = \frac{1}{n} \sum_{i=1}^{i=n} (x_{i0}^2 + y_{i0}^2)$$

**Inertie statistique  $I_s$**

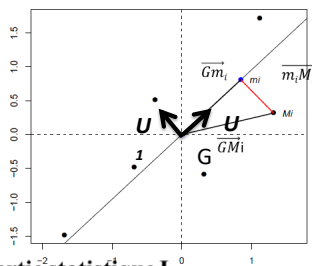
Du nuage de points dans  $R^p$  par rapport à une direction  $\Delta$  de  $R^p$  définie par un vecteur unitaire  $u$ :

$$I_s(u) = \frac{1}{n} \sum_{i=1}^{i=n} \|\overrightarrow{Gm_i}\|^2$$

Avec  $\overrightarrow{Gm_i}$  le projeté orthogonal de  $\overrightarrow{GM_i}$  sur  $D$

29

## 4.8 Décomposition de l'Inertie du nuage (b)



**Inertie Totale du nuage  $I_T$**

$$I_T = \frac{1}{n} \sum_{i=1}^{i=n} \|\overrightarrow{GM_i}\|^2 = \frac{1}{n} \sum_{i=1}^{i=n} (x_{i0}^2 + y_{i0}^2)$$

**Inertie statistique  $I_s$**

Du nuage de points dans  $R^p$  par rapport à une direction  $\Delta$  de  $R^p$  définie par un vecteur unitaire  $u$ : avec  $\overrightarrow{Gm_i}$  le projeté orthogonal de  $\overrightarrow{GM_i}$  sur  $u$

$$I_s(u) = \frac{1}{n} \sum_{i=1}^{i=n} \|\overrightarrow{Gm_i}\|^2$$

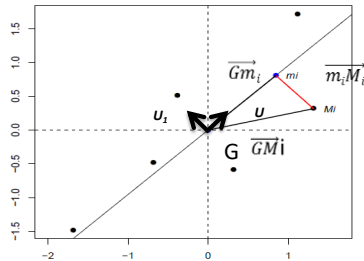
**Inertie Mécanique  $I_M$**

Du nuage de points dans  $R^p$  par rapport à une direction  $\Delta$  de  $R^p$  définie par un vecteur unitaire  $u$ : avec  $\overrightarrow{Gm_i}$  le projeté orthogonal de  $\overrightarrow{GM_i}$  sur  $u$

$$I_M(u) = \frac{1}{n} \sum_{i=1}^{i=n} \|\overrightarrow{m_i M_i}\|^2$$

30

## 4.8 Décomposition de l'Inertie totale (c)



**Inertie Totale du nuage  $I_T$**

$$I_T = \frac{1}{n} \sum_{i=1}^n \|\vec{GM}_i\|^2 = \frac{1}{n} \sum_{i=1}^n (x_{i0}^2 + y_{i0}^2)$$

**Inertie statistique  $I_S$**

Du nuage de points dans  $\mathbb{R}^p$  par rapport à une direction  $\Delta$  de  $\mathbb{R}^p$  définie par un vecteur unitaire  $u$ :

avec  $\vec{Gm}_i$  le projeté orthogonal de  $\vec{GM}_i$  sur  $u$

**Inertie Mécanique  $I_M$**

$$I_S(u) = \frac{1}{n} \sum_{i=1}^n \|\vec{Gm}_i\|^2$$

**Théorème de Pythagore :**

$$\|\vec{GM}_i\|^2 = \|\vec{Gm}_i\|^2 + \|\vec{m}_i \vec{M}_i\|^2$$

$$I_M(u) = \frac{1}{n} \sum_{i=1}^n \|\vec{m}_i \vec{M}_i\|^2$$

$$\underbrace{I_M(u)}_{\text{à minimiser}} = I_T - \underbrace{I_S(u)}_{\text{à maximiser}}$$

**à minimiser**

**à maximiser**

31

## Liens entre géométrie et variables statistiques (dans $\mathbb{R}^2 \rightarrow \mathbb{R}^p$ ) (1)

une variable statistique quantitative bidimensionnelle  $(X, Y)$  à valeurs dans  $\mathbb{R}^2$ , définie dans une population  $\Omega$  de taille  $n$

$\mathbb{R}^2$  est l'**espace des individus**

La variable statistique est représentée par un nuage de points dans  $\mathbb{R}^2$  et chaque point du nuage statistique représente un individu de la population  $\Omega$ .

Les  $n$  valeurs  $X(\omega)$  de  $X$  pour les  $n$  individus de la population peuvent être considérées comme les  $n$  valeurs  $X(\omega)$  de  $X$  pour les  $n$  individus de la population peuvent être considérées comme les coordonnées d'un vecteur de  $\mathbb{R}^n$ .

$$X = \begin{pmatrix} x_1 \\ \vdots \\ x_i \\ \vdots \\ x_n \end{pmatrix}$$

$$Y = \begin{pmatrix} y_1 \\ \vdots \\ y_i \\ \vdots \\ y_n \end{pmatrix}$$

L'espace  $E = \mathbb{R}^n$  : l'**espace des variables**

Chaque élément de  $E$  peut être considéré comme les valeurs d'une variable statistique quantitative réelle définie sur  $\Omega$ .

32

## Liens entre géométrie et variables statistiques (dans $\mathbb{R}^2 \rightarrow \mathbb{R}^p$ ) (2)

### Produit scalaire

Dans cet espace des variables, la matrice  $D = In$ , où  $In$  est la matrice unité à  $n$  lignes et  $n$  colonnes, définit un **produit scalaire** :

$$\langle X | Y \rangle_{D \frac{1}{n}} = \langle X | D \frac{1}{n} | Y \rangle = \sum_i \frac{1}{n} x_i y_i = \frac{1}{n} \sum_i x_i y_i = \frac{1}{n} \langle X | Y \rangle$$

### Moyenne d'une variable statistique.

La moyenne de la variable statistique  $X$  :

$$\bar{X} = \frac{1}{n} \sum_{\omega} X(\omega) = \frac{1}{n} \sum_i x_i = \frac{1}{n} \sum_i x_i \times 1 = \langle X | D \frac{1}{n} | \mathbf{1}_n \rangle = \langle X | \mathbf{1}_n \rangle_{D \frac{1}{n}}$$

$$\mathbf{1}_n = \begin{pmatrix} 1 \\ \vdots \\ 1 \\ \vdots \\ 1 \end{pmatrix}$$

Soit  $X_0$  la variable centrée correspondant à  $X$  : pour chaque individu  $\omega$  de population, sa valeur est  $X(\omega) - \bar{X}$  :

$$X_0 = \begin{pmatrix} x_1 - \bar{X} \\ \vdots \\ x_i - \bar{X} \\ \vdots \\ x_n - \bar{X} \end{pmatrix} = \begin{pmatrix} x_1 \\ \vdots \\ x_i \\ \vdots \\ x_n \end{pmatrix} - \bar{X} \begin{pmatrix} 1 \\ \vdots \\ 1 \\ \vdots \\ 1 \end{pmatrix} = X - \bar{X} \mathbf{1}_n$$

$$X = X_0 + \bar{X} \mathbf{1}_n = X_0 + \langle X | \mathbf{1}_n \rangle_{D \frac{1}{n}} \mathbf{1}_n$$

33

## Liens entre géométrie et variables statistiques (dans $\mathbb{R}^2 \rightarrow \mathbb{R}^p$ ) (3)

### Variance d'une variable statistique

$$s^2(X) = \overline{X_0^2} = \frac{1}{n} \sum_i (x_i - \bar{X})^2 = \langle X_0 | D \frac{1}{n} | X_0 \rangle = \|X_0\|^2$$

$$s^2(X) = \|X_0\|^2$$

### Covariance

$$\text{Cov}(X, Y) = \frac{1}{n} \sum_i (x_i - \bar{X})(y_i - \bar{Y}) = \langle X_0 | D \frac{1}{n} | Y_0 \rangle = \langle X_0 | Y_0 \rangle_{D \frac{1}{n}}$$

$$\text{Cov}(X, Y) = \langle X_0 | D \frac{1}{n} | Y_0 \rangle = \langle X_0 | Y_0 \rangle_{D \frac{1}{n}}$$

### Coefficient de corrélation linéaire

$$r_{XY} = \frac{\text{Cov}(X, Y)}{s(X) s(Y)} = \frac{\langle X_0 | D \frac{1}{n} | Y_0 \rangle}{\|X_0\|_{\Phi} \|Y_0\|_{\Phi}} = \cos(X_0, Y_0)$$

$$r_{XY} = \cos(X_0, Y_0)$$

34



## 4.9 Inertie statistique (= inertie projetée)

Pour le vecteur unitaire  $u$  qui définit l'axe factoriel D  $u = \begin{pmatrix} \beta \\ -\alpha \end{pmatrix}$ .

**Inertie statistique est :**  $I_S(u) = \frac{1}{n} \sum_{i=1}^{i=n} \|\vec{Gm}_i\|^2$

avec  $\vec{Gm}_i = \langle \vec{Gm}_i | u \rangle u = (\beta x_{i0} - \alpha y_{i0}) \begin{pmatrix} \beta \\ -\alpha \end{pmatrix}$

$$I_S(u) = \frac{1}{n} \sum_{i=1}^{i=n} (\beta x_{i0} - \alpha y_{i0})^2 = \beta^2 \times \frac{1}{n} \sum_{i=1}^{i=n} x_{i0}^2 + \alpha^2 \times \frac{1}{n} \sum_{i=1}^{i=n} y_{i0}^2 - 2 \alpha \beta \times \frac{1}{n} \sum_{i=1}^{i=n} x_{i0} y_{i0}$$

**OR :**  $\frac{1}{n} \sum_{i=1}^{i=n} x_{i0}^2 = s^2(X), \frac{1}{n} \sum_{i=1}^{i=n} y_{i0}^2 = s^2(Y), \frac{1}{n} \sum_{i=1}^{i=n} x_{i0} y_{i0} = \text{Cov}(X, Y),$

$$\begin{aligned} I_S(u) &= \beta^2 s^2(X) + \alpha^2 s^2(Y) - 2 \alpha \beta \text{Cov}(X, Y) \\ &= (\beta - \alpha) \begin{pmatrix} s^2(X) & \text{Cov}(X, Y) \\ \text{Cov}(X, Y) & s^2(Y) \end{pmatrix} \begin{pmatrix} \beta \\ -\alpha \end{pmatrix} \\ &= {}^t u A u \end{aligned}$$

35

## 4.10. Inertie statistique $I_S$ et matrice de variance-covariance $A$ et recherche du vecteur $u$

$$I_S(u) = \frac{1}{n} \sum_{i=1}^{i=n} (\beta x_{i0} - \alpha y_{i0})^2 = \beta^2 \times \frac{1}{n} \sum_{i=1}^{i=n} x_{i0}^2 + \alpha^2 \times \frac{1}{n} \sum_{i=1}^{i=n} y_{i0}^2 - 2 \alpha \beta \times \frac{1}{n} \sum_{i=1}^{i=n} x_{i0} y_{i0}$$

**Et :**

$$\begin{aligned} I_S(u) &= \beta^2 s^2(X) + \alpha^2 s^2(Y) - 2 \alpha \beta \text{Cov}(X, Y) \\ &= (\beta - \alpha) \begin{pmatrix} s^2(X) & \text{Cov}(X, Y) \\ \text{Cov}(X, Y) & s^2(Y) \end{pmatrix} \begin{pmatrix} \beta \\ -\alpha \end{pmatrix} \\ &= {}^t u A u \end{aligned}$$

36

## 4.11 Zoom sur la matrice de covariance A

### La matrice de variance-covariance

$$A = \begin{pmatrix} s^2(X) & \text{Cov}(X, Y) \\ \text{Cov}(X, Y) & s^2(Y) \end{pmatrix}$$

$$A = \frac{1}{n} {}^tZ Z = {}^tZ D \frac{1}{n} Z$$

Matrice symétrique réelle qui si est diagonalisable sera décomposable



Par :  $P^{-1} A P = D \Leftrightarrow A = P D P^{-1}$  et  $A P = P D$

Pour lesquels il existe des valeurs propres  $\lambda$  et vecteurs propres  $V$  associés avec

$$A V = \lambda V \Leftrightarrow (A - \lambda I) V = 0$$

Donc un polynôme caractéristique :  $\det((A - \lambda I)) = 0$

Qui représente (application linéaire de  $\mathbb{R}^2$  dans  $\mathbb{R}^2$  dont le noyau de l'endomorphisme est défini par la matrice  $A - \lambda I$ ).

37

## 4.11 Valeurs propres et vecteurs propres de A (1)

Les valeurs propres sont solution de :

$$\begin{aligned} \det((A - \lambda I_2)) &= 0 \\ \lambda^2 - (s^2(X) + s^2(Y)) \lambda + s^2(X) s^2(Y) - (\text{Cov}(X, Y))^2 &= 0 \\ P_A(\lambda) &= (-1)^d \lambda^d + (-1)^{d-1} \text{tr}(A) \lambda^{d-1} + \det(A) \end{aligned}$$

Le discriminant de cette équation est :

$$(s^2(X) + s^2(Y))^2 - 4(s^2(X) s^2(Y) - (\text{Cov}(X, Y))^2) = (s^2(X) - s^2(Y))^2 + 4(\text{Cov}(X, Y))^2 \geq 0$$

Le polynôme caractéristiques (pour  $n=2$ ) a 2 racines donc 2 valeurs propres distinctes ( $\lambda_1$  et  $\lambda_2$ ) avec :

- la somme des valeurs propres est la **trace** de la matrice, somme des éléments de la première diagonale  $\lambda_1 + \lambda_2 = s^2(X) + s^2(Y) \geq 0$
- le produit de ces valeurs propres est le **déterminant** de la matrice

$$\lambda_1 \lambda_2 = s^2(X) s^2(Y) - (\text{Cov}(X, Y))^2 \geq 0$$

38

## 4.12 Valeurs propres et vecteurs propres de A (2)

$\lambda_1$  et  $\lambda_2$  les **valeurs propres de la matrice des variances - covariances**,  
rangées par ordre décroissant :

$$\lambda_1 > \lambda_2 > 0$$

$$\lambda_1 = \frac{1}{2} \left( s^2(X) + s^2(Y) + \sqrt{[s^2(X) - s^2(Y)]^2 + 4[\text{Cov}(X, Y)]^2} \right)$$

$$\lambda_2 = \frac{1}{2} \left( s^2(X) + s^2(Y) - \sqrt{[s^2(X) - s^2(Y)]^2 + 4[\text{Cov}(X, Y)]^2} \right)$$

$\mathbb{R}^2$  possède une **base propre orthonormée**, une base  $\{u_1, u_2\}$  orthonormée

$$A u_1 = \lambda_1 u_1 \text{ et } A u_2 = \lambda_2 u_2$$

et  $\|u_1\|^2 = 1, \|u_2\|^2 = 1, \langle u_1 | u_2 \rangle = 0$ .

Pour une valeur propre  $\lambda$  : on peut trouver un vecteur propre et un vecteur propre normé :

$$u = \frac{1}{\sqrt{[s^2(Y) - \lambda]^2 + [\text{Cov}(X, Y)]^2}} \begin{pmatrix} s^2(Y) - \lambda \\ -\text{Cov}(X, Y) \end{pmatrix}$$

39

## 4.11 Valeurs propres et vecteurs propres de A (3)

Pour une valeur propre  $\lambda$  : on peut trouver un vecteur propre normé :

$$u = \frac{1}{\sqrt{[s^2(Y) - \lambda]^2 + [\text{Cov}(X, Y)]^2}} \begin{pmatrix} s^2(Y) - \lambda \\ -\text{Cov}(X, Y) \end{pmatrix}$$

Et définir deux vecteurs  $u_1$  et  $u_2$  formant une base orthonormée dans  $\mathbb{R}^2$  :

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

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

La matrice  $V$  formée par les deux vecteurs propres  $u_1$  et  $u_2$  est une matrice orthogonale ou :  $V^{-1} = {}^tV$

On peut alors obtenir :

$$\begin{cases} \Lambda = V A {}^tV \\ A = {}^tV \Lambda V \end{cases} \quad (\text{diagonalisation de la matrice } A)_{40}$$

## 4.12 Recherche des axes principaux (1)

Posons, pour un vecteur normé  $u$ ,  $v = V u$

Dans  $R^2$ , rapporté à la base  $\{u_1; u_2\}$  notons  $v = \begin{pmatrix} v_1 \\ v_2 \end{pmatrix}$

On a :  ${}^t v = {}^t u {}^t V$

$$\|v\|^2 = {}^t v v = {}^t u {}^t V V u = {}^t u u = \|u\|^2 = 1$$

$u$  est normé et  $v$  aussi

$$Is(u) = {}^t u A u = {}^t u {}^t V \Lambda V u = {}^t v \Lambda v$$

$$Is(u) = {}^t v \Lambda v = \begin{pmatrix} v_1 & v_2 \end{pmatrix} \begin{pmatrix} \lambda_1 & 0 \\ 0 & \lambda_2 \end{pmatrix} \begin{pmatrix} v_1 \\ v_2 \end{pmatrix} = \lambda_1 v_1^2 + \lambda_2 v_2^2$$

**Recherche de la droite D orthogonale se ramène à**

**Maximiser  $\lambda_1 v_1^2 + \lambda_2 v_2^2$  sous la contrainte  $v_1^2 + v_2^2 = 1$  avec  $\lambda_1 > \lambda_2 > 0$**

41

## 4.12 Recherche des axes principaux (2)

**Recherche de la droite D orthogonale :**

**Maximiser  $\lambda_1 v_1^2 + \lambda_2 v_2^2$  sous la contrainte  $v_1^2 + v_2^2 = 1$  avec  $\lambda_1 > \lambda_2 > 0$**

$$\begin{aligned} Is(u) &= \lambda_1 v_1^2 + \lambda_2 v_2^2 \\ &= \lambda_1 (1 - v_2^2) + \lambda_2 v_2^2 \\ &= \lambda_1 - (\lambda_1 - \lambda_2) v_2^2 \end{aligned}$$

**La valeur  $\lambda_1 - (\lambda_1 - \lambda_2) v_2^2$  avec  $\lambda_1 > \lambda_2$  atteint son maximum  $\lambda_1$  lorsque l'on prend  $v_2 = 0$  donc  $|v_1| = 1$**

La direction du premier axe factoriel est prise par le vecteur  $v$  de coordonnées  $\begin{pmatrix} 1 \\ 0 \end{pmatrix}$  dans la base  $(u_1; u_2)$  (vecteurs propres de  $A$ ):  $v = u_1$

Et  $Is(u_1) = \lambda_1$

**Théorème :** la direction du premier axe factoriel est définie par le vecteur propre associé à la plus grande valeur propre de la matrice de variances-covariances

42

## 4.12 Recherche des axes principaux (3)

### Recherche des autres axes factoriels :

**Corollaire :** la direction perpendiculaire au premier axe factoriel définit le deuxième axe factoriel il est défini par le deuxième vecteur propre associé à la deuxième valeur propre de la matrice variance-covariance

La direction du deuxième axe factoriel est prise par le vecteur  $v$  de coordonnées  $\begin{pmatrix} 0 \\ 1 \end{pmatrix}$  dans la base  $(u_1; u_2)$  (vecteurs propres de  $A$ ):  $v = u_2$

Et  $Is(u_2) = \lambda_2$

La somme des valeurs propres est égale à la somme des variances de chaque variable = trace de la matrice de variance -covariance  $A$

$$I_T = Is(u_1) + Is(u_2)$$

Pour dimension  $p$  :

Le premier vecteur propre normé  $u_1$  est un vecteur de  $\mathbb{R}^p$  qui maximise l'inertie projetée

Le 2<sup>ème</sup> vecteur propre normé  $u_2$  est un vecteur de  $\mathbb{R}^p$  orthogonal à  $u_1$  qui maximise à nouveau l'inertie projetée. Ainsi de suite pour  $u_3 \dots u_r$  axes suivants où  $r$  représente le rang de la matrice diagonalisée

43

## 4.13 En résumé l' ACP est basée sur:

- Une matrice de covariance/ corrélation
- Une décomposition de la matrice  $A$  en :
  - une matrice de valeurs propres ordonnées et
  - une matrice de vecteurs propres associés, ordonnés selon un critère de conservation maximale de l'inertie du nuage
  - Une visualisation dans un espace  $\mathbb{R}^k$   $k \ll p$

Lors de cette réduction de dimension : choix de la dimension  $k \ll p$  ?

44

## 5. Choix des $k$ premières composantes principales ( $\dim < p$ )

**Choix des  $r$  premières composantes principales**

$k \ll p$  réduction de la dimension

objectif : garder un maximum d'information des données initiales

**Mesure de cette information : le % de variance expliquée**

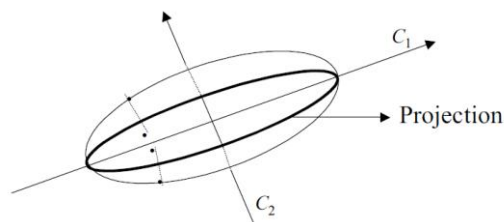
**Pourquoi ce critère :**

- La conservation de la variance initiale du nuage comme hypothèse de la décomposition
- Mais une interprétation géométrique →

Géométriquement : revient à projeter les données dans un sous-espace de dimension  $k$ , centré sur  $g$ , reprenant les  $k$  premiers axes principaux d'allongement du nuage

45

### 5.1 Exemple : données initiales à 3 dimensions distribuées dans un 'ballon de rugby'



Plus le nuage est aplati sur  $C_1, C_2 \Rightarrow$  moins de variance sur la 3<sup>ème</sup> dimension.  
 $\Rightarrow$  % de variance expliquée par  $C_1, C_2$

En général :

- Le % de variance expliquée par  $C_1, C_2, \dots, C_r =$  mesure d'aplatissement du nuage sur le sous-espace des composantes (à  $k$  dim.).

Plus ce % est grand, meilleure est la représentation des données dans le sous-espace

- Les composantes principales sont entièrement déterminées par la matrice **A**  
**variance-covariance**

$\Rightarrow$  toute modification de  $A \rightarrow$  **modification des composantes**

46

## 5.2 Axes factoriels et taux d'inertie expliquée

Le taux d'inertie expliqué par chaque composante  $i$

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

Le taux d'inertie expliqué par  $k$  composantes ( soit par l'hyperplan de dimension  $k$  )

$$PI_{1..k} = \frac{\sum_{i=1}^k \lambda_i}{\sum_{i=1}^p \lambda_i}$$

$\frac{\lambda_1 + \lambda_2}{\sum_{i=1}^p \lambda_i}$  est la part d'inertie expliquée par le premier plan principal.

47

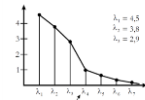
## 5.3 Différentes règles existes (a)

1 ) Un seuil de pourcentage d'inertie est fourni

Composante	Valeur propre = Eigenvalue	Pourcentage de variance	Pourcentage de variance cumulée
C1	3.91217	47.3	47.3
C2	2.61602	31.4	78.7
C3	.57780	8.3	87.0
C4	.23840	3.4	90.4
C5	.13528	1.9	92.3
C6	.08297	1.2	93.5
C7	.03878	.5	99.0
Total :	7.00000	100.00	

- L'inertie cumulée permet de choisir par rapport à un seuil fourni entre 80 à 90%
  - Souffre du caractère a priori

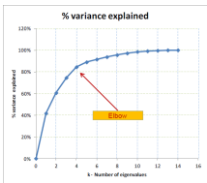
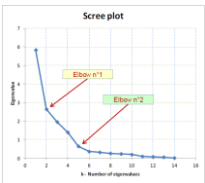
Construire l'histogramme des inerties



2) la courbe de décroissance des valeurs propres  $\lambda_k$  ( règle de Cattell ).

L'idée est de détecter les « coudes » (les « cassures ») signalant un changement de structure :

- intéressante car permet de dépasser l'arbitraire purement numérique.
- mais compliquée à mettre en œuvre, car soumise à l'appréciation



Combinaisons des deux indicateurs pour choisir ( ici  $k=4$  ):

{Scree plot + % explained variance}

## 5.3 Différentes règles existes (b)

### 3 ) Règle de Kaiser - Guttman

pour ACP normée : données initiales sont centrées et réduites

ACP normée, la somme des valeurs propres est égale au nombre de variables

Un axe est retenu si sa valeur propre est  $\geq 1$

Kaiser-Meyer-Olkin			
Valeur			
1	0.882116	81.76%	47.76%
Axes	Eigen value	% explained	% cumulated
1	2.440156	18.86%	80.56%
2	1.389336	10.95%	74.51%
3	1.389336	10.95%	64.41%
4	0.514465	4.33%	88.94%
5	0.353217	2.82%	91.47%
6	0.310082	2.21%	93.68%
7	0.252763	1.81%	95.49%
8	0.228203	1.63%	97.12%
9	0.189241	1.39%	98.47%
10	0.092301	0.66%	99.13%
11	0.047970	0.34%	99.96%
12	0.009091	0.04%	100.00%
13			
14			
Tot.	14		

### 4) Règle de Karlis - Saporta - Spinaki

$$\lambda > 2 \sqrt{\frac{p-1}{n-1}}$$

Seuil = la moyenne des valeurs  $\lambda$  propres + 2 fois leur écart-type

Une règle plus restrictive

n	47		
p	14		
KARLIS et al			
critical value			
2.063			
Kaiser-Meyer-Olkin			
Valeur			
1	0.882116	41.70%	41.70%
2	2.440156	18.86%	60.56%
3	1.351464	13.95%	74.51%
4	1.385435	9.95%	84.46%
5	0.634400	4.53%	88.94%
6	0.353217	2.52%	91.47%
7	0.310082	2.21%	93.68%
8	0.252763	1.81%	95.49%
9	0.228203	1.63%	97.12%
10	0.189341	1.32%	98.47%
11	0.092301	0.66%	99.13%
12	0.047970	0.49%	99.96%
13	0.047970	0.34%	99.96%
14	0.005091	0.04%	100.00%

## 6. ACP → construire de nouvelles variables (a)

Une projection : construire de nouvelles variables

### Recherche des p composantes principales

Composantes :  $C_1, C_2, \dots, C_k, \dots, C_p$

$C_k$  = nouvelle variable = combinaison linéaire des variables d'origine  $X_1, \dots, X_p$ :

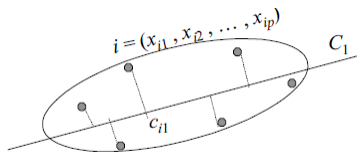
$C_k = a_{1k} X_1 + a_{2k} X_2 + \dots + a_{pk} X_p$  coefficients  $a_{jk}$  à déterminer

tel que les  $C_k$  soient:

- 2 à 2 non corrélées,
- de variance maximale,
- d'importance décroissante.

$C_1$  = 1ère composante principale doit être de variance maximale

Géométriquement :  $C_1$  détermine une nouvelle direction dans le nuage de points qui suit l'axe d'allongement (étirement) maximal du nuage.



$c_{i1}$  = coordonnée du point  $i$  sur l'axe  $C_1$   
projection de  $\mathbf{x}_i$  sur  $C_1$

$$c_{i1} = \sum_{j=1}^p a_{1j} x_{ij}$$

$C_1$  de variance maximale les projections  $c_{i1}$  sont les plus dispersées possible.

50



## 6. ACP → construire de nouvelles variables (b)

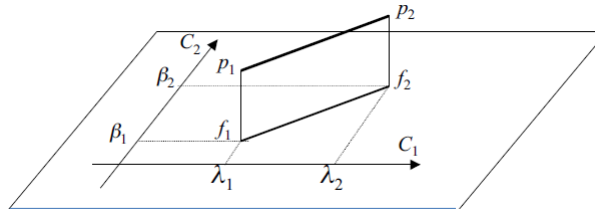
$C_1$  = droite passant par le centre de gravité réalisant le meilleur ajustement possible du nuage c-à-d : qui conserve au mieux la distance entre les points (après projection)

=> droite de projection assurant une distorsion minimale.

$C_2$  = 2ème composante, orthogonale à  $C_1$  et de variance maximale.

Géométriquement :  $C_2$  détermine une droite perpendiculaire à  $C_1$  (au point  $g$ ), suivant un axe (perpendiculaire au 1er) d'allongement maximum.

=>  $C_1$  et  $C_2$  déterminent le plan principal : le meilleur plan de projection (de distorsion minimum).



$C_1$  est telle que la moyenne des  $d_2(\lambda_i, \lambda'_i)$  max.

$C_2$  est  $\perp$  à  $C_1$  et telle que la moyenne des  $d_2(\beta_i, \beta'_i)$  max.

=>  $C_1$  et  $C_2$  déterminent le plan tel que  $d_2(f_i, f'_i)$  soit maximum.

=>  $C_3$  est la droite  $\perp$  à  $C_1$  et  $C_2$  (par  $g$ ) telle que la variance des coord. soit maximum ...

51

## 7. Quelques termes de l'ACP et notations

Nouvelles variables seront appelées : « **composantes principales** » : max p  
A chaque axe est associé 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.

$$C_k = a_{1k} X_1 + a_{2k} X_2 + \dots + a_{pk} X_p \text{ coefficients } a_{jk}$$

les axes qu'elles déterminent : « **axes principaux** » : il y a p axes, de direction principale donnée par le vecteur propre associé: le premier est associé à l'axe principal de direction donnée par  $u_1$  et associé à la plus grande valeur propre  $\lambda_1$ , ...

les formes linéaires associées : « **facteurs principaux** »

52

## 8.1 Projection des individus selon leurs coordonnées - les composantes principales

La  $k$  ème composante principale  $c_k$  fournit les coordonnées des  $n$  individus sur le  $k$  ème axe principal. Si on désire une représentation plane des individus, la meilleure sera celle réalisée grâce aux deux premières composantes principales.

Si  $u_k$  est le vecteur propre de rang  $k$ , les coordonnées des projections des  $n$  points sont obtenus :

$$l_k = \begin{bmatrix} l_{1k} \\ \vdots \\ l_{nk} \end{bmatrix} = \begin{bmatrix} \langle M_1 | u_k \rangle \\ \vdots \\ \langle M_n | u_k \rangle \end{bmatrix} = \begin{bmatrix} \sum_{j=1}^p (x_{1j} - \bar{x}_j) u_{jk} \\ \vdots \\ \sum_{j=1}^p (x_{nj} - \bar{x}_j) u_{jk} \end{bmatrix}$$

soit en écriture matricielle :  $l_k = Z u_k$

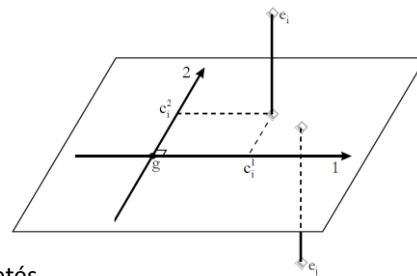
$V$  : matrice de vecteurs propres

(chaque colonne est un vecteur propre  $u_k$ )

$Z$  matrice des coordonnées des variables centrées

$L_i^k$  : coordonnées de l'individu  $M_i$  sur la composante  $C_k$  ( $X$ : matrice des données dans  $R^p$ )

$l_k$  axe principal composé des  $n$  individus projetés



## 8. 2 Axes principaux et composante principale $l_k$

$u_k$  est appelé **axe principal** de rang  $k$ .

$l_k$  est appelé vecteur des **coordonnées** sur l'axe principal.

C'est une variable artificielle de moyenne nulle et de variance  $\lambda_k$ .

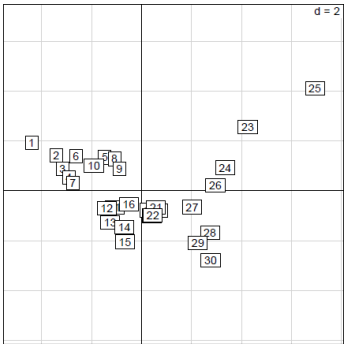
$$\begin{aligned} m(l_k) &= \frac{1}{n} \sum_{i=1}^n l_{ik} = \frac{1}{n} \sum_{i=1}^n \sum_{j=1}^p (x_{ij} - \bar{x}_j) u_{jk} \\ &= \sum_{j=1}^p u_{jk} \frac{1}{n} \sum_{i=1}^n (x_{ij} - \bar{x}_j) = 0 \end{aligned}$$

$$\begin{aligned} v(l_k) &= \frac{1}{n} \sum_{i=1}^n l_{ik}^2 = \frac{1}{n} (\mathbf{X} u_k)^T \mathbf{X} u_k = u_k^T \mathbf{C} u_k \\ &= \lambda_k u_k^T u_k = \lambda_k \end{aligned}$$

### 8.3 individus projetés sur couple d’axes principaux

La représentation du nuage projeté sur un couple d'axes principaux est appelée carte factorielle. C'est une manière de voir l'information multidimensionnelle.

La carte factorielle des axes 1 et 2 est dite premier plan factoriel et représente la part maximale de la variabilité. Chaque point  $i$  est positionné par ses deux coordonnées  $(l_{i1}, l_{i2})$ .



55

### 9.1 Qualité de la projection des individus : representation des individus

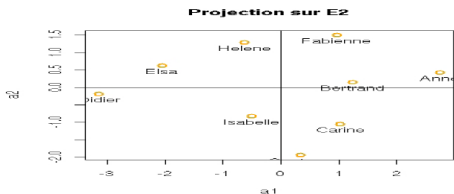
Qualité de la représentation de l'individu  $i$  sur l'espace  $E$  à  $k$  dimensions :

$$Q_{ik} = \frac{\sum_{j=1}^k l_{ij}^2}{\sum_{j=1}^p l_{ij}^2}$$

**Remarque :** c'est le cosinus carré de l'angle entre  $x_i$  et sa projection sur l'axe  $j$ .

Qualité de la représentation sur  $E_k$

Contribution de l'inertie d'un axe par un individu



56

### 9.2 Contribution d'un individu i à une composante j

La contribution d'un individu i , observation, quantifie l'importance de i dans la définition d'un vecteur propre j

$$ct_{i\ j} = \frac{n^{-1}(l_i^j)^2}{\lambda_j}$$

57

### 9.3 Variables et espace réduit R<sup>p</sup>

2/ Interprétation des composantes principales

corrélations avec les variables initiales

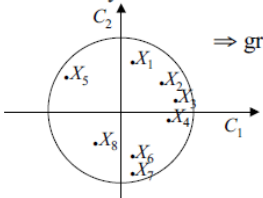
	C <sub>1</sub>	C <sub>2</sub>	C <sub>3</sub>	...
X <sub>1</sub>	r <sub>11</sub>	r <sub>12</sub>	r <sub>13</sub>	...
X <sub>2</sub>	r <sub>21</sub>	r <sub>22</sub>	r <sub>23</sub>	...
⋮	⋮	⋮	⋮	...
X <sub>p</sub>	r <sub>p1</sub>	r <sub>p2</sub>	r <sub>p3</sub>	...

repérer les variables très corrélées  
( r ≈ 1 ou r ≈ -1 )

Interprétation des 2 premières composantes C<sub>1</sub> , C<sub>2</sub> : cercle des corrélations :

C<sub>1</sub> et C<sub>2</sub> étant non-corrélées, on a  $r^2 ( c_1, x_j ) + r^2 ( c_2, x_j ) \leq 1$

=> chaque variable représentée par les coordonnées : ( r ( c<sub>1</sub>, x<sub>j</sub> ) , r ( c<sub>2</sub>, x<sub>j</sub> ) ) est dans un cercle de rayon 1



=> groupes de variables liées ou opposées

⚠ si proches de la circonférence, bien représentées par les 2 composantes !

## 9.4 Qualité et Contribution d'un variable à une composante

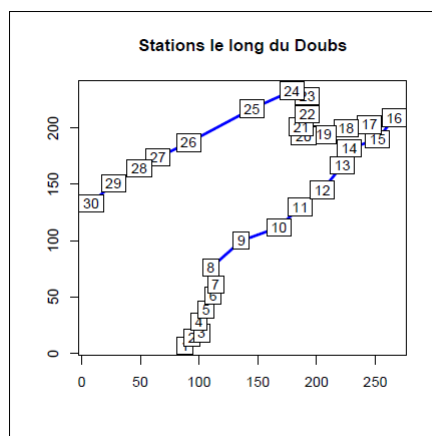
**La qualité la représentation** de la variable  $i$  sur un axe  $k$  est égal à la corrélation entre la variable  $i$  et le vecteur propre  $k$   
 Pour une ACP réduite = au carré de la coordonnée

$$Q_{vik} = (c_{vik})^2 / (l_{vk})^2$$

**La contribution d'une variable**  $i$  à un axe  $k$  est l'importance dans la définition du vecteur propre associé à l'axe  $k$  et de valeur propre  $\lambda_k$   
 La contribution d'une variable  $i$  est égale à la valeur au carré des éléments des vecteurs propres

59

## 10. Un exemple (1): données initiales

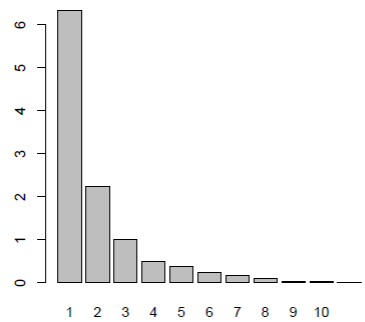


11 variables physico-chimiques :  
 distance à la source, altitude,  
 pente  
 débit, pH, dureté de l'eau  
 phosphate, nitrate, ammoniacque  
 oxygène, demande biologique en  
 oxygène

60

# 10. Un exemple(2): valeurs propres

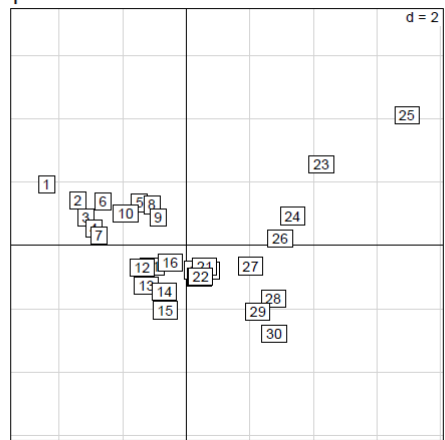
Graphe des valeurs propres



61

# 10. Un exemple(3): projections des individus sur les 2 axes factoriels

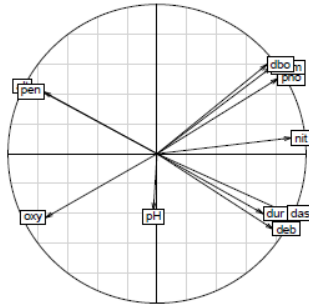
Représentation des individus - stations



62

## 10. Cas de données normées corrélation entre une variable $X_j$ et une composante $k$ ...

Représentation des variables physico-chimiques



63

## 11. Interprétation : un fil conducteur

- Etapes d'interprétation: **Considérer l'analyse dans son ensemble**  
Soit répondre aux questions suivantes :
  - A-t-on pu réduire la dimension du problème ?
  - Combien de vecteurs propres sont importants ?
  - Pourcentage de l'inertie initiale retenue ?
  - Quelles sont les observations( variables) qui contribuent le plus à définir une composante ?
  - Quelles sont les observations( variables) qui sont bien représentées ?
  - Quels sont les signes affectant les coordonnées des observations (variables)
  - Quelle est la position des différents facteurs sur les différents plans
  - Y a t il des regroupements des tendances au niveau des observations si oui quelles en sont les explications ?
  - Identifier de individus bizarres (*un point extrême , écarté du nuage de points peut avoir une forte contribution*) :individu standard ou à exclure ?

*Quelque recommandation dans la partie visualisation : les distances sur les plans factoriels sont des distances, projetées et après réduction !!*

64

## 11. Interprétation : quelques repères

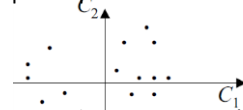
**Représentation** observations (variables) n'est valable que si le % de variance expliquée par  $C_1, \dots, C_k$  est suffisamment grand

- vérifier si les proximités se maintiennent dans d'autres plans de projection:

$C_1 - C_3, C_2 - C_3, \dots$

- individus les mieux représentés: points proches du plan (projection peu importante des groupes d'individus présentant des similitudes:

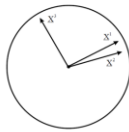
*Attention* aux proximités abusives dues aux projections



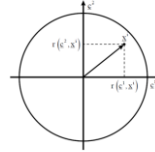
### Représentation des variables

Les « proximités » entre les composantes principales et les variables initiales sont mesurées par les covariances, et surtout les corrélations,  $r(C_j, x_i)$  coefficient de corrélation linéaire entre la composante  $C_j$  et l'individu  $Z_i$ .

$r(C_j, Z_i)$  est le coefficient de corrélation linéaire  $C_j$  et l'individu  $Z_i$ .



$X^1$  et  $X^2$  ont une corrélation proche de 1.  
 $X^1$  et  $X^3$  ont une corrélation proche de 0.



CERCLE DES CORRÉLATIONS

65

