

ANALYSE DES DONNEES  
Master MSS semestre1  
Année universitaire2021-2022

A.RAHMOUNE  
Département de Maths  
Faculté des sciences

Université de Boumerdes M'hamed Bougara  
UMBB

Novembre 2021

# Support du cours

## Chapitre 1

On distingue deux grandes parties complémentaires en Analyse des données (AdD):

**I: Analyse Factorielle** (en abréviation **AF**): C'est une méthode de **réductions** des données, projections sur des axes bien choisis (Axes factoriels)

**II. Classification automatique (CA)**: C'est une méthode de **groupement** de données par classes (arboréscence)

**Les différentes méthodes de l'analyse factorielle sont:**

1. Analyse Factorielle générale (AFG) d'un nuage de points quelconques: La démarche théorique de toutes analyses.
2. Analyse en composantes principales (ACP) (On distingue 3 variantes)
3. Analyse factorielle des correspondances (AFC): La plus puissante dans les applications, et la plus intéressante du point de vue logique.
4. Analyse des correspondances multiples (ACM), généralise en quelque sorte de l'AFC.
5. Analyse factorielle discriminante (AFD).
6. Analyse Canonique simple (ACS)- Utilité plus que théorique que pratique.
7. Analyse canonique généralisée (ACG): La plus généraliste, de cette méthode on peut retrouver toutes les méthodes citées plus haut.

**Concernant les méthodes des CA, on peut citer:**

1. Les Méthodes Hierarchiques Ascendantes et descendantes(CHA, CHD).
2. Les méthodes de partitionnements: Méthode des nuées dynamique (Méthode des centres mobiles)
3. Les méthodes classées dans un espace métrique quelconques.
4. Les méthodes classées à partir d'observations qualitatives.

**Remarques:**

Les mot clés de **I: Analyse Factorielle ( AF) sont:**

Le nuage de points (individus, variables), les masses affectées aux points (pondération); la métrique, pour les << **inputs**>>

Pour les << **outputs**>> sont: Les axes d'inertie-Axes factorielles-,les coordonnées des points sur ces axes-composantes principales-, et divers indications aidants à l'interprétation.

D'une méthode d'analyse factorielle à l'autre, seuls varient les <<inputs>>

Le but et la démarche essentielle de l'Analyse des Données est:

**Consentir une perte en information<sup>\*</sup> afin d'obtenir un gain en signification<sup>\*\*</sup>.**

Illustrons cette idée par un exemple:

Dans une étude d'une série statistique( p.ex: Etude des revenus d'une population de 1000 ménages), on peut remplacer un diagramme en bâton (données discrètes)par un histogramme (données groupées en classes).

Cette opération- Le groupement des données en classes- nous fait **perdre** l'information et la précision ( faire l'hypothèse d'uniformité à l'intérieur de chaque classe), en contre partie, nous fait **gagner** en lecture, lisibilité, et facilité d'interprétation.

Certains mots: *Information, et signification* doivent être bien précisés, en effet:

L' information doit être utilisée dans le sens de la théorie de l'information ( Voir Information au sens de **Kullback**)

<b>Analyse Factorielle générale(AFG) et analyse en composantes principales(ACP)</b>
---

<b>Notes historiques et bibliographiques</b>
--

L'analyse de données au départ était destinée à répondre aux besoins des Sciences humaines et plus spécialement la psychologie (Spearman-psychologue américain, 1928- s'est penché sur l'étude des profils psychologiques des individus, tel que: Déterminer une variable explicative cachée dans des batteries de tests présentées aux sujets, détecter le caractère dominant d'un sujet,trouver des corrélations entre variables,...,etc.), Burt (connu par ses tableaux, ..)

**Remarquer** la terminologie de l'analyse des données est issue du langage des psychologues: profil, caractère,...,etc.

L'école française durant les années soixante sous légide de J.P. Benzécri a développé cette partie de statistique exploratoire en mettant l'accent sur l'aspect géométrique ( centre de gravité, inertie,...,etc)

Après la diffusion de la micro-informatique et la vulgarisation de l'aspect logiciel, l'analyse des données a pris de l'ampleur et la diversité.

Actuellement elle touche la quasi totalité des domaines: les domaines de la biologie (biométrie, épidémiologie) la météorologie, l'économie, la gestion(marketing: Analyse d'enquêtes d'opinion et segmentation des marchés-objet de la classification hiérarchique-

Succinctement:

Les méthodes d'analyse des données concernent les domaines où une masse importante, voire grandiose de données et de variables sont en jeu.

Le recours à la statistique descriptive uni et bidimensionnelle (représentations graphiques et calcul des caractéristiques ) ne suffisent guère pour répondre aux divers questions, tels que:

Quelle est la variable dominante?.

Quel est le profil d'un individu vérifiant certaine modalité?

et ne peut nullement mettre en évidence l'aspect structurel important tels que les liaisons, les proximités entre les variables, corrélation,...,etc.

La Théorie de L'analyse des données est divisée en deux grandes parties complémentaires, les méthodes d'analyse factorielle et les méthodes de classification.

Si La première est une méthode de **réduction** des données initiales, basée sur le concept de **projection**.

La deuxième est celle de **groupement** par classes homogènes (arboréscence).

Comme la projection sur des espaces (Axe facoriels) nous fait perdre certaine objectivité dans l'interprétation des résultats, on complète généralement notre étude par les méthodes de classification automatique

Commençons par la description de l'analyse générale d'un nuage de points quelconques

### **Remaquer:**

les principales méthodes d'analyse factorielle (ACP, AFC),...,etc) en découlent

L'analyse en composantes principales (ACP): C'est la première méthode d'analyse des données apparue, elle traite les données croisant: individus et variable quantitatives ou qualitatives rendues quantitatives par codage numérique (Cette notion de codage sera expliciter plus tard)

L'analyse factorielle des correspondances (AFC) :C'est la plus importante méthode, privilégiée par l'école française.

- **Remarques**

les différentes méthodes de l'analyse factorielle (leurs applications) dépend de **la nature**

## des données traitées et de leurs représentations sous forme de tableaux

### Analyse générale d'un nuage de points quelconques

On désigne par *l'analyse générale*:

Les résultats classiques qui sont à la base des diverses méthodes d'analyse plus spécialisées. la plupart des méthodes (ACP, AFCorrépondance,...etc) peuvent se ramener au modèle général via une transformation sur les données (matrice de base) tel que (centrer, réduire-normaliser-les observations,...,etc)

Cette méthode générale est basée sur la notion de projection sur des espaces métriques muni de leurs produit scalaire (tel que l'espace euclidien  $\mathbb{R}^n$  muni du produit scalaire usuel) sur les propriétés des matrices symétriques définies positives (type:  ${}^tXX$ ).

L'objectif est d'exhiber (trouver) un sous espace de dimension minimum engendré par une base (formée par des axes qu'on appellera axes factoriels) tel que le nuage ramené à cette base est visible-facile à interpréter tout en minimisant la perte d'information (Principe de toute analyse statistique), cette opération est appelée **ajustement**.

### Ajustement par un sous espace de $\mathbb{R}^p$

Pour simplifier, plaçons nous dans  $\mathbb{R}^p$  muni d'un produit scalaire usuelle-métrique  $M = I_p$ , cherchons à décrire la position par rapport à l'origine  $0_{\mathbb{R}^p}$  et la forme du nuage dans un espace de dimension aussi faible que possible.

Commençons donc par chercher la droite  $F_1$  passant par l'origine de vecteur directeur  $u \in \mathbb{R}^p$  qui ajuste au mieux le nuage.

Ainsi  $F_1 = D_u = \left\{ \alpha u, \text{ avec } \|u\|_{I_p} = 1, \alpha \in \mathbb{R} \right\}$  droite engendrée par le vecteur  $u$  (avec  $u$  unitaire)

La projection du vecteur  $X_i = {}^t(x_{i1}, \dots, x_{ip})$  où ( $i=1, \dots, n$ ) sur la droite  $F_1 = D_u$  s'écrit:  
 $p_{r_u}(X_i) = \alpha.u$  où  $\alpha = \langle u, X_i \rangle_{I_p} = {}^t u X_i$

$${}^t u X_i = (u_1, \dots, u_p) \begin{pmatrix} x_{i1} \\ \vdots \\ x_{ip} \end{pmatrix} = \sum_{j=1}^p u_j x_{ij} \quad i \text{ étant fixé (l'individu } n^{\circ} i)$$

Or le carré de la distance à l'origine d'un point (individu) du nuage qui est une quantité fixe se décompose en carré de la projection sur  $F_1$  soit  $p_r(u)$  et en carré de la distance à  $F_1$  notons la  $d(u)$  (Le fameux Théorème de Pythagore)

Il est clair que:

Rendre minimum  $d(u)$  est équivalent à rendre maximum  $p_u(u)$ .

Maximiser la somme des carrés des projections (pour tous les individus): C'est rendre maximum  $S_1^2$  où

$$S_1^2 = \sum_{i=1}^n \{{}^t u X_i\}^2$$

Or  $\{{}^t u X_i\}^2 = \{{}^t u X_i\}^t \{{}^t u X_i\} = \{{}^t u X_i\} \{{}^t X_i u\} = {}^t u (X_i^t X_i) u$  remarquer le transposé d'un scalaire est lui même.

la quantité à maximiser s'écrit alors

$$S_1^2 = {}^t u \left\{ \sum_{i=1}^n X_i^t X_i \right\} u = {}^t u^t X X u$$

avec la contrainte  ${}^t u u = 1$  (vecteur unitaire)

Remarquer  ${}^t X_i = (x_{i1}, \dots, x_{ip})$  et  ${}^t X = ({}^t X_1, \dots, {}^t X_n)$

$$\text{Résolution du système: } \begin{cases} \max_u S_1^2 = \max_u \{{}^t u \{\sum_{i=1}^n X_i^t X_i\} u\} = \max_u \{{}^t u^t X X u\} \\ \text{Avec la contrainte } {}^t u u = 1 \end{cases}$$

Une des méthodes pour résoudre le système est celle ci:

faisons intervenir l'opérateur de Lagrange  $l$  avec comme multiplicateur  $\lambda$

$$l(u) = {}^t u X X u - \lambda ({}^t u u - 1)$$

Dérivons  $l(u)$  par rapport  $u$  :  $\frac{dl(u)}{du} = 2 {}^t X X u - 2 \lambda u$

On a:

$$\frac{dl(u)}{du} = 0 \Leftrightarrow {}^tXXu = \lambda u$$

Cela signifie que:

**u est un vecteur propre de la matrice  ${}^tXX$  de valeur propre  $\lambda$ .**

**Le maximum cherché est donc une valeur propre de  ${}^tXX$  (la plus grande)**

En effet:

${}^tXXu = \lambda u$  prémultiplions par  ${}^tu$ , cela donne:  ${}^tu{}^tXXu = \lambda{}^tuu$  sous la contrainte  ${}^tuu = 1$   
maximiser  ${}^tu{}^tXXu$  ça revient à exhiber la plus grande  $\lambda$ .

Si nous cherchons l'espace à deux dimensions qui ajuste le mieux le nuage, on doit chercher une deuxième droite de direction  $v$  passant par l'origine et maximise  $v{}^tX{}^tXv$  où  $v$  étant unitaire et orthogonale au vecteur  $u$  ( ${}^tvu=0$ : Le produit scalaire des  $\mathbb{R}^p$  est nul)

On a la relation matricielle  $v{}^tX{}^tXv - 2\lambda v - \beta u = 0$  (généralisation du lagrangien avec deux contraintes)

$\lambda$  et  $\beta$  étant les multiplicateurs de Lagrange, prémultiplions par  ${}^tu$

et que  ${}^tvu=0$  il vient alors  ${}^tXXv = \lambda v$

$v$  sera donc le second vecteur propre, associé à la seconde valeur propre de  ${}^tXX$ , en généralisant,

on obtient le résultat suivant:

Une base orthonormée du sous-espace vectorielle à  $q$  dimensions ajustant au sens des moindres carrées le nuage- maximisations des projections(  $p_r$ )ou minimisation des(d)- est constituée par:

**Les  $q$  vecteurs propres correspondants aux  $q$  plus grandes valeurs propres de la matrice symétrique  ${}^tXX$ .**

### Ajustement par un sous espace de $\mathbb{R}^n$

Dans l'espace  $\mathbb{R}^n$  on dispose de  $p$  (variables)  $X^j$

Soit  $w$  le vecteur des cosinus directeurs d'une droite  $G_1$  passant par l'origine, la condition



que  $G_1$  ajuste au mieux le nuage de  $\mathbb{R}^n$  est: la somme des carrées des projections-  ${}^t w X^j$ —soit maximale, qui se traduit par:

$$\max_v \sum_{j=1}^p \{ {}^t w X^j \}^2$$

or

$$\sum_{j=1}^p \{ {}^t w X^j \}^2 = \sum_{j=1}^p \{ {}^t w X^{jt} X^j w \}$$

Soit

$$S_1^2 = \sum_{j=1}^p \{ {}^t w X^{jt} X^j w \} = {}^t w \sum_{j=1}^p \{ X^{jt} X^j \} {}^t w = {}^t w X^t X w$$

Le vecteur cherché **w à n** composantes rend donc maximum la forme quadratique  $S_1^2 = {}^t w X^t X w$  avec la contrainte  ${}^t v v = 1$

d'après ce qui précède  $w$  est un vecteur propre de la matrice  $(X^t X)_{n \times n}$  (remarquer la dimension:  $(n \times p) \times (p \times n) = n \times n$ ) relatif à la plus grande valeur propre.

## Relation entre l'étude sur $\mathbb{R}^p$ et $\mathbb{R}^n$

Trouvons les relations entre le vecteur  $w_{(n)}$  (matrice  $X^t X$ ) et  $u_{(p)}$  (matrice  ${}^t X X$ )

On a la relation:

$$w_q \text{ est le } q\text{-ième vecteur propre de la matrice } X^t X \text{ de valeur propre } \lambda_q : X^t X w_q = \lambda_q w_q \quad (*)$$

il y'a  $m$  valeurs propres non nulles,  $m$  étant le rang de  $X^t X$ , donc le rang de  $X$  (voir propriété du rang)

E prémultiplions par  ${}^t X$  l'égalité:  $X^t X w_q = \lambda_q w_q$

$${}^t X X ({}^t X w_q) = \lambda_q ({}^t X w_q)$$

or cette égalité signifie: le vecteur  $({}^t X w_q)$  est un vecteur propre de  ${}^t X X$  relative à  $\lambda_q$

A chaque vecteur propre  $w_q$  ( $q \leq m$ ) de  $X^t X$  correspond donc un vecteur propre

$$u_q = {}^t X w_q$$

de  ${}^t X X$  relative à la même valeur propre .

Toute valeur propre non nulle de  $X^t X$  est donc valeur propre de  ${}^t X X$  (et réciproquement) et ainsi les vecteurs propres correspondants sont liés par la relation  $u_q = {}^t X w_q$  avec  ${}^t w w = 1$  (remarquer les dimensions:  $u_q$  est un vecteur de  $p$  lignes)

Les vecteurs propres sont cependant définis à un coefficient près: si  $w_q$  est unitaire alors  $u_q (= {}^t X w_q)$  n'est pas forcément unitaire, en effet:  ${}^t u_q u_q = {}^t ({}^t X w_q) ({}^t X w_q) = {}^t w_q X^t X w_q = \lambda_q$  d'après (\*)

Exprimons  $w_q$  en fonction de  $u_q$  en prémultipliant  $u_q = {}^t X w_q$  par  $X \Rightarrow X u_q = X {}^t X w_q$  d'après  ${}^t u_q u_q = \lambda_q$

donc

$$w_q = \frac{1}{\lambda_q} X u_q$$

Afin de donner à ces relations une forme symétrique

remplaçons l'égalité  $u_q = {}^t X w_q$  par

$$u_q = \frac{1}{\sqrt{\lambda_q}} {}^t X w_q \quad \text{où } q=1, \dots, m$$

$i - e$  : On rend les vecteurs  $u_q$  unitaires

On a dans ce cas

$$w_q = \frac{1}{\sqrt{\lambda_q}} X u_q \quad \text{où } q=1, \dots, m$$

avec  $u_q, w_q$  unitaires

La coordonnée d'un point-individu- de  $\mathbb{R}^p$  sur l'axe factoriel  $F_q$  de cosinus directeur  $u_q$  est  ${}^t u_q X_i$  (ou encore  ${}^t X_i u_q$ )

les coordonnées des  $n$  individus-nuage de l'espace  $\mathbb{R}^p$ -sont donc:

$$({}^t X_1 u_q, {}^t X_2 u_q, \dots, {}^t X_n u_q)$$

Ce sont donc les composantes du vecteur

$$X u_q$$

c'est à dire au coefficient  $\frac{1}{\sqrt{\lambda_q}}$  près, les composantes de  $w_q$  d'après la relation  $w_q = \frac{1}{\sqrt{\lambda_q}} X u_q$ .

De même , les composantes des  $p$  variables de  $\mathbb{R}^n$  sur l'axe factoriel  $G_q$  de cosinus directeur  $w_q$  sont:

$$({}^t X^1 w_q, {}^t X^2 w_q, \dots, {}^t X^p w_q)$$

Ce sont donc les composantes du vecteur

$${}^t X w_q$$

et au même coefficient  $\frac{1}{\sqrt{\lambda_q}}$  près, celles du vecteur  $u_q$  d'après l'égalité:  $u_q = \frac{1}{\sqrt{\lambda_q}} {}^t X w_q$ .

### Résumé

Les cosinus directeurs du  $q$ -ième axe factoriel dans un espace: Sont les coordonnées des points sur le  $q$ -ième axe factoriel de **l'autre** espace, multiplié par

$$\frac{1}{\sqrt{\lambda_q}}.$$

### Reconstitution du tableau X des données initiales

Un repère formé par les  $s$  premiers axes factoriels permet de reconstituer les positions des points avec une précision dépendant du rapport

$$\frac{\lambda_1 + \dots + \lambda_s}{\sum_{i=1}^m \lambda_i} \quad \text{où } s \leq m$$

(on suppose qu'il y'a  $m$  valeurs propres)

Remarque

$$\sum_{i=1}^m \lambda_i = \text{Trace } {}^t X X = \sum_j \sum_i x_{ij}^2$$

Si on veut reconstituer les valeurs numériques initiales

La relation

$$w_q = \frac{1}{\sqrt{\lambda_q}} X u_q \Leftrightarrow X u_q = \sqrt{\lambda_q} w_q \text{ posmultiplions les deux membres par } {}^t u_q : X u_q {}^t u_q = \sqrt{\lambda_q} {}^t u_q w_q$$

faisons la somme sur l'indice q:  $X \sum_{q=1}^p \{u_q^t u_q\} = \sum_{q=1}^p \sqrt{\lambda_q} w_q^t u_q$

La quantité  $\sum_{q=1}^p \{u_q^t u_q\} = I_{p \times p}$  la matrice unite de  $\mathbb{R}^p$  (c'est le produit de la matrice orthogonale par sa transposée, qui est aussi son inverse.d'où:

$$X_{n \times p} = \sum_{q=1}^p \sqrt{\lambda_q} w_q^t u_q$$

Remarquer les dimensions :  $n \times p = (n \times 1)(1 \times p)$

On a donc une reconstitution approchée du tableau X (données initiales)

En se limitant aux **s premiers** axes factoriels, et en négligeant les termes  $\sqrt{\lambda_{s+1}}, \dots, \sqrt{\lambda_p}$

on a:

$$X_{n \times p} = \sum_{q=1}^s \sqrt{\lambda_q} w_q^t u_q$$

On remplace ainsi les  $n \times p$  nombres du tableau X par  $s \times (n + p + 1)$  nombres

Si X est tableau de  $1000 \times 80$  (1000 individus sur lesquels on fait passer 80 variables) et si  $s=10$  (on se contente des dix premiers facteurs), on remplace 80000 données par 10810 nombres(méthode de réduction de données)

En pratique on ne cherche pas à retrouver intégralement le tableau initiale mais un tableau approché.

### L'analyse en composantes principales (ACP)

L'ACP s'applique aux donnée numériques **quantitatives**, présentées sous forme de tableaux:

Individus  $\times$  variables

Elle comporte 3 variantes :ACP générale, ACP centrée réduite et ACP des rangs. on présente par l'ACP centrée réduite dite ACP faite sur la matrice de corrélation, c'est la plus courante.

#### Nature des données:

Les données traitées par l'ACP doivent être quantitatives (modalités mesurables, à titre d'exemples: le nombre d'enfants, l'âge, le poids, la taille,...,etc.)

Ces variables quantitatives sont isolées (**discrètes**) ou groupées en classes (**continues**)

Il convient de noter les données suivant leurs rôle dans l'analyse, ainsi on distingue **les données actives**, celles qui s'impliquent dans toutes les phases de l'analyse: De la détermination de l'espace sur lequel on projette jusqu'à la visualisation et l'interprétation des résultats , les données **passives (illustrative ou encore supplémentaires)** elle interviennent uniquement dans la phase finale de visualisation et permettent de mieux visualiser et de bien caractériser des groupes d'individus et les relations entre les variables initiales.

#### Présentation:

Tableau des données.

Les données soumises à l'ACP sont présentées sous la forme:

Tableau croisant n individus et p variables

X (Matrice des données)					
Individus/Variables	$X^1$	$X^2 \dots$	$X^j \dots$	$X^{p-1}$	$X^p$
1	$x_{11}$	$x_{12} \dots$	$x_{1j} \dots$		$x_{1p}$
.....	.....	.....	.....	.....	.....
i	$x_{i1}$	$x_{i2} \dots$	$x_{ij}$ notée parfois $(x_i^j) ..$		$x_{ip}$
.....	.....	.....	.....	.....	.....
n	$x_{n1}$	$x_{n2} \dots$	$x_{nj} \dots$	.....	$x_{np}$

Ainsi  $x_{ij}$  désigne l'observation de la j-ième variable faite sur le i-ième individu (parfois notée

$x_i^j$ , remarquer le numéro de l'individu  $i$  est en bas)

Les calculs statistiques (évaluation des paramètres) peuvent être effectués pour **chaque colonne**, quant aux lignes si les variables sont **homogènes exprimés dans la même unité**, c'est affirmatif, exemple chiffre d'affaire durant les  $p$  années, revenu durant les  $p$  mensualités, sinon, si les variables sont hétérogènes et hétéroclites, unités différentes, c'est négatif, exemple pour une région déterminée (individu) il est aberrant de parler d'un nombre moyen de ses caractéristiques (variables) tel que la superficie, le nombre d'habitants, infrastructure urbaine (le nombre de logements), nombre de lignes téléphoniques.

Comme la perception visuelle humaine est dans un espace n'excédant pas 3 dimensions ( $\mathbb{R}$  : La droite réelle,  $\mathbb{R}^2$  : le plan, et l'espace :  $\mathbb{R}^3$ ) dès lors que le nombre des données est au delà, des méthodes d'étude s'imposent, l'ACP réduite en est un.

### Présentation de l'ACP réduite

La présentation des données dans l'espace:

Décomposons l'ensemble en deux représentations: Celle des individus et celles des variables

Le tableau des données décrit ci-dessus les  $n$  individus à l'aide des  $p$  variables

**Chaque individu  $i$**  peut donc être représenté par un point de l'espace à  **$p$  dimensions** appelé **espace des variables**, l'ensemble des points représentés dans cet espace est appelé **nuage des individus**

La figure ci-dessous correspond à un nuage des individus

avec espace des variables à 3 dimensions (variable 1, 2, 3)

les coordonnées de  $i$  dans cet espace sont  $(x_{i1}, x_{i2}, x_{i3})$

**De même chaque variable  $X_j$**  peut être représentée par un point dans un espace à  **$n$  dimensions** appelé **espace des individus**, l'ensemble de ces points forme **le nuage des variables**

La figure ci-dessous représente la visualisation de la variable  $X_j$  dans l'espace de 3 individus

Les coordonnées sont(  $x_{1j}, x_{2j}, x_{3j}$ )

L'analyse des individus permet la recherche des ressemblances, d'éventuelle existence homogénéité, de corrélation.

Quant à l'analyse des variables dont le principe consiste outre à déterminer celles qui sont liées et celles qui s'opposent, et d'éventuelles corrélation.

Remarquer la correspondance entre les deux analyses est très forte, en effet: Il s'agit d'une lecture différente du même texte(les données, le corpus statistique ).Chaque individu est décrit par une variable et chaque variable est affecté à un individu. Autrement dit, ils'agit d'une traduction doublée de la même information.

Pour parer à l'hétérogénéité des variables, la solution consiste à les centrer et à les réduire.

### Transformation des données (centrage et normalisation)

Chaque variable  $X^j = {}^t(x_{1j}, \dots, x_{ij}, \dots, x_{nj})$  est remplacée par  $X^{*j} = {}^t(x_{1j}^*, \dots, x_{ij}^*, \dots, x_{nj}^*)$  avec

$$x_{ij}^* = \frac{x_{ij} - \bar{X}^j}{\sigma_j}$$

Les variables se retrouvent ainsi toutes exprimées en nombre d'écart types, ce qui permet de considerer chacune d'elles indépendamment du choix de son unité.

Analyse du nuage des individus					
Individus/Variables	$X^1$	$X^2 \dots$	$X^j \dots$	$X^{p-1}$	$X^p$
1	$x_{11}$	$x_{12} \dots$	$x_{1j} \dots$		$x_{1p}$
.....	.....	.....	.....	.....	.....
i	$x_{i1}$	$x_{i2} \dots$	$x_{ij}$		$x_{ip}$
.....	.....	.....	.....	.....	.....
n	$x_{n1}$	$x_{n2} \dots$	$x_{nj} \dots$	.....	$x_{np}$

Point moyen ou centre de gravité:  $G = {}^t(\bar{X}^1, \dots, \bar{X}^p)$

On a avec  $\mathbf{1}_n = {}^t(1, \dots, 1)$

$$G = {}^t X D_p \mathbf{1}_n$$

Remarquer l'égalité des dimensions:  $p \times 1 = (p \times n) \times (n \times n) \times (n \times 1)$

Soit la matrice (données centrées)

$$Y = X - 1^t G$$

Matrice variance covariance  $V_{p \times p} = [cov(X^j, X^{j'})_{j,j'=1,\dots,p}]$  V s'écrit:

$$V = {}^t Y D_p Y$$

$$\text{où la matrice de poids } D_p = \begin{bmatrix} p_1 & \dots & 0 \\ . & . & \\ 0 & \dots & p_n \end{bmatrix} \quad \text{où } p_i > 0, \text{ et } \sum_{i=1}^n p_i = 1$$

généralement les n individus sont choisis aléatoirement ( $p_i = \frac{1}{n}$  d'où  $D_p = \frac{1}{n} I_n$  avec  $I_n$  matrice unité de  $\mathbb{R}^n$ )

Dans ce cas

$$V = \frac{1}{n} {}^t Y Y \quad \text{avec} \quad Y = X - 1^t G$$

Recherche des axes factoriels  $u_q$  dans  $\mathbb{R}^p$

Comme le vecteur propre de la matrice variance covariance  $V = \frac{1}{n} {}^t Y Y$  correspond à la  $q$ -ième valeur propre  $\lambda_q$

$u_q$  est donc vecteur propre de  ${}^t Y Y$  relative à la valeur propre  $n \times \lambda_q$

Les formules précédentes de l'analyse générale deviennent:

$$w_q = \frac{1}{\sqrt{n\lambda_q}} Y u_q. \text{ La } i\text{-ème composante de } w_q \text{ vaut } w(i)_q = \frac{1}{\sqrt{n\lambda_q}} {}^t Y_i u_q.$$