

# Analyse de données

L. BELLANGER

Master 1 Ingénierie Statistique  
Dpt de Mathématiques - Université de Nantes

## Plan

- 0. Introduction
- I. Outils de représentation d'un échantillon**
- II. Analyse en Composantes Principales (*ACP*)
- III. Analyse Factorielle des Correspondances (*AFC*)
- IV. Classification et Classement
- V. Conclusion

## **I. Outils de représentation d'un échantillon**

### Plan du ch. I

1. *Structure d'un tableau*
2. *Résumés unidimensionnels*
3. *Résumés multidimensionnels*
  - Rappels d'algèbre linéaire
  - Espaces de représentation
4. *Dualité en analyse de données*

# 1. Structure d'un tableau de données

## Questions préalables

Avant toute analyse de données recueillies après *expérience* ou *enquête*, il faut se poser un certain nombre de questions préalables:

- **Que sont les données ?** Sont-elles en nombre important, quelle est leur précision ?
- **D'où proviennent-elles ?** Les valeurs sont-elles raisonnables ?
- **Comment et quand ont-elles été mesurées ?** Des biais sont-ils possibles ? L'observateur sait-il arrondir les nombres ?
- **Comment ont-elles été acquises ?** A-t-on mesuré toutes les unités de la population, ou bien a-t-on fait un échantillonnage ?
- **Existe-t-il une structure logique entre les observations ?** Les observations ont-elles en commun certains facteurs d'environnement ?

Après avoir répondu à ces questions seulement, il est possible de pousser l'analyse plus loin !

5

# 1. Structure d'un tableau de données

## La forme du tableau de données

### • Tableau individus×caractères

- Représentation des données sous la forme d'un **tableau rectangulaire** à :
  - $I$  lignes et  $J$  colonnes ;
  - Formellement, un tableau de données  $X$  est de la forme :

$$X = [x_{ij}] = \begin{bmatrix} x_1^T \\ \vdots \\ x_I^T \end{bmatrix} = [x^1 \quad \dots \quad x^J] \quad i = 1, \dots, I \text{ (ou } n); j = 1, \dots, J \text{ (ou } p)$$

- Les deux indices  $i$  et  $j$  repèrent la **valeur** ligne  $i$ , colonne  $j$ .
- **Mais la représentation et le rôle d'une ligne (individu) ou d'une colonne (variable) sont ≠ :**
  - Le vecteur ligne  $x_i \in \mathbb{R}^J$  **mis en colonne** représente la  $i^{\text{ème}}$  **observation élémentaire** sur laquelle ont été mesurées les  $J$  caractéristiques ou variables.
  - Le vecteur colonne  $x^j \in \mathbb{R}^I$  représente la  $j^{\text{ème}}$  **variable ou indicateur**.

[1] Dans la terminologie du modèle linéaire, on emploie le mot **facteur**.

6

# 1. Structure d'un tableau de données

## La forme du tableau de données

### • Tableau individus×caractères

Exemple (Bouroche & Saporta (1992))

		Caractères					
		Age $x^1$	Revenu imposable $x^2$	...	Salaire brut $x^j$	...	Ancien- neté $x^p$
Individus	1	$x_1^1$	$x_1^2$	...	$x_1^j$	...	$x_1^p$
	2	$x_2^1$	$x_2^2$	...	$x_2^j$	...	$x_2^p$
	...	...	...	...	...	...	...
	i	$x_i^1$	$x_i^2$	...	$x_i^j$	...	$x_i^p$
	...	...	...	...	...	...	...
	n	$x_n^1$	$x_n^2$	...	$x_n^j$	...	$x_n^p$

7

# 1. Structure d'un tableau de données

## La forme du tableau de données

### • Tableau individus×caractères

Exemple suite (Bouroche & Saporta (1992))

Sur les mêmes individus, on aurait aussi pu observer les caractères « **sexe** » et « **situation matrimoniale** ». Pour être traitées numériquement, ces caractères doivent être représentés sous la forme d'un **tableau de variables indicatrices** prenant les valeurs 0 ou 1.

On dit alors que les données sont représentées sous **forme disjonctive complète**

		Caractères					
		Sexe		Situation matrimoniale			
		F	M	Marié	Pacsé	Céliba- taire	Veuf, divorcé
Individus	1	1	0	1	0	0	0
	2	0	1	0	1	0	0
	...	...	...	...	...	...	...
	i	1	0	0	0	1	0
	...	...	...	...	...	...	...
	n	0	1	0	0	0	1

8

# 1. Structure d'un tableau de données

## La forme du tableau de données X

### • Tableau individus×caractères

- X est souvent constitué par **deux sous-tableaux** :  $X = [X^1 \quad X^2]$ 
  - $X^1$  (resp.  $X^2$ ) est un tableau à  $I$  lignes et  $J = p$  (resp.  $q$ ) colonnes.
  - L'individu  $i$  (ligne  $i$  de X), est formée de  $p$  variables et  $q$  indicateurs ; les  $p$  variables définissent le **vecteur observation** lui-même « repéré » par un **vecteur de contrôle**.
  - L'ensemble des méthodes statistiques se rattachent :
    - soit à des **méthodes unidimensionnelles** pour l'analyse séparée de chaque colonne de X.
    - soit à des **méthodes multidimensionnelles** pour l'analyse simultanée de l'ensemble des colonnes de  $X^1$ .
  - L'existence d'un vecteur de contrôle induit une **structure dans le tableau de données** :
    - glt utilisée dans l'analyse des données ou
    - uniq rôle d'aide à l'interprétation.

9

# 1. Structure d'un tableau de données

## La forme du tableau de données X

### • Tableau individus×caractères

- X peut être **complété** par d'autres :
  - **observations** sur les mêmes  $J$  colonnes ; on parlera :
    - **d'observations supplémentaires**, si ces individus ne sont pas pris en compte dans l'analyse initiale
    - par opposition aux premières qui sont des **observations actives**.
  - **colonnes** sur les mêmes individus ; on parlera de :
    - **variables supplémentaires** (les secondes) et
    - de **variables actives** (les premières).

10

# 1. Structure d'un tableau de données

## La forme du tableau de données X

### • Tableau individus×caractères

- Schématiquement le tableau X **peut donc être complété** de la manière suivante :

$$\begin{bmatrix} X & X^{Csup} \\ X_{Lsup} & / \end{bmatrix}$$

le symbole « / » peut être remplacé par **des variables supplémentaires associées à des individus supplémentaires**.

11

# 1. Structure d'un tableau de données

## La forme du tableau de données X

### • Tableau de contingence

- Un **tableau de contingence** est le croisement de 2 variables qualitatives  $A$  et  $B$  :  
le coefficient  $n_{ij}$  du tableau = l'**effectif** (nb d'individus ) présentant à la fois la modalité  $i$  de  $A$  et la modalité  $j$  de  $B$ .
- Dans un tel tableau, les individus ont été regroupés et ne peuvent plus être distingués.
- **Autre représentation** : à chacun des caractères nominaux, on associe un tableau de variables indicatrices (une variable par modalité).

12

# 1. Structure d'un tableau de données

## La forme du tableau de données X

### • Tableau de contingence

- Autre représentation : exemple

$$X_1 = \begin{bmatrix} 100 \\ 010 \\ 010 \\ 100 \\ 001 \\ 001 \\ 010 \\ 001 \end{bmatrix} \quad X_2 = \begin{bmatrix} 10 \\ 10 \\ 01 \\ 10 \\ 01 \\ 01 \\ 10 \\ 10 \end{bmatrix} \quad , \text{ on obtient le tableau de contingence :}$$

$$(X_1)^T X_2 = \begin{bmatrix} 20 \\ 21 \\ 12 \end{bmatrix}$$

13

# 1. Structure d'un tableau de données

## La forme du tableau de données X

### • Tableau de proximité

- On dispose de mesures de ressemblance ou de dissemblance entre tous les objets pris 2 à 2.
  - ex : tableau des distances entre les principales villes de France
- Tableau glt **symétrique** contenant des **nbs**  $\geq 0$  analogues à des distances (ou à des inverses de distances) ; mais n'en possédant pas tjs toutes les propriétés axiomatiques (inégalité triangulaire)
  - Rappel : d est une **distance** si
    - (a)  $d(a, b) = 0 \Leftrightarrow a = b$  ;
    - (b)  $d(a, b) = d(b, a)$  (**symétrie**) ;
    - (c)  $d(a, b) \leq d(a, c) + d(b, c)$  (**inégalité triangulaire**)
- Si (c) n'est pas vérifiée, on parlera plutôt de **dissimilarité**.

14

# 1. Structure d'un tableau de données

## Notion de type de variables

### • Variables quantitatives

- continues** (ou d'échelle) : dont les valeurs forment un sous-ensemble infini de  $\mathbb{R}$   
(ex : poids, taille, revenu)
- discrètes** : dont les valeurs forment un sous-ensemble fini ou infini de  $\mathbb{N}$   
(ex : nombre d'enfants)

### • Variables catégorielles (ou qualitatives)

- dont l'ensemble des valeurs est fini — ces valeurs sont numériques ou alphanumériques, mais quand elles sont numériques, ce ne sont que des codes et non des quantités (ex : n° de département)

15

# 1. Structure d'un tableau de données

## Passage d'un codage simple à un codage disjonctif complet

Qd la valeur d'une variable ne correspond pas à une structure numérique, on peut s'y ramener par un **simple codage**.

**Exemple** : Couleur présentant trois caractéristiques {vert, jaune, marron}

On peut soit :

- établir la correspondance :  
 $\{\text{vert, jaune, marron}\} \Leftrightarrow \{1, 2, 3\}$
- utiliser pour chaque réalisation une variable à deux valeurs slt, 0 ou 1. Chaque modalité de la variable définit une **variable dichotomique** ; la correspondance est alors :

Tableau 1 - Exemple de passage d'un codage simple à un codage disjonctif complet.

	Codage initial	vert	jaune	marron
Vert	1	1	0	0
Jaune	2	0	1	0
marron	3	0	0	1

Pour un traitement statistique, on utilise ce dernier codage appelé **codage disjonctif complet**. Ainsi un objet vert sera analysé par l'ensemble des trois valeurs {1,0,0}, un jaune par {0,1,0} et un marron par {0,0,1}.

16

# 1. Structure d'un tableau de données

## En résumé :

Tableau 2 - Les différents types de variables ou d'indicateurs, nature et exemples.

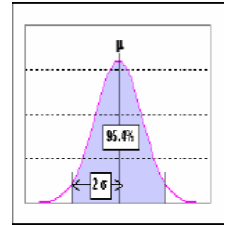
Variable	Type	Nature	Exemples
quantitative	continu		poids, taille, revenu
	entier	discrète	nombre d'enfants
qualitative	polytomie ordonnée	ordinaire	(mauvais, moyen, bon),
	polytomie non ordonnée	nominale	profession, couleur, variété végétale
	dichotomie	binaire	(vrai, faux), (oui, non), (présence, absence)

17

# 2. Résumés unidimensionnels

## Pourquoi ?

- **Explorer la distribution** de chaque variable prise une à une pour :
  - vérifier la fiabilité des variables ;
  - détecter :
    - les valeurs incohérentes ou **manquantes** => **imputation** ou **suppression**
    - les valeurs extrêmes ou aberrantes (**outliers**) à éliminer ?
- **Variables continues**
  - tester la Normalité des variables (surtout si petits effectifs) et les transformer pour augmenter la Normalité.
- **Variables discrètes**
  - regrouper certaines modalités trop nombreuses ou avec des effectifs trop petits (poids trop grand).



18

# 2. Résumés unidimensionnels

## Graphiques

**Histogramme**, **boîte à moustaches** (ou box-plot), **graphique de densité** (après estimation de celle-ci) et **graphique de Normalité** (qq-plot ou pp-plot).

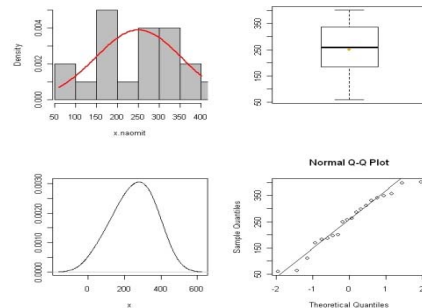


Figure 1. **Eaux1** : Graphiques variable  $HCO_3$

19

# 2. Résumés unidimensionnels

## Graphiques

### Histogramme

- Parlant uniquement si  $n$  est suffisamment élevé ( $>100$ ),
- souvent utile de le remplacer ou ajouter le graphique de densité, version lissée de l'histogramme.

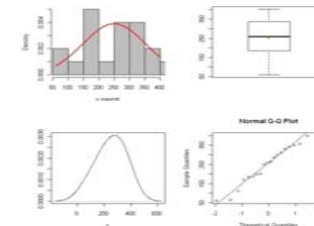


Figure 1. **Eaux1** : Graphiques variable  $HCO_3$

Sous R : `hist()`

20

## 2. Résumés unidimensionnels

### Graphiques

**Boîte à moustaches (ou boxplot)** - Sous R: `boxplot()`

soit  $Q_{0.25}$ ,  $Q_{0.50}$  (la **médiane**) et  $Q_{0.75}$  les valeurs tq 25%, 50% et les 75% des observations leur soient inférieures (i.e. les **quartiles**).

On trace une boîte comme celle de la Fig 2 :

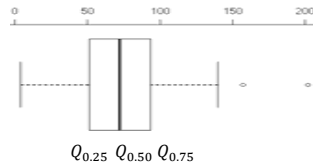


Figure 2. Eaux1 : Boxplot de Ca

- Les extrémités indiquées par « | » sont à :  
 $\max(\min, Q_{0.25} - 1.5(Q_{0.75} - Q_{0.25}))$  et  $\min(\max, Q_{0.75} + 1.5(Q_{0.75} - Q_{0.25}))$
- les valeurs extérieures (« ° ») : **suspectes** ou **extrêmes**
- ↪ Très facile à interpréter ! On voit très rapidement si la distribution est symétrique.

21

## 2. Résumés unidimensionnels

### Paramètres numériques classiques

• **Position** : moyenne  $\bar{x}$ , médiane  $Me$ , mode  $Mo$

- La **moyenne**, notée  $\bar{x}$ , est définie par :

$$\bar{x} = \sum_{i=1}^n x_i / n$$

Si on a  $n_i$  valeurs  $x_i$  identiques, et si le nombre total est  $n = \sum n_i$ , on associe à chaque observation une **pondération**  $p_i = n_i/n$ , et donc  $\sum p_i = 1$  ; alors on a :

$$\bar{x} = \sum_{i=1}^n p_i x_i$$

Sous R: `mean()` et `summary(base)`

On appelle **variable centrée** la variable  $x_c$  de coordonnées :  $x_i - \bar{x}$  ;

↪ la moyenne de cette variable est donc nulle.

22

## 2. Résumés unidimensionnels

### Paramètres numériques classiques

• **Position** : moyenne  $\bar{x}$ , médiane  $Me$ , mode  $Mo$

- La **médiane**, notée  $Me$ , est la valeur qui divise les  $n$  obs. en 2 parties égales (50 % lui sont inférieures et 50 % supérieures)
- paramètre plus stable que  $\bar{x}$  et moins sensible qu'elle aux valeurs suspectes. On dit que c'est un paramètre **robuste**.

Sous R: `median()` et `summary(base)`

23

## 2. Résumés unidimensionnels

### Paramètres numériques classiques

• **Position** : moyenne  $\bar{x}$ , médiane  $Me$ , mode  $Mo$

- Le **mode**, notée  $Mo$ , est la valeur ou la classe de fréquence maximale. Pour une variable quantitative, c'est la valeur la plus probable.
  - contrairement aux autres paramètres de position, le mode peut ne pas être unique : on parlera alors de **distribution multimodale**.
  - Pour une distribution Normale :  $\bar{x} = Me = Mo$ .

24

## 2. Résumés unidimensionnels

### Paramètres numériques classiques

- **Dispersion** : ces paramètres définissent la variabilité de la distribution.

- Le plus utilisé, la **variance**, définie par :

$$s^2 = \text{var}(x) = \sum_{i=1}^n (x_i - \bar{x})^2 / (n - 1)$$

- division par  $(n - 1)$  justifiée pour des raisons statistiques ; quelquefois la division est faite par  $n$  (cas de l'analyse des données ...après) ;
- la variance est alors la moyenne des carrés des écarts à la moyenne.

On note  $s$  l'**écart-type** qui est  $\sqrt{\text{var}(x)}$ . Il s'exprime dans la même unité que la variable étudiée.

On appelle **variable centrée réduite** notée  $x_{CR}$  la variable de coordonnées :  $(x_i - \bar{x})/s$  ; la moyenne de  $x_{CR}$  est donc 0 et son écart-type vaut 1.

Sous R : `var(base)` et `sd(base)`  
avec ....division par  $(n - 1)$

25

## 2. Résumés unidimensionnels

### Paramètres numériques classiques

- **Dispersion** :

- **Étendue** (souvent peu significative à cause des extrêmes) :  
 $\max_i x_i - \min_i x_i$

➢ paramètre très sensible aux **valeurs extrêmes**.

- **Écart interquartile** :  $Q_{0.75} - Q_{0.25}$

➢ plus stable que l'étendue.

- **Coefficient de variation** :  $CV(\%) = 100s/\bar{x}$

➢  $X$  dispersée si  $CV > 25\%$

➢ grandeur sans unité  $\Rightarrow$  utile pour comparer la dispersion de deux échantillons d'une même variable.

26

## 2. Résumés unidimensionnels

### Paramètres numériques classiques

#### • Centrer et réduire une variable

- Pour **centrer** une variable, on ôte à chaque individu la valeur moyenne  $\bar{x}$ .
- Pour **réduire** une variable, on divise chaque individu par l'écart-type  $s_x$ .
- La moyenne devient égale à 0. L'écart-type devient égal à 1.
- On peut ramener ainsi les moyennes et écart-types des variables quantitatives d'un tableau de données à des valeurs identiques.

$$x_{CR} = \frac{x - \bar{x}}{s_x}$$

Sous R : `scale(x, center = TRUE, scale = TRUE)`

27

## 2. Résumés unidimensionnels

### Paramètres numériques classiques

#### • Paramètres de forme

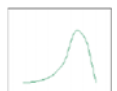
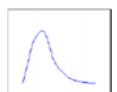
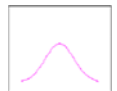
- **Coefficient d'asymétrie** (« **skewness** »)

$$b_1 = \frac{m_3}{m_2^{3/2}} = \frac{1}{n} \left[ \sum_{i=1}^n (x_i - \bar{x})^3 \right] / \left( \frac{1}{n} \left[ \sum_{i=1}^n (x_i - \bar{x})^2 \right] \right)^{3/2}$$

➢ = 0 si la série de données est **symétrique**

➢ > 0 si elle est **allongée vers la droite**  
- fréquent dans les données économiques

➢ < 0 si elle est **allongée vers la gauche**



Sous R : `skewness(e1071 ; agricolae ; moments)`

28



## 2. Résumés unidimensionnels

### Paramètres numériques classiques

#### • Paramètres de forme (suite)

##### ▪ Coefficient d'aplatissement (« kurtosis »)

$$b_2 = \frac{m_4}{m_2^2} = \frac{1}{n} \left[ \sum_{i=1}^n (x_i - \bar{x})^4 \right] / \left( \frac{1}{n} \left[ \sum_{i=1}^n (x_i - \bar{x})^2 \right] \right)^2$$

➤ = 3 si aplatie comme Gauss



➤ > 3 si plus concentrée que Gauss



➤ < 3 si plus aplatie que Gauss



On **normalise** souvent le kurtosis en soustrayant **3** ; SAS et SPSS le font.

Sous R aussi : `kurtosis(e1071 ; agricolae ; moments)` ...

29

## 2. Résumés unidimensionnels

### Les variables qualitatives

#### Pour chaque modalité :

- **Effectif** :  $n_i$  nombre d'individus présentant la modalité  $i$
- **Fréquence** :

$$f_i = \frac{n_i}{\sum_{i=1}^I n_i} = \frac{n_i}{n_+}$$

#### Pour une variable qualitative ordinale

- **Fréquence cumulée** :  $F_i = \sum_{j=1}^i f_j$

Niveau d'études				
	Effectifs	Pourcentage	Pourcentage valide	Pourcentage cumulé
Valide <=Bac	12	12,0	12,0	12,0
IUT	35	35,0	35,0	47,0
Bac+4	26	26,0	26,0	73,0
Ecoles	16	16,0	16,0	89,0
Doctorat	11	11,0	11,0	100,0
Total	100	100,0	100,0	

30

## 2. Résumés unidimensionnels

### Les variables qualitatives

#### Représentations graphiques

Les représentations classiques pour une variable qualitative sont :

- Les **diagrammes en bâtons**
- Les **graphiques circulaires** (camemberts)

Pour chacun de ces graphiques on peut choisir d'afficher les effectifs ou les pourcentages.

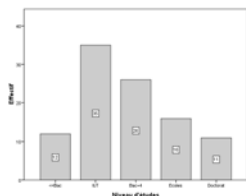
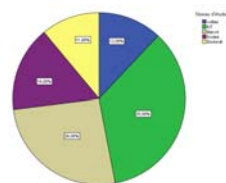


Diagramme en bâtons



Graphique circulaire

Sous R : `barplot()`, `pie()`, `mosaicplot()`

31

## 3. Résumés multidimensionnels

### La forme du tableau de données

un tableau de données  $X$  est une **représentation rectangulaire** à  $n$  **lignes** et  $p$  **colonnes** de la forme :

$$X = \begin{bmatrix} x_1^1 & \dots & x_1^j & \dots & x_1^p \\ \vdots & \ddots & \vdots & \ddots & \vdots \\ x_i^1 & \dots & x_i^j & \dots & x_i^p \\ \vdots & \ddots & \vdots & \ddots & \vdots \\ x_n^1 & \dots & x_n^j & \dots & x_n^p \end{bmatrix}$$

← individu  $i$  : vecteur ligne dans  $\mathbb{R}^p$

↑  
variable  $j$  : vecteur colonne dans  $\mathbb{R}^n$

#### 2 nuages de pts :

- nuage des  $n$  individus dans  $\mathbb{R}^p$
- nuage des  $p$  variables dans  $\mathbb{R}^n$

32



### 3. Résumés multidimensionnels

#### Espaces de représentation

Interprétation **géométrique** des lignes et les colonnes du tableau  $X \in \mathcal{M}_{n \times p}$  par des points dans 2 espaces différents:

l'espace des individus et l'espace des variables :

- **les espaces de représentation :**
  - **celui des  $n$  individus**, de dim  $p$ , noté  $\mathcal{R}^p \subset \mathbb{R}^p$ 
    - Les  $n$  lignes **mis en colonne** sont considérées comme  $n$  pts de l'**espace des individus** à  $p$  dimensions.
    - 2 points sont très proches si les  $p$  coord. de ces 2 pts sont très proches (mêmes valeurs pour les différentes variables).
  - **celui des  $p$  variables**, de dim  $n$ , noté  $\mathcal{R}^n \subset \mathbb{R}^n$ 
    - Les  $p$  colonnes sont considérées comme  $p$  pts de l'**espace des variables** à  $n$  dimensions.
    - 2 variables sont proches si leurs  $n$  coordonnées sont très voisines (i.e. ces variables mesurent la même chose ou sont liées par une relation particulière).

33

### 3. Résumés multidimensionnels

Un jeu de données est constitué par un **triplet**  $(X, Q, D)$  défini par les 3 éléments suivants:

1.  $X = [x_i^j]$  matrice des données brutes  $n$  mesures de  $p$  variables, quantitatives ou non
2.  $Q, p \times p$ , métrique Euclidienne sur l'espace  $\mathbb{R}^p$  des lignes  $x_i$  de  $X$  (**transformées en colonne**)
3.  $D, n \times n$ , métrique Euclidienne sur l'espace  $\mathbb{R}^n$  des colonnes  $x^j$  de  $X$ , **tjs diagonale**.  $D = \text{diag}(p_1, \dots, p_n)$

Les espaces Euclidiens  $(\mathbb{R}^n, D)$  et  $(\mathbb{R}^p, Q)$  sont resp. les **espaces des variables** et **des individus**.

#### Notation

$$r = \text{rg}(X) \leq \min(n, p)$$

34

### 3. Résumés multidimensionnels

#### Rappels d'algèbre linéaire (cf. WikiStat)

1. Notations
2. Matrices
3. Espaces euclidiens
4. Éléments propres
5. Optimisation

35

### 3. Résumés multidimensionnels

#### Rappels d'algèbre linéaire (cf. WikiStat)

- Notations et rappels d'algèbre linéaire de niveau L.
- Introduction des principaux théorèmes d'approximation matricielle par décomposition en valeurs singulières, à la base des méthodes statistique factorielles.

#### 1. Notations

$E$  et  $F$ , 2 espaces vectoriels réels munis respectivement des bases canoniques  $\mathcal{E} = \{e_j; j = 1, \dots, p\}$  et  $\mathcal{F} = \{f_i; i = 1, \dots, n\}$ .

36

### 3. Résumés multidimensionnels

#### Rappels d'algèbre linéaire

##### 2. Matrices

###### 2.1 Notations et définitions

La matrice d'ordre  $(n \times p)$  associée à une application linéaire de  $E$  dans  $F$  est décrite par un tableau :

$$A = \begin{bmatrix} a_1^1 & \dots & a_1^j & \dots & a_1^p \\ \vdots & & & & \\ a_i^1 & \dots & a_i^j & \dots & a_i^p \\ \vdots & & & & \\ a_n^1 & \dots & a_n^j & \dots & a_n^p \end{bmatrix} \in \mathcal{M}_{n \times p}$$

On note par la suite :

- $a_i^j = [A]_i^j$  le terme général de la matrice  $A$  ;
- $a_i = [a_i^1 \dots a_i^p]^T$  un vecteur-ligne *mis en colonne* ;
- $a^j = [a_1^j \dots a_n^j]^T$  un vecteur-colonne.

37

### 3. Résumés multidimensionnels

#### Rappels d'algèbre linéaire

##### 2. Matrices

###### 2.2 Définitions

Une matrice est dite :

- **vecteur-ligne (colonne)** si  $n = 1$  ( $p = 1$ ) ;
- **vecteur-unité d'ordre  $p$**  si elle vaut  $1_p = [1 \dots 1]^T$  ;
- **scalaire** si  $n = 1$  et  $p = 1$  ;
- **carrée** si  $n = p$ .

Une matrice **carrée d'ordre  $p$**  est dite :

- **identité** notée  $I_p$  si  $a_i^j = \delta_i^j = \begin{cases} 0 & \text{si } i \neq j \\ 1 & \text{si } i = j \end{cases}$  ;
- **diagonale** si  $a_i^j = 0$  lorsque  $i \neq j$  ; sous R : **diag()**
- **symétrique** si  $a_i^j = a_j^i \forall (i, j)$  ;
- **triangulaire supérieure** (resp. **inférieure**) si  $a_i^j = 0$  lque  $i > j$  (resp.  $i < j$ ).

38

### 3. Résumés multidimensionnels

#### Rappels d'algèbre linéaire

##### 2. Matrices

###### 2.2 Définitions

Une matrice est dite **partitionnée en blocs** si ses éléments sont eux-mêmes des matrices.

Par exemple :

$$A_{n \times p} = \begin{bmatrix} A_{1(r \times s)}^1 & A_{1(r \times (p-s))}^2 \\ A_{2((n-r) \times s)}^1 & A_{2((n-r) \times (p-s))}^2 \end{bmatrix}$$

Sous R :

Créer une matrice **matrix( )**

39

### 3. Résumés multidimensionnels

#### Rappels d'algèbre linéaire

##### 2. Matrices

###### 2.2 Opérations sur les matrices

- **Somme** :  $[A + B]_i^j = a_i^j + b_i^j$ , pour  $A$  et  $B$  de même ordre  $n \times p$  ;
- **Multiplication par un scalaire** :  $[\alpha A]_i^j = \alpha a_i^j, \forall \alpha \in \mathbb{R}$  ;
- **Transposition** : notation  $A^T$  ou  $A'$ 
  - $[A^T]_i^j = a_j^i, A^T \in \mathcal{M}_{p \times n}$
  - $[A^T]^T = A$  ;
  - $[A + B]^T = A^T + B^T$  ;
  - $(AB)^T = B^T A^T$
- et
  - $\begin{bmatrix} A_{1(r \times s)}^1 & A_{1(r \times (p-s))}^2 \\ A_{2((n-r) \times s)}^1 & A_{2((n-r) \times (p-s))}^2 \end{bmatrix}^T = \begin{bmatrix} [A_{1(r \times s)}^1]^T & [A_{1(r \times (p-s))}^2]^T \\ [A_{2((n-r) \times s)}^1]^T & [A_{2((n-r) \times (p-s))}^2]^T \end{bmatrix}$
  - Sous R : fonction **t( )**

40

### 3. Résumés multidimensionnels

#### Rappels d'algèbre linéaire

#### 2. Matrices

##### 2.2 Opérations sur les matrices

- **Produit scalaire élémentaire** :  $\mathbf{a}^T \mathbf{b} = \sum_{i=1}^n a_i b_i$  où  $\mathbf{a}$  et  $\mathbf{b}$  sont des vecteurs colonnes de  $\mathbb{R}^n$  ;
- **Produit** :  $[\mathbf{AB}]_i^j = \mathbf{a}_i^T \mathbf{b}^j$  avec  $\mathbf{A} \in \mathcal{M}_{n \times p}$ ,  $\mathbf{B} \in \mathcal{M}_{p \times q}$  et  $\mathbf{AB} \in \mathcal{M}_{n \times q}$

Et pour des matrices blocs :

$$\begin{bmatrix} \mathbf{A}_1^1 & \mathbf{A}_1^2 \\ \mathbf{A}_2^1 & \mathbf{A}_2^2 \end{bmatrix} \begin{bmatrix} \mathbf{B}_1^1 & \mathbf{B}_1^2 \\ \mathbf{B}_2^1 & \mathbf{B}_2^2 \end{bmatrix} = \begin{bmatrix} \mathbf{A}_1^1 \mathbf{B}_1^1 + \mathbf{A}_1^2 \mathbf{B}_2^1 & \mathbf{A}_1^1 \mathbf{B}_1^2 + \mathbf{A}_1^2 \mathbf{B}_2^2 \\ \mathbf{A}_2^1 \mathbf{B}_1^1 + \mathbf{A}_2^2 \mathbf{B}_2^1 & \mathbf{A}_2^1 \mathbf{B}_1^2 + \mathbf{A}_2^2 \mathbf{B}_2^2 \end{bmatrix}$$

Sous réserve de compatibilité des dimensions !

Sous R : `A %*% B`

⇒ Faire exercice 3 feuille de Td/TP.

41

### 3. Résumés multidimensionnels

#### Rappels d'algèbre linéaire

#### 2. Matrices

##### 2.3 Propriétés des matrices carrées

###### • Trace

Soit  $\mathbf{A} \in \mathcal{M}_{p \times p}$  on définit la trace de  $\mathbf{A}$  par :  $\text{tr} \mathbf{A} = \sum_{j=1}^p a_j^j$

###### Propriétés de la trace :

- $\text{tr} \alpha = \alpha$  ;
- $\text{tr} \mathbf{A} = \text{tr} \mathbf{A}^T$  ;
- $\text{tr}(\mathbf{A} + \mathbf{B}) = \text{tr} \mathbf{A} + \text{tr} \mathbf{B}$  ;
- $\text{tr} \mathbf{AB} = \text{tr} \mathbf{BA}$ , reste vrai si  $\mathbf{A} \in \mathcal{M}_{n \times p}$  et  $\mathbf{B} \in \mathcal{M}_{p \times n}$  ;
- $\text{tr} \mathbf{CC}^T = \text{tr} \mathbf{C}^T \mathbf{C} = \sum_{i=1}^n \sum_{j=1}^p (c_i^j)^2$  avec  $\mathbf{C} \in \mathcal{M}_{n \times p}$ .
- Sous R : `sum(diag( ))`

42

### 3. Résumés multidimensionnels

#### Rappels d'algèbre linéaire

#### 2. Matrices

##### 2.3 Propriétés des matrices carrées

###### • Déterminant :

« Initialement introduit en [algèbre](#), pour résoudre un [système d'équations linéaires](#) comportant autant d'équations que d'inconnues. Il se révèle un outil très puissant dans de nombreux domaines. Il intervient ainsi dans l'étude des [endomorphismes](#), la recherche de leurs [valeurs propres](#), les propriétés d'[indépendance linéaire](#) de certaines [familles](#) de [vecteurs](#), ... » <http://fr.wikipedia.org/>

On note  $\det(\mathbf{A})$  ou  $|\mathbf{A}|$ , le déterminant de la matrice carrée  $\mathbf{A} \in \mathcal{M}_{p \times p}$

###### Propriétés du déterminant :

- $\det(\mathbf{A}) = \prod_{j=1}^p a_j^j$  si  $\mathbf{A}$  est triangulaire ou diagonale ;
- $\det(\alpha \mathbf{A}) = \alpha^p \det(\mathbf{A})$  ;  $\det(\mathbf{AB}) = \det(\mathbf{A}) \det(\mathbf{B})$  ;
- $\det \begin{bmatrix} \mathbf{A} & \mathbf{B} \\ \mathbf{0} & \mathbf{C} \end{bmatrix} = \det(\mathbf{A}) \det(\mathbf{C})$  ;
- $\det \begin{bmatrix} \mathbf{A}_1^1 & \mathbf{A}_1^2 \\ \mathbf{A}_2^1 & \mathbf{A}_2^2 \end{bmatrix} = \det(\mathbf{A}_1^1) \det(\mathbf{A}_2^2 - \mathbf{A}_2^1 (\mathbf{A}_1^1)^{-1} \mathbf{A}_1^2) = \det(\mathbf{A}_2^2) \det(\mathbf{A}_1^1 - \mathbf{A}_1^2 (\mathbf{A}_2^2)^{-1} \mathbf{A}_2^1)$   
si  $(\mathbf{A}_1^1)^{-1}$  et  $(\mathbf{A}_2^2)^{-1}$  existent
- Sous R : `det( )`

43

### 3. Résumés multidimensionnels

#### Rappels d'algèbre linéaire

#### 2. Matrices

##### 2.3 Propriétés des matrices carrées

###### • Inverse

On note  $\mathbf{A}^{-1}$ , la matrice unique lorsqu'elle existe, inverse de  $\mathbf{A}$ , tq :

$$\mathbf{A}^{-1} \mathbf{A} = \mathbf{A} \mathbf{A}^{-1} = \mathbf{I} ;$$

elle existe si et seulement si  $\det(\mathbf{A}) \neq 0$ .

###### Propriétés de l'inverse :

- $(\mathbf{A}^{-1})^T = (\mathbf{A}^T)^{-1}$  ;
- $(\mathbf{AB})^{-1} = \mathbf{B}^{-1} \mathbf{A}^{-1}$  ;
- $\det(\mathbf{A}^{-1}) = 1/\det(\mathbf{A})$
- Sous R : `solve( )`

44

### 3. Résumés multidimensionnels

#### Rappels d'algèbre linéaire

#### 2. Matrices

##### 2.3 Propriétés des matrices carrées

###### Définitions :

Une matrice carrée est dite:

- **symétrique** si  $A^T = A$  ;
- **singulière** si  $\det(A) = 0$  ;
- **régulière** si  $\det(A) \neq 0$  ;
- **idempotente** si  $AA = A^2 = A$  ;
- **définie-positif** si,  $\forall x \in \mathbb{R}^p, x^T A x \geq 0$  et  $x^T A x = 0 \Rightarrow x = 0$  ;
- **positive**, ou **semi-définie-positif** si,  $\forall x \in \mathbb{R}^p, x^T A x \geq 0$  ;
- **orthogonale** si  $AA^T = A^T A = I$  (ie  $A^T = A^{-1}$ ).

45

### 3. Résumés multidimensionnels

#### Rappels d'algèbre linéaire

#### 3. Espaces euclidiens

Soit  $E$  un espace vectoriel réel de dimension  $p$  isomorphe à  $\mathbb{R}^p$

Rappels succincts :

un **isomorphisme** entre deux **ensembles structurés** est une **application bijective** qui préserve la structure et dont la **réciproque** préserve aussi la structure.

un **espace vectoriel** est un ensemble muni d'une **structure** permettant d'effectuer des **combinaisons linéaires**.

##### 3.1 Sous-espaces (exercice 8 feuille de Td/TP)

- Un sous-ensemble  $E_q$  de  $E$  est un **sous-espace vectoriel** (s.e.v.) de  $E$  s'il est non vide et stable par combinaisons linéaires :  
 $\forall (x, y) \in E_q, \forall \alpha \in \mathbb{R}, \alpha(x + y) \in E_q$ .
- Le  $q$ -uplet  $\{x_1, \dots, x_q\}$  de  $E$  constitue un **système linéairement indépendant** si et slt si :  $\sum_{i=1}^q \alpha_i x_i = 0 \Rightarrow \alpha_1 = \dots = \alpha_q = 0$ .  
Dans le cas contraire ils sont dits **linéairement dépendants**.
- Un système linéairement indépendant  $\varepsilon_q = \{e_1, \dots, e_q\}$  qui engendre le s.e.v.  $E_q = \text{vect}\{e_1, \dots, e_q\}$  en constitue une **base** i.e. tout vecteur de  $E_q$  s'exprime de manière unique comme combinaison linéaire des éléments de  $\varepsilon_q$  ; on a alors :  $\dim(E_q) = \text{card}(\varepsilon_q) = q$ .

46

### 3. Résumés multidimensionnels

#### Rappels d'algèbre linéaire

#### 3. Espaces euclidiens

##### 3.2 Rang d'une matrice $A \in \mathcal{M}_{n \times p}$

Soit  $A$  la matrice d'une application linéaire de  $E = \mathbb{R}^p$  dans  $F = \mathbb{R}^n$ .

Le **rang d'une matrice** : rang de l'**application linéaire qu'elle représente** ou encore rang de la famille de ses vecteurs colonnes (ie la **dimension** du **sous-espace vectoriel** engendré par cette famille).

Le **rang d'une application linéaire**  $f$  de  $E$  dans  $F$  : dimension de son **image**, s.e.v. de  $F$ .

###### Définition

- $\text{Im}(A) = \text{vect}\{a^1, \dots, a^p\}$  est le s.e.v. de  $F$  **image** de  $A$  ;
- $\text{Ker}(A) = \{x \in E ; Ax = 0\}$  est le s.e.v. de  $E$  **noyau** de  $A$  ;

**Théorème du rang** : Il relie  $\dim(E)$  à celle du **noyau** de  $f$  et au rang de  $f$ .

$$\dim(E) = \dim(\text{Im}(A)) + \dim(\text{Ker}(A))$$

et  $\text{rg}(A) = \dim(\text{Im}(A))$  est le rang de  $A$

47

### 3. Résumés multidimensionnels

#### Rappels d'algèbre linéaire

#### 3. Espaces euclidiens

##### 3.2 Rang d'une matrice $A \in \mathcal{M}_{n \times p}$

Le rang d'une matrice  $A$  est :

- le nombre maximal de vecteurs lignes (ou colonnes) linéairement indépendants ;
- la dimension du s.e.v. engendré par les vecteurs lignes (ou colonnes) de  $A$ .

###### Propriétés :

- $\text{rg}(A) = \dim(\text{Im}(A)) ; 0 \leq \text{rg}(A) \leq \min(n, p) ;$
- $\text{rg}(A) = \text{rg}(A^T) ;$
- $\text{rg}(A + B) \leq \text{rg}(A) + \text{rg}(B) ;$
- $\text{rg}(AB) \leq \min(\text{rg}(A), \text{rg}(B)) ;$
- $\text{rg}(A) = \text{rg}(AA^T) = \text{rg}(A^T A) ;$
- Si  $B \in \mathcal{M}_{p \times q}$  de rang  $q < p$  et  $A \in \mathcal{M}_{p \times p}$  de rang  $p$  alors la matrice  $B^T A B$  est de rang  $q$ .

48

### 3. Résumés multidimensionnels

#### Rappels d'algèbre linéaire

#### 3. Espaces euclidiens

##### 3.2 Rang d'une matrice $A \in \mathcal{M}_{n \times p}$

##### Remarques :

- La **Décomposition en Valeurs Singulières** (DVS) de  $A$  est un outil efficace pour déterminer  $\text{rg}(A)$ .
- $A$  est carrée d'ordre  $p$  inversible  $\Leftrightarrow \det(A) \neq 0$   
 $\Leftrightarrow A$  de plein rang (i.e.  $\text{rg}(A) = p$ ).
- Sous R : `qr( )`

49

### 3. Résumés multidimensionnels

#### Rappels d'algèbre linéaire

#### 3. Espaces euclidiens

##### 3.3 Métrique euclidienne (exercice 9 feuille de Td/TP)

Soit  $M \in \mathcal{M}_{p \times p}$ , matrice **carrée**, **symétrique**, **définie-positive** ;  $M$  définit sur l'espace  $E$  :

- un **produit scalaire** :  $\langle x, y \rangle_M = x^T M y$  ;
- une **norme** :  $\|x\|_M = \langle x, x \rangle_M^{1/2}$  ;
- une **distance** :  $d_M(x, y) = \|x - y\|_M$  ;
- des **angles** :  $\cos \theta_M(x, y) = \frac{\langle x, y \rangle_M}{\|x\|_M \|y\|_M}$ .

La matrice  $M$  étant donnée, on dit que :

- une matrice  $A$  est **M-symétrique** si  $(MA)^T = MA$  ;
- deux vecteurs  $x$  et  $y$  sont **M-orthogonaux** si  $\langle x, y \rangle_M = 0$  ;
- un vecteur  $x$  est **M-normé** si  $\|x\|_M = 1$  ;
- une base  $\varepsilon_q = \{e_1, \dots, e_q\}$  est **M-orthonormée** si  $\forall (i, j), \langle e_i, e_j \rangle_M = \delta_i^j$

50

### 3. Résumés multidimensionnels

#### Rappels d'algèbre linéaire

#### 3. Espaces euclidiens

##### 3.3 Métrique euclidienne

**Remarque** : **Obtention** de  $u_1$  **M-normé**

Soient  $u$  vecteur de  $\mathbb{R}^p$  et  $M \in \mathcal{M}_{p \times p}$ , matrice carrée, symétrique, définie-positive :

1. On commence par calculer  $\|u\|_M = \sqrt{u^T M u}$ ,  $M$ -norme de  $u$  ;
2. On définit ensuite :

$$u_1 = \frac{u}{\|u\|_M} \text{ M-normé (ie } \|u_1\|_M = 1)$$

51

### 3. Résumés multidimensionnels

#### Rappels d'algèbre linéaire

#### 3. Espaces euclidiens

##### 3.4 Factorisation de Cholesky

Soit  $M \in \mathcal{M}_{p \times p}$ , matrice **carrée**, **symétrique**, **définie-positive**,

$\exists T$  matrice triangulaire supérieure avec  $t_i^i > 0$  tq

$$M = T^T T ; \text{ cette décomposition est unique !}$$

Sous R : `chol( )`

##### Intérêt :

$M$  carrée, symétrique, définie-positive est souvent utilisée en statistique comme métrique sur  $\mathbb{R}^p$ . On définit ainsi le produit scalaire :  $\langle x, y \rangle_M = x^T M y$  ;  $\forall x, y \in \mathbb{R}^p$

La décomposition de Cholesky de  $M$  donne :  $\langle x, y \rangle_M = \langle Tx, Ty \rangle_{I_p}$

Changer de métrique

$\Leftrightarrow$

effectuer une transformation linéaire des données.

$\Rightarrow$  Voir plus loin en ACP ....

52

### 3. Résumés multidimensionnels

#### Rappels d'algèbre linéaire

#### 3. Espaces euclidiens

##### 3.5 Projection

Soit  $W$  un s.e.v. de  $E$  et  $\mathcal{B} = \{b^1, \dots, b^q\}$  une base de  $W$  ;

$P \in \mathcal{M}_{p \times p}$  est une **matrice de projection M-orthogonale sur  $W$**

$\Leftrightarrow$

$$\forall y \in E, Py \in W \text{ et } \langle Py, y - Py \rangle_M = 0$$

##### Propriété :

Toute matrice  $P$  *idempotente* ( $P^2 = P$ ) et  $M$ -symétrique ( $(MP)^T = MP$ ) est une matrice de projection  $M$ -orthogonale et réciproquement.

53

### 3. Résumés multidimensionnels

#### Rappels d'algèbre linéaire

#### 3. Espaces euclidiens

##### 3.5 Projection

##### Autres propriété :

- Les valeurs propres de  $P$  sont 0 ou 1 :
  - $u \in W$  ;  $Pu = u$  ;  $\lambda = 1$ , de multiplicité  $\dim(W)$  ;
  - $v \perp W$  ; (on note  $v \in W^\perp$ ) ;  $Pv = 0$ ,  $\lambda = 0$ , de multiplicité  $\dim(W^\perp)$ .
- $\text{tr}P = \dim(W)$  ;
- $P = B(B^TMB)^{-1}B^TM$  où  $B = [b^1, \dots, b^q]$
- Dans le cas particulier où
  - les  $b^j$  sont  $M$ -orthonormés alors  $P = BB^TM$
  - $q = 1$  alors  $P = \frac{bb^T}{b^T M b} M$
- Si  $P_1, \dots, P_q$  sont des matrices de projection  $M$ -orthogonales alors la somme  $P_1 + \dots + P_q$  est une matrice de projection  $M$ -orthogonale  $\Leftrightarrow P_k P_j = \delta_j^k P_j$ .
- La matrice  $I_p - P$  est la matrice de projection  $M$ -orthogonale sur  $W^\perp$ .

54

### 3. Résumés multidimensionnels

#### Rappels d'algèbre linéaire

#### 4. Eléments propres : décomposition d'une matrice

Pour §4.1 à §4.3, soit  $A \in \mathcal{M}_{p \times p}$  une matrice **carrée**

##### 4.1 Définitions (exercices 4, 5, 10, 11 feuille de Td/TP)

- Un vecteur  $v \neq 0 \in \mathbb{R}^p$  est appelé **vecteur propre** de  $A$  s'il existe  $\lambda \in \mathbb{C}$  tq :

$$Av = \lambda v$$

le scalaire  $\lambda$  est appelé **valeur propre** de  $A$  associé à  $v$

- Si  $\lambda$  est une valeur propre de  $A$ , le noyau  $\ker(A - \lambda I)$  est un s.e.v. de  $E$ , appelé **sous-espace propre**. Sa dimension est majorée par l'ordre de multiplicité de  $\lambda$ .
- Le **spectre** (ens. des valeurs propres) de  $A$  correspond aux racines, avec leur multiplicité, du **polynôme caractéristique** :

$$P_A(\lambda) = \det(A - \lambda I) = 0$$

ensuite pour trouver  $v$  associé à  $\lambda$  résoudre :  $(A - \lambda I)v = 0$ .

**Remarque :** les vecteurs propres  $v \in \ker(A - \lambda I)$  n'étant définis qu'à une homothétie près, le plus souvent, on choisit de les normer à 1 en utilisant la norme euclidienne.

55

### 3. Résumés multidimensionnels

#### Rappels d'algèbre linéaire

#### 4. Eléments propres

##### 4.2 Propriétés

Soit  $\{\lambda_1, \dots, \lambda_p\}$ , le spectre de  $A$  tq  $\lambda_1 \geq \dots \geq \lambda_p$  alors

1. si  $\lambda_k \neq \lambda_j$ ,  $v^k \perp_M v^j$  ;
2.  $\det(A) = \prod_{i=1}^p \lambda_i$  et  $\text{tr}(A) = \sum_{i=1}^p \lambda_i$

##### Par conséquent :

- $A$  est régulière  $\Leftrightarrow \forall k, \lambda_k \neq 0$  (aucune v.p. nulle);
- $A$  est positive (resp. définie-positive)  $\Leftrightarrow \lambda_p \geq (\text{resp. } >) 0$  ;
- si  $(\lambda, v)$  est un élément propre de  $A$  alors
 
$$A^k v = \lambda^k v, \forall k \in \mathbb{N}^* \text{ et si } A \text{ est régulière } A^{-1}v = \lambda^{-1}v$$

56

### 3. Résumés multidimensionnels

#### Rappels d'algèbre linéaire

#### 4. Éléments propres

##### 4.2 Propriétés

###### Décomposition spectrale :

Soit  $\{\lambda_1, \dots, \lambda_p\}$ , le spectre de  $A \in \mathcal{M}_{p \times p}$  tq  $\lambda_1 \geq \dots \geq \lambda_p$  alors

$$\Lambda = V^{-1}AV \Leftrightarrow A = V\Lambda V^{-1}$$

où

$\Lambda = \text{diag}(\lambda_1, \dots, \lambda_p)$  et

$V = [v^1 \ \dots \ v^p]$  (matrice de passage de l'ancienne base à la base diagonale) matrice

contenant les vecteurs propres de  $A$  associés aux valeurs propres rangées par ordre décroissant dans  $\Lambda$ .

*Quelques matrices diagonalisables particulières :*

symétriques, M-symétriques, définies-positives, de projection M-orthog.

57

### 3. Résumés multidimensionnels

#### Rappels d'algèbre linéaire

#### 4. Éléments propres

##### 4.3 Théorèmes

###### Théorème de Cayley-Hamilton (pour ex. 11) :

Si le polynôme caractéristique  $P_A(\lambda) = \lambda^p + P_{p-1}\lambda^{p-1} + \dots + P_1\lambda + P_0$  alors  $P_A(A) = A^p + P_{p-1}A^{p-1} + \dots + P_1A + P_0I = 0$

###### Théorème 1. — Soit deux matrices $A \in \mathcal{M}_{n \times p}$ et $B \in \mathcal{M}_{p \times n}$

Les valeurs propres non nulles de  $AB$  et  $BA$  sont identiques avec le même degré de multiplicité.

Si  $u$  est vecteur propre de  $BA$  associé à la valeur propre différente de zéro, alors  $v = Au$  est vecteur propre de la matrice  $AB$  associé à la même valeur propre.

*Décomposition spectrale* pour des types particuliers de matrices.

58

### 3. Résumés multidimensionnels

#### Rappels d'algèbre linéaire

#### 4. Éléments propres

##### 4.3 Théorèmes

**Théorème 2.** — Une matrice  $A$  *symétrique réelle* admet  $p$  valeurs propres *réelles*. Ses vecteurs propres peuvent être choisis pour constituer une **base orthonormée** de  $E$  ;  $A$  se décompose en :

$$A = V\Lambda V^T = \sum_{k=1}^p \lambda_k v^k (v^k)^T$$

où

$\Lambda = \text{diag}(\lambda_1, \dots, \lambda_p)$  et

$V = [v^1 \ \dots \ v^p]$  matrice orthogonale contenant les vecteurs propres de  $A$  associés aux valeurs propres réelles rangées par ordre décroissant dans  $\Lambda$ .

(pour ex. 10)

59

### 3. Résumés multidimensionnels

#### Rappels d'algèbre linéaire

#### 4. Éléments propres

##### 4.3 Théorèmes

**Théorème 3.** — Une matrice  $A$  *M-symétrique réelle* admet  $p$  valeurs propres réelles. Ses vecteurs propres peuvent être choisis pour constituer une **base M-orthonormée** de  $E$  ;  $A$  se décompose en :

$$A = V\Lambda V^T M = \sum_{k=1}^p \lambda_k v^k (v^k)^T M$$

où

$\Lambda = \text{diag}(\lambda_1, \dots, \lambda_p)$  et

$V = [v^1 \ \dots \ v^p]$  matrice  $M$ -orthogonale ( $V^T M V = I_p$  et  $V V^T = M^{-1}$ ) contenant les vecteurs propres de  $A$  associés aux valeurs propres réelles rangées par ordre décroissant dans  $\Lambda$ .

Par définition, si  $A$  est aussi *def. positive*, on note la *racine carrée* de  $A$  :

$$A^{1/2} = V\Lambda^{1/2}V^T M = \sum_{k=1}^p \sqrt{\lambda_k} v^k (v^k)^T M$$

60



### 3. Résumés multidimensionnels

#### Rappels d'algèbre linéaire

##### 4. Éléments propres

###### Remarques :

- Les décompositions ne sont pas uniques :
  - pour une valeur propre simple (de multiplicité 1) le vecteur propre normé est défini à un signe près, tandis que
  - pour une valeur propre multiple, une infinité de bases M-orthonormées peuvent être extraites du sous-espace propre unique associé.
- $rg(A) = rg(\Lambda)$  où  $rg(\Lambda)$  n'est autre que le nombre (répétées avec leurs multiplicités) de valeurs propres non nulles.
- Sous R : `eigen()`

61

### 3. Résumés multidimensionnels

#### Rappels d'algèbre linéaire

##### 4. Éléments propres

###### 4.4 Décomposition en valeurs singulières (DVS)

↪ Construction de la décomposition d'une matrice **rectangulaire**  $X \in \mathcal{M}_{n \times p}$

Pour une matrice rectangulaire, la notion de valeur propre n'a pas de sens ; néanmoins :

- les matrices carrées  $X^T X$  et  $XX^T$  sont sym., semi-définie positives ;
- $rg(X) = rg(X^T X) = rg(XX^T) = r$ , les valeurs propres  $\neq 0$  (positives) de  $XX^T$  et  $X^T X$ .

**Définition :** On appelle **valeurs singulières** de  $X$ , les racines carrées des valeurs propres  $\neq 0$  de  $XX^T$  et  $X^T X$ .

62

### 3. Résumés multidimensionnels

#### Rappels d'algèbre linéaire

##### 4.4 Décomposition en valeurs singulières (DVS)

**Version « maigre » :** Toute matrice  $X_{n \times p}$  peut être décomposée de la façon suivante :

$$X_{n \times p} = U_{n \times r} \Theta_{r \times r} (V_{p \times r})^T$$

Où

- $\Theta_{r \times r} = \Lambda_r^{1/2}$  matrice **diagonale** contenant les  $r$  **valeurs singulières** de la matrice  $X$  ; i.e. les racines carrées des v.p. **positives et non nulles**  $\lambda_i$  de  $X^T X$  ou de  $XX^T$  ;  $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_r > 0$
- $r = rg(X) \leq \min(n, p)$ . Si  $r = p$ , la matrice  $X$  est dite de **plein rang**.
- $U_{n \times r} = [u^1 \dots u^r]$  unitaire  $n \times r$  et tq  $u^i$  est vecteur propre de  $XX^T$  associé à la valeur propre non nulle  $\lambda_i$   
⇒ Vecteurs principaux de l'espace des colonnes
- $V_{p \times r} = [v^1 \dots v^r]$  unitaire  $p \times r$  et tq  $v^i$  est vecteur propre de  $X^T X$  associé à la valeur propre non nulle  $\lambda_i$   
⇒ Vecteurs principaux de l'espace des lignes

**U et V sont des matrices colonnes-orthonormales** (i.e. matrices contenant des vecteurs colonnes unitaires et orthogonaux 2 à 2).

63

### 3. Résumés multidimensionnels

#### Rappels d'algèbre linéaire

##### 4.4 Décomposition en valeurs singulières (DVS)

###### Remarques DVS :

- $[u^1 \dots u^r]$  vecteurs principaux de l'espace des colonnes ;  $\|u^i\| = 1$  et si  $i \neq j$  alors  $\langle u^i, u^j \rangle = 0$ .
- $[v^1 \dots v^r]$  vecteurs principaux de l'espace des lignes ;  $\|v^i\| = 1$  et si  $i \neq j$  alors  $\langle v^i, v^j \rangle = 0$ .
- $\sqrt{\lambda_1} \geq \sqrt{\lambda_2} \geq \dots \geq \sqrt{\lambda_r} > 0$  valeurs singulières.
- Cette décomposition **n'est pas unique**.
- Sous R : `svd()`

64

### 3. Résumés multidimensionnels

#### Rappels d'algèbre linéaire

##### 4.4 Décomposition en valeurs singulières (DVS)

###### En pratique :

- Le nb de valeurs singulières fournit le  $rg$  de la matrice  $X$ . Ces valeurs singulières sont ordonnées, ce qui induit un ordre sur les colonnes de  $U$  et  $V$ .
- Calcul de  $U$  et  $V$  : Calcul des vecteurs propres de  $X^T X$  ou de  $XX^T$  (prendre la matrice de + petite dim). Les vecteurs propres de l'autre se déduisent par les *formules de transition* :

$$\begin{cases} U = X V \Lambda_r^{-1/2} \\ V = X^T U \Lambda_r^{-1/2} \end{cases}$$

Ex 12 feuille de Td/TP

65

### 3. Résumés multidimensionnels

#### Rappels d'algèbre linéaire

##### 4.4 Décomposition en valeurs singulières (DVS)

Sous R : `svd()`

###### Applications : Rang d'une matrice

$$\begin{bmatrix} 2 & 1 & 3 & 4 \\ -1 & 6 & -3 & 0 \\ 1 & 20 & -3 & 8 \end{bmatrix} = \begin{bmatrix} -0.08682 & 0.84068 \\ -0.26247 & -0.53689 \\ -0.96103 & 0.07069 \end{bmatrix} \begin{bmatrix} 22.650 & 0 \\ 0 & 6.081 \end{bmatrix} \begin{bmatrix} -0.03851 & -0.92195 & 0.15055 & -0.35477 \\ 0.37642 & -0.15902 & 0.64476 & 0.64600 \end{bmatrix}$$

$$r = rg(X) = 2 \Rightarrow X \text{ n'est pas de plein rang}$$

66

### 3. Résumés multidimensionnels

#### Rappels d'algèbre linéaire

##### 4.4 Décomposition en valeurs singulières (DVS)

###### Applications : Décomposition d'une matrice carrée

Considérons une matrice de variance-covariance  $S = X^T D X$

$$S_{2 \times 2} = U_{2 \times 2} \Theta_{2 \times 2} V_{2 \times 2}^T$$

$$\begin{bmatrix} 9.2 & 1.6 \\ 1.6 & 5.8 \end{bmatrix} = \begin{bmatrix} -0.8944 & -0.4472 \\ -0.4472 & 0.8944 \end{bmatrix} \begin{bmatrix} 9 & 0 \\ 0 & 5 \end{bmatrix} \begin{bmatrix} -0.8944 & -0.4472 \\ -0.4472 & 0.8944 \end{bmatrix}$$

- Les valeurs singulières de  $S$  obtenues sur la diagonale de la matrice  $\Theta$  sont égales aux *valeurs propres* de  $S$  qui sont réelles : ceci est vrai pour toute matrice carrée symétrique.
- La matrice  $U$  (resp.  $V$ ) contient des vecteurs colonnes identiques aux vecteurs propres obtenus par diagonalisation de  $S$ , leurs signes peuvent varier en fonction du programme utilisé.

67

### 3. Résumés multidimensionnels

#### Rappels d'algèbre linéaire

##### 4.4 Décomposition en valeurs singulières (DVS)

###### Applications : Inverse généralisée d'une matrice

Inversion de  $X^T X$  quand  $X$  n'est pas de plein  $rg$  (pb en régression quand pb de colinéarité).

On peut montrer que :

$$(X^T X)^- = V \Theta^- U^T$$

où

$\Theta^-$  est l'inverse généralisée de  $\Theta$  obtenue en prenant dans cette dernière matrice l'inverse des éléments diagonaux non nuls.

68

### 3. Résumés multidimensionnels

#### Rappels d'algèbre linéaire

#### 4.4 Décomposition en valeurs singulières (DVS)

##### Généralisation :

En Analyse des Données (AD) : définition de métriques sur **chacun** des espaces associés aux lignes et aux colonnes d'une matrice réelle.

On définit le triplet  $(\mathcal{M}_{n \times p}, \mathbf{Q}, \mathbf{D})$  par la donnée de :

- $\mathcal{M}_{n \times p}$  espace vectoriel des matrices réelles  $n \times p$ , de dim  $np$
- $\mathbf{Q} \in \mathcal{M}_{p \times p}$  métrique sur l'espace des lignes  $\mathcal{R}^p \subset \mathbb{R}^p$  ;
- $\mathbf{D} \in \mathcal{M}_{n \times n}$  métrique sur l'espace des colonnes  $\mathcal{R}^n \subset \mathbb{R}^n$ .

DVS de  $\mathbf{X} \in \mathcal{M}_{n \times p}$  relativement à 2 matrices symétriques et positives  $\mathbf{Q}$  et  $\mathbf{D}$

69

### 3. Résumés multidimensionnels

#### Rappels d'algèbre linéaire

#### 4.4 Décomposition en valeurs singulières (DVS)

##### Généralisation : Version « maigre »

**Théorème 4.** — Une matrice  $\mathbf{X} \in \mathcal{M}_{n \times p}$  de rang  $r$  peut s'écrire :

$$\mathbf{X}_{n \times p} = \mathbf{U}_{n \times r} \mathbf{\Theta}_{r \times r} (\mathbf{V}_{p \times r})^T = \sum_{k=1}^r \sqrt{\lambda_k} \mathbf{u}^k (\mathbf{v}^k)^T$$

- $\mathbf{U}_{n \times r}$  contient les vecteurs propres  $\mathbf{D}$ -orthonormés ( $\mathbf{U}^T \mathbf{D} \mathbf{U} = \mathbf{I}_r$ ) de la matrice  $\mathbf{D}$ -symétrique positive  $\mathbf{X} \mathbf{Q} \mathbf{X}^T \mathbf{D}$  associés aux  $r$  v.p.  $\lambda_k \neq 0$  rangées par ordre décroissant dans la matrice diagonale  $\mathbf{\Theta}_{r \times r}^2 = \mathbf{\Lambda}_r$  ;
- $\mathbf{V}_{p \times r}$  contient les vecteurs propres  $\mathbf{Q}$ -orthonormés ( $\mathbf{V}^T \mathbf{Q} \mathbf{V} = \mathbf{I}_r$ ) de la matrice  $\mathbf{Q}$ -symétrique positive  $\mathbf{X}^T \mathbf{D} \mathbf{X} \mathbf{Q}$  associés aux mêmes valeurs propres.

De plus,

$$\mathbf{U} = \mathbf{X} \mathbf{Q} \mathbf{V} \mathbf{\Lambda}_r^{-1/2} \text{ et } \mathbf{V} = \mathbf{X}^T \mathbf{D} \mathbf{U} \mathbf{\Lambda}_r^{-1/2}$$

70

### 3. Résumés multidimensionnels

#### Rappels d'algèbre linéaire

#### 5. Optimisation

##### 5.1 Norme d'une matrice

Soit le triplet  $(\mathcal{M}_{n \times p}, \mathbf{Q}, \mathbf{D})$  par la donnée de :

- $\mathcal{M}_{n \times p}$  espace vectoriel des matrices réelles  $n \times p$ , de dim  $np$
- $\mathbf{Q} \in \mathcal{M}_{p \times p}$  métrique sur l'espace des lignes  $\mathcal{R}^p \subset \mathbb{R}^p$  ;
- $\mathbf{D} \in \mathcal{M}_{n \times n}$  métrique sur l'espace des colonnes  $\mathcal{R}^n \subset \mathbb{R}^n$ .

**Définition :** On munit  $\mathcal{M}_{n \times p}$  du **produit scalaire de Hilbert-Schmidt** :

$$\langle \mathbf{X}, \mathbf{Y} \rangle_{\mathbf{Q}, \mathbf{D}} = \text{tr}(\mathbf{X} \mathbf{Q} \mathbf{Y}^T \mathbf{D}) ; \forall \mathbf{X}, \mathbf{Y} \in \mathcal{M}_{n \times p}$$

Cas particulier :  $\mathbf{Q} = \mathbf{I}_p$  et  $\mathbf{D} = \mathbf{I}_n$ , on a :

$$\langle \mathbf{X}, \mathbf{Y} \rangle_{\mathbf{I}_p, \mathbf{I}_n} = \text{tr}(\mathbf{X} \mathbf{Y}^T) = \sum_{i=1}^n \sum_{j=1}^p x_i^j y_i^j$$

71

### 3. Résumés multidimensionnels

#### Rappels d'algèbre linéaire

#### 5. Optimisation

##### 5.1 Norme d'une matrice

- La norme associée au *produit scalaire de Hilbert-Schmidt* est encore appelée **norme trace** ou **norme de Frobenius** :

$$\|\mathbf{X}\|_{\mathbf{Q}, \mathbf{D}}^2 = \text{tr}(\mathbf{X} \mathbf{Q} \mathbf{X}^T \mathbf{D}) ; \forall \mathbf{X} \in \mathcal{M}_{n \times p}$$

Cas particulier :  $\mathbf{Q} = \mathbf{I}_p$  et  $\mathbf{D} = \mathbf{I}_n$ , on a :

$$\|\mathbf{X}\|_{\mathbf{I}_p, \mathbf{I}_n}^2 = \text{tr}(\mathbf{X} \mathbf{X}^T) = \sum_{i=1}^n \sum_{j=1}^p (x_i^j)^2 = \text{SSQ}(\mathbf{X})$$

SSQ : « sum of squares ».

- Si  $\mathbf{D} = \text{diag}(p_1, \dots, p_n)$ , la distance associée à cette norme s'écrit :

$$d^2(\mathbf{X}, \mathbf{Y}) = \|\mathbf{X} - \mathbf{Y}\|_{\mathbf{Q}, \mathbf{D}}^2 = \sum_{i=1}^n p_i \|\mathbf{x}_i - \mathbf{y}_i\|_{\mathbf{Q}}^2$$

appelé **critère des moindres carrés ordinaires**.

72

### 3. Résumés multidimensionnels

#### Rappels d'algèbre linéaire

#### 5. Optimisation

##### 5.2 Approximation d'une matrice

Soient  $X \in \mathcal{M}_{n \times p}$  de rang  $r$ ,  $Q \in \mathcal{M}_{p \times p}$  et  $D \in \mathcal{M}_{n \times n}$

**Objectif:** Trouver la matrice  $Z_q$  de rang  $q < r$ , qui soit la + proche possible de  $X$ .

**Théorème 5.** La solution du problème :

$$\min_Z \{ \|X - Z\|_{Q,D}^2 ; Z \in \mathcal{M}_{n \times p}, rg(Z) = q < r \}$$

est donnée par la somme des  $q$  1<sup>ers</sup> termes de la DVS (Théo 4) de  $X$  :

$$Z_q = \sum_{k=1}^q \sqrt{\lambda_k} u^k (v^k)^T = U_{n \times q} \Theta_{q \times q} (V_{p \times q})^T$$

- Le minimum atteint est :  $\|X - Z_q\|_{Q,D}^2 = \sum_{k=q+1}^r \lambda_k$
- Les matrices  $U_{n \times q}$  et  $V_{p \times q}$  contiennent les  $q$  1<sup>ers</sup> vecteurs et valeurs propres donnés par la DVS de  $X$ .

$Z_q$  est appelée **approximation de rang  $q$**  de  $X$ .

73

### 3. Résumés multidimensionnels

#### Espaces de représentation

Interprétation **géométrique** des lignes et les colonnes du tableau  $X \in \mathcal{M}_{n \times p}$  par des points dans 2 espaces différents:

l'espace des individus et l'espace des variables :

- **les espaces de représentation :**

- celui des  $n$  individus**, de dim  $p$ , noté  $\mathcal{R}^p \subset \mathbb{R}^p$ 
  - Les  $n$  lignes **mis en colonne** sont considérées comme  $n$  pts de l'espace des individus à  $p$  dimensions.
  - 2 points sont très proches si les  $p$  coord. de ces 2 pts sont très proches (mêmes valeurs pour les différentes variables).
- celui des  $p$  variables**, de dim  $n$ , noté  $\mathcal{R}^n \subset \mathbb{R}^n$ 
  - Les  $p$  colonnes sont considérées comme  $p$  pts de l'espace des variables à  $n$  dimensions.
  - 2 variables sont proches si leurs  $n$  coordonnées sont très voisines (i.e. ces variables mesurent la même chose ou sont liées par une relation particulière).

74

### 3. Résumés multidimensionnels

Un jeu de données est constitué par un **triplet**  $(X, Q, D)$  défini par les 3 éléments suivants:

- $X = [x_i^j]$  matrice des données brutes  $n$  mesures de  $p$  variables, quantitatives ou non
- $Q$ ,  $p \times p$ , métrique Euclidienne sur l'espace  $\mathbb{R}^p$  des  $n$  lignes  $x_i$  de  $X$  (**transformées en colonne**)
- $D$ ,  $n \times n$ , métrique Euclidienne sur l'espace  $\mathbb{R}^n$  des  $p$  colonnes  $x^j$  de  $X$ , **tjs diagonale**.  $D = \text{diag}(p_1, \dots, p_n)$

Les espaces Euclidiens  $(\mathbb{R}^n, D)$  et  $(\mathbb{R}^p, Q)$  sont resp. les **espaces des variables** et **des individus**.

#### Notation

$$r = rg(X) \leq \min(n, p)$$

75

### 3. Résumés multidimensionnels

#### Espaces de représentation

- Le **point moyen** du nuage des  $n$  points de  $\mathbb{R}^p$ , aussi appelé **centre de gravité ou centroïde**, a comme coordonnées celles du **vecteur moyenne**, calculées par :

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i = \frac{1}{n} X^T \mathbf{1}_{n \times 1} \in \mathbb{R}^p$$

où  $\mathbf{1}_{n \times 1}$  vecteur à  $n$  lignes formées de 1

i.e.  $\bar{x} = (\bar{x}^1 \quad \dots \quad \bar{x}^p)$  où  $\bar{x}^j$  est la **moyenne empirique** de la variable  $j$  (de même  $s^j$  sera son **écart-type**)

76

### 3. Résumés multidimensionnels

#### Espaces de représentation

En retranchant à chaque colonne sa moyenne, on obtient la **matrice des données centrées** :

$$\begin{aligned} \mathbf{X}_C^T &= [x_1 - \bar{x} \quad \dots \quad x_n - \bar{x}] = \mathbf{X}^T - \bar{x} \mathbf{1}_{1 \times n} \\ &= \mathbf{X}^T \left[ \mathbf{I}_n - \frac{1}{n} \mathbf{1}_{n \times 1} (\mathbf{1}_{n \times 1})^T \right] = \mathbf{X}^T \left[ \mathbf{I}_n - \frac{1}{n} \mathbf{1}_{n \times n} \right] \end{aligned}$$

- où  $\mathbf{1}_{n \times n}$  matrice formée de 1 et  $\mathbf{I}_n - \frac{1}{n} \mathbf{1}_{n \times n}$  sont des matrices carrées de dim  $n$  ; la 1<sup>ère</sup> est de rang 1 et la 2<sup>nde</sup> de rang  $n - 1$ .
- Ces matrices carrées sont **idempotentes** ( $\mathbf{A} \times \mathbf{A} = \mathbf{A}$ ) et symétriques ( $\mathbf{A}^T = \mathbf{A}$ ) : ce sont des **matrices de projection**.
- En transposant la dernière équation, on a :

$$\mathbf{X}_C = \left[ \mathbf{I}_n - \frac{1}{n} \mathbf{1}_{n \times n} \right] \mathbf{X}$$

77

### 3. Résumés multidimensionnels

#### Première transformation : matrice de dispersion

- **Matrice de variances-covariances**

$$\mathbf{S}_{p \times p} = \mathbf{X}_C^T \mathbf{D} \mathbf{X}_C = \begin{bmatrix} s_1^2 & s_{12} & \dots \\ & s_2^2 & \dots \\ \vdots & & \ddots & \dots \\ & & & s_p^2 \end{bmatrix}$$

où  $\mathbf{D}$  est la **matrice carrée** ( $n \times n$ ) **diagonale des poids**.

- Si toutes les obs ont la même importance, elles ont même poids, i.e.  $1/n$ . On a bien sûr :  $\mathbf{D} = \mathbf{I}_n/n$
- Quelquefois, il peut être utile de leur donner des poids  $\neq$ ,  $p_i$  pour l'obs  $i$ . Ces poids sont des nbs  $> 0$  de somme à 1 :

$$\mathbf{D} = \begin{bmatrix} p_1 & 0 & \dots & 0 & 0 \\ 0 & p_2 & 0 & \dots & 0 \\ \vdots & & \ddots & & \vdots \\ 0 & 0 & \dots & 0 & p_n \end{bmatrix} \text{ avec } \sum_{i=1}^n p_i = 1$$

- Avec pondération, le **vecteur moyenne** est alors calculé par :

$$\bar{x} = \mathbf{X}^T \mathbf{D} \mathbf{1}_{n \times 1} \in \mathbb{R}^p \quad \text{et} \quad \mathbf{X}_C = [\mathbf{I}_n - \mathbf{1}_{n \times n} \mathbf{D}] \mathbf{X}$$

78

### 3. Résumés multidimensionnels

#### Première transformation : matrice de dispersion

- **Matrice de variances-covariances**

- On a alors la **matrice centrée réduite** par :

$$\mathbf{X}_{CR} = \mathbf{X}_C (\text{diag}(\mathbf{S}))^{-1/2}$$

- **Matrice des coeff. de corrélation linéaire**  $\mathbf{R}$ , carrée de dim  $p$

$$\mathbf{R} = \mathbf{D}_{1/s} \mathbf{S} \mathbf{D}_{1/s} = \mathbf{X}_{CR}^T \mathbf{D} \mathbf{X}_{CR}$$

Où

$\mathbf{D}_{1/s}$  est la matrice diagonale définie par  $\mathbf{D}_{1/s} = (\text{diag}(\mathbf{S}))^{-1/2}$ .

- Sous R : avec  $1/(n-1)$  pour le calcul et non  $1/n$  comme en AD
  - `var( ) ; cor( )`
  - `scale(x, center = TRUE, scale = TRUE)`

79

### 3. Résumés multidimensionnels

- **Matrice des coeff. de corrélation linéaire**  $\mathbf{R}$ , de dim  $p$

$$\mathbf{R} = \mathbf{D}_{1/s} \mathbf{S} \mathbf{D}_{1/s} = \mathbf{X}_{CR}^T \mathbf{D} \mathbf{X}_{CR}$$

- Le **degré de liaison** entre 2 variables quantitatives se mesure à l'aide du **coefficient de corrélation**. D'autant plus voisin de +1 ou -1 que la liaison est étroite.
- **Attention ! Une corrélation n'est pas forcément une relation de cause à effet.**
- **Attention ! L'absence de corrélation ne signifie pas que les variables ne sont pas liées.** Il peut exister des relations non linéaires.
- **Remarque sur**  $\mathbf{S}$  et  $\mathbf{R}$  : matrices carrées d'ordre  $p$ , symétriques et def. positives  $\Rightarrow$  **valeurs propres  $\geq 0$** .

80

### 3. Résumés multidimensionnels

#### • Liaison entre 2 variables qualitatives : exemple tiré de Analyses multidimensionnelles, INRA formation permanente, Janvier 2006, André Bouchier

- Nous avons artificiellement découpé la population de l'INRA en 4 types professionnels :

A pour administratifs  
I pour ingénieurs  
S pour scientifiques  
T pour techniciens

- Nous avons croisé ce type professionnel avec le sexe des agents INRA (ramené à 50/50 pour faciliter la compréhension). Nous avons obtenu le **tableau de contingence** suivant :

	A	I	S	T	Total
F	15	9	2	24	50
H	0	5	15	30	50
Total	15	14	17	54	100

81

### 3. Résumés multidimensionnels

#### • Liaison entre 2 variables qualitatives : exemple

**Exemple** : homogénéité de 2 échantillons de la population de l'INRA

Tableau des effectifs

	A	I	S	T	Total
F	15	9	2	24	50
H	0	5	15	30	50
Total	15	14	17	54	100

Tableau des effectifs ajustés (fréquences théoriques)

	A	I	S	T	Total
F	7.5	7	8.5	27	50
H	7.5	7	8.5	27	50
Total	15	14	17	54	100

82

### 3. Résumés multidimensionnels

#### • Liaison entre 2 variables qualitatives : exemple

- Le tableau des écarts à l'indépendance, (fréquences observées - fréquences théoriques) contient l'information « **intéressante** »

Tableau des écarts à l'indépendance  
(fréq obs - fréq théorique)

	A	I	S	T	Total
F	7.5	+2	-6.5	-3	0
H	-7.5	-2	+6.5	+3	0
Total	0	0	0	0	0

83

### 3. Résumés multidimensionnels

#### • Liaison entre 2 variables qualitatives : exemple

- On calcule la contribution du khi2 pour chaque cellule du tableau.

$$K_{hi2} = \frac{(FréqObservée - FréqThéorique)^2}{FréqThéorique}$$

Tableau des contributions au Khi2

	A	I	S	T
F	7.5	0.57	4.97	0.33
H	7.5	0.57	4.97	0.33

- Le Khi2 est ici égal à 26.74. Mais ce chiffre n'a pas de signification en lui même. En effet, plus le tableau sera grand, plus le Khi2 sera élevé.

84



### 3. Résumés multidimensionnels

#### • Liaison entre 2 variables qualitatives : exemple

- On calcule les contributions en pourcentage de chaque cellule au  $\chi^2$

	A	I	S	T	Total
F	28,05	2,13	18,59	1,23	50%
H	28,05	2,13	18,59	1,23	50%
Total	56,10	4,26	37,17	2,47	100%

- On remarque par exemple que la contribution maximale est apportée par la catégorie A.
- Plus de 93% de l'inertie du tableau est due aux catégories A et S. Concrètement, ça veut dire quoi ?

85

### 3. Résumés multidimensionnels

#### Espaces de représentation

Interprétation **géométrique** des lignes et des colonnes du tableau  $X \in \mathcal{M}_{n \times p}$  par des points dans 2 espaces différents:

l'espace des individus et l'espace des variables :

#### • les espaces de représentation :

- celui des  $p$  variables, de dim  $n$ , noté  $\mathcal{R}^n \subset \mathbb{R}^n$ 
  - Les  $p$  colonnes sont considérées comme  $p$  pts de l'espace des variables à  $n$  dimensions.
  - 2 variables sont proches si leurs  $n$  coordonnées sont très voisines (i.e. ces variables mesurent la même chose ou sont liées par une relation particulière).
- celui des  $n$  individus, de dim  $p$ , noté  $\mathcal{R}^p \subset \mathbb{R}^p$ 
  - Les  $n$  lignes **mis en colonne** sont considérées comme  $n$  pts de l'espace des individus à  $p$  dimensions.
  - 2 points sont très proches si les  $p$  coord. de ces 2 pts sont très proches (mêmes valeurs pour les différentes variables).

86

### 3. Résumés multidimensionnels

#### L'espace des individus : regard sur les lignes, inertie

L'analyse du nuage de points utilise la notion fondamentale de distance.

- En physique, on munit l'espace des individus de la distance euclidienne classique:

$$d^2(x_1, x_2) = (x_1^1 - x_2^1)^2 + \dots + (x_1^p - x_2^p)^2 = (x_1 - x_2)^T (x_1 - x_2)$$

- En statistique, pb souvent un peu plus compliqué car les unités de mesure sont  $\neq$  ; on peut avoir un mélange de mensurations, de poids et d'âge... Il faut donc **pondérer**, d'une certaine manière, les différents carrés de l'expression précédente.

On choisit  $\mathbb{R}^p$  muni de la **métrique**  $Q \in \mathcal{M}_{p \times p}$  (matrice carrée, sym. def. pos.), lui conférant une structure d'espace euclidien :

**Produit scalaire** :  $\langle x_1, x_2 \rangle_Q = x_1^T Q x_2$

**Carré de la distance** :  $d(x_1, x_2) = (x_1 - x_2)^T Q (x_1 - x_2)$

**Q-norme de  $x_1$**  :  $\|x_1\|_Q^2 = x_1^T Q x_1$

87

### 3. Résumés multidimensionnels

#### L'espace des individus : regard sur les lignes, inertie

#### • Matrice produit scalaire entre individus : $W = X Q X^T$

- $Q = I_p$  : **produit scalaire usuel**. Les unités initiales de chaque variable sont conservées
  - revient à rendre dominantes les variables de plus grande variabilité.
- $Q = (\text{diag}(S))^{-1}$  : **métrique inverse des variances**.
  - Chaque var est divisée par son écart-type ; son importance devient indep. de sa dispersion.
  - distances calculées non plus sur  $X$ , mais sur la matrice centrée réduite  $X_{CR}$ .
- Utiliser une métrique diagonale quelconque  $Q = T^T T$  car sym. pos....

88



### 3. Résumés multidimensionnels

*L'espace des individus : regard sur les lignes, inertie*

*Matrice produit scalaire entre individus :  $W = XQX^T$*

**Propriété (Factorisation de Cholesky) :**  $Q$  étant symétrique, def. positive, on peut écrire  $Q = T^T T$  avec  $T$  matrice triangulaire sup., d'où

$$\langle x_i; x_{i'} \rangle_Q = (Tx_i)^T Tx_{i'}$$

Tout se passe comme si on utilisait la métrique  $Q = I_p$  sur le tableau  $XT^T$ .

Utiliser la métrique  $Q = (\text{diag}(S))^{-1}$ , sur le tableau centré  $X_C$

$\Leftrightarrow$

utiliser la métrique  $Q = I_p$  sur le tableau centré réduit  $X_{CR}$ .

89

### 3. Résumés multidimensionnels

*L'espace des individus : regard sur les lignes, inertie*

• *L'inertie totale du nuage de points :  $I_g$*

$\Rightarrow$  L'inertie  $I_g$  permet de quantifier la **variabilité** contenue dans un tableau de données.

- On appelle **inertie** la quantité d'information contenue dans un tableau de données.
- Une inertie nulle signifie que tous les individus sont presque identiques.
- L'inertie du nuage sera égale à la somme des variances des  $p$  variables.
- Si les  $p$  variables sont centrées-réduites, l'inertie sera égale à  $p$ .

*Soit sous forme mathématique !*

90

### 3. Résumés multidimensionnels

*L'espace des individus : regard sur les lignes, inertie*

• *L'inertie totale du nuage de points :  $I_g$*

**Définition :** L'**inertie totale**, notée  $I_g$ , est la *moyenne des carrés des distances pondérées* des  $n$  pts au centre de gravité  $\bar{x}$ . Soit :

$$I_g = \sum_{i=1}^n p_i \|x_i - \bar{x}\|_Q^2 = \text{tr} \left( \sum_{i=1}^n p_i (x_i - \bar{x})^T Q (x_i - \bar{x}) \right) = \text{tr}(QS)$$

➤ Cette quantité caractéristique du nuage mesure **l'éloignement des points par rapport à leur centre de gravité** :

Ex :  $I_g$  proche de 0  $\Rightarrow$  les individus sont identiques ou presque et sont confondus avec leur centre de gravité.

91

### 3. Résumés multidimensionnels

*L'espace des individus : regard sur les lignes, inertie*

• *L'inertie totale du nuage de points :  $I_g$*

**Propriétés :**

$$I_g = \text{tr}(QS) = \text{tr}(SQ) = \text{tr}(WD) = \text{tr}(DW)$$

où

$W = XQX^T$ , matrice produit scalaire entre individus et  $S = X_C^T D X_C$ , matrice des variances-covariances.

D'où si :

- $Q = I_p$ , l'inertie est égale à la **somme des  $p$  variances**,
- $Q = (\text{diag}(S))^{-1}$ , l'inertie est égale à  **$p$** , le nombre de variables.

92

### 3. Résumés multidimensionnels

#### L'espace des individus : regard sur les lignes, inertie

##### Matrice d'inertie

##### Définition :

Pour une métrique qcque  $Q \in \mathcal{M}_{p \times p}$ , on appelle **matrice d'inertie** de  $X$ , la matrice

$$S_I = X^T D X Q$$

- Pour  $Q = I_p$ , on a :  $S_I = S = X^T D X$ , matrice de variances-covariance, si  $X = X_C$  est centrée.

##### Interprétation de $S_I$ :

$S_I$  est symétrique et en gal def. positive.

Elle définit une forme quadratique et pour  $u$  donné, tq  $u^T u = 1$ , on a  $u^T X^T D X u$  qui vaut **l'inertie du nuage projeté sur  $\mathbb{R}u$**  (cf cours ACP,  $X = X_C$  et  $Q = I_p$ ).

93

### 3. Résumés multidimensionnels

#### L'espace des individus : regard sur les lignes, inertie

##### L'inertie totale du nuage de points : Propriétés

**Propriété :** L'inertie totale  $I_g$  est la moitié de la moyenne pondérée des carrés des distances entre les individus:

$$2I_g = \sum_{i=1}^n \sum_{i' \neq i} p_i p_{i'} \|x_i - x_{i'}\|_Q^2$$

**Dem :** Soit  $N$  un nuage de points  $\{(x_i, p_i)_{i=1, \dots, n}\}$  de c.d.g.  $\bar{x}$

$$\|x_i - x_{i'}\|_Q^2 = \|x_i - \bar{x} + \bar{x} - x_{i'}\|_Q^2$$

$$= \|x_i - \bar{x}\|_Q^2 + \|\bar{x} - x_{i'}\|_Q^2 + 2 \langle x_i - \bar{x}, \bar{x} - x_{i'} \rangle_Q$$

$$\text{D'où comme } \sum_i p_i (x_i - \bar{x}) = 0$$

$$\sum_{i=1}^n \sum_{i' \neq i} p_i p_{i'} \|x_i - x_{i'}\|_Q^2 = 2 \sum_{i=1}^n p_i \|x_i - \bar{x}\|_Q^2 = 2I_g.$$

94

### 3. Résumés multidimensionnels

#### L'espace des individus : regard sur les lignes, inertie

##### L'inertie totale du nuage de points : Propriétés

##### Propriété:

$$I_g = tr(X^T D X Q) = tr(S_I) = \|X\|_{Q \times D} \text{ norme de Hilbert-Schmidt de } X$$

$$= tr(QS) = tr(SQ)$$

$$= tr(WD) = tr(DW) \text{ si } X \text{ est centrée}$$

où

$W = X Q X^T$ , matrice produit scalaire entre individus et

$S = X^T D X$ , matrice des var.- cov.

D'où si :

- $Q = I_p$ , l'inertie est égale à la **somme des  $p$  variances**,
- $Q = \text{diag}\{S\}^{-1}$ , l'inertie est égale à  **$p$** , le nbre de variables.

**dem :** si  $X$  est centrée

95

### 3. Résumés multidimensionnels

#### L'espace des individus : regard sur les lignes, inertie

##### L'inertie totale du nuage de points : Propriétés

**Dem :** si  $X$  est centrée i.e.  $X = X_C$

$$I_g = \sum_{i=1}^n p_i \|x_i\|_Q^2 = \sum_{i=1}^n p_i (x_i)^T Q (x_i) : \text{si } \lambda \text{ nb réel positif } tr(\lambda) = \lambda \text{ d'où}$$

$$= \sum_{i=1}^n tr(p_i (x_i)^T Q (x_i)) = \sum_{i=1}^n tr(p_i Q (x_i) (x_i)^T) : \text{or } tr(AB) = tr(BA) \text{ d'où ;}$$

$$= tr\left(\sum_{i=1}^n p_i (x_i)^T Q (x_i)\right) = tr(Q X^T D X) = tr(X^T D X Q).$$

96

### 3. Résumés multidimensionnels

#### L'espace des individus : regard sur les lignes, inertie

- L'inertie totale du nuage de points : Propriétés

##### Théorème de Huyghens :

Si on note:

$$I_a = \sum_{i=1}^n p_i (\mathbf{x}_i - \mathbf{a})^T \mathbf{Q} (\mathbf{x}_i - \mathbf{a}) = \sum_{i=1}^n p_i \|\mathbf{x}_i - \mathbf{a}\|_Q^2$$

l'inertie du nuage de points par rapport à un pt  $\mathbf{a}$  quelconque

$\Rightarrow$  l'inertie totale du nuage de pts  $I_g$  vérifie :

$$I_a = I_g + \|\bar{\mathbf{x}} - \mathbf{a}\|_Q^2$$

97

### 3. Résumés multidimensionnels

#### L'espace des individus : regard sur les lignes, inertie

- L'inertie totale du nuage de points

$\Rightarrow$  projection des individus de  $\mathbb{R}^p$  dans un ss-espace de  $\dim < p$

**Principe ACP** (cf ch. 2) = obtenir une *représentation approchée* du nuage des  $n$  individus dans un ss-espace de faible  $\dim k < p$ .

$\Rightarrow$  Pour y parvenir, **projection** des individus sur un ss-espace de dimension plus faible ( $= k < p$ ), ss-espace de projection choisi selon le critère suivant:

Les distances en projection devront être le moins déformées possible.

98

### 3. Résumés multidimensionnels

#### L'espace des variables : regard sur les colonnes

- Pour étudier la proximité des variables, il faut munir cet espace d'une **métrique** ie trouver une matrice carrée d'ordre  $n$  sym. def. pos.
- On choisit la **matrice diagonale des poids D** (matrice  $n \times n$ , sym. def. pos.), lui conférant une structure d'espace euclidien :  $\mathbf{x}^j \in \mathbb{R}^n$
- Nous supposons les  $\mathbf{x}^j$  centrées :

**Produit scalaire entre 2 variables** :  $\langle \mathbf{x}^j; \mathbf{x}^k \rangle_D = \mathbf{x}^{jT} \mathbf{D} \mathbf{x}^k = \sum_{i=1}^n p_i x_i^j x_i^k$

**Norme d'une variable** :  $\|\mathbf{x}^j\|_D^2 = s_j^2$

ie la « longueur » d'une variable est égale à son écart-type.

Dans un espace euclidien, on définit l'angle  $\theta$  entre 2 vecteurs par son cosinus de la façon suivante:

$$\cos(\theta_{jk}) = \frac{\langle \mathbf{x}^j; \mathbf{x}^k \rangle_D}{\|\mathbf{x}^j\|_D \|\mathbf{x}^k\|_D} = \frac{s_{jk}}{s_j s_k} = r_{jk}$$

Le cosinus de l'angle entre 2 variables centrées est égal à leur **coefficient de corrélation linéaire**.

99

### 3. Résumés multidimensionnels

#### Définition d'une étude

On appelle **étude** un **triplet**  $(\mathbf{X}, \mathbf{Q}, \mathbf{D})$  où

- $\mathbf{X}$  est **tableau de données** croisant individus et variables ;
- $\mathbf{Q}$  est la métrique permettant le calcul des distances entre individus ;
- $\mathbf{D}$  est la métrique des poids, permettant le calcul des distances entre variables.

$\hookrightarrow$  Une étude peut être caractérisée ou représentée par deux « **objets** » différents :

$$\mathbf{W} = \mathbf{X} \mathbf{Q} \mathbf{X}^T \text{ ou } \mathbf{S} = \mathbf{X}^T \mathbf{D} \mathbf{X}$$

100

### 3. Résumés multidimensionnels

#### Transformation linéaire d'un tableau de données

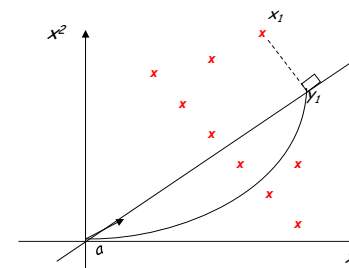
- Les analyses que nous allons étudier ultérieurement consistent, dans la majorité des cas, à **trouver des combinaisons linéaires des  $p$  variables** définissant les colonnes de  $X$ .
- Nous chercherons des **combinaisons linéaires** qui ont des **propriétés particulières**.
  - Ces combinaisons linéaires seront des transformations des vecteurs  $x^i (i = 1, \dots, p)$  qui vont fournir de nouvelles variables.

101

### 3. Résumés multidimensionnels

#### Transformation linéaire d'un tableau de données

- A chaque vecteur d'observation  $x \in \mathbb{R}^p$ , on peut associer un **axe** de l'espace des observations  $\mathbb{R}^p$  ; une **transformation linéaire** consiste à **projeter les observations sur un nouvel axe de vecteur unitaire  $a \in \mathbb{R}^p$** .



102

### 3. Résumés multidimensionnels

#### Transformation linéaire d'un tableau de données

- Si on a  $i$  ( $i = 1, \dots, n$ ) observations, la **coordonnée de l'obs  $i$  sur cet axe** est définie par :

$$y_i = a^T Q(x_i - \bar{x}) = \langle a; x_i - \bar{x} \rangle_Q \in \mathbb{R}$$

- L'ensemble des  $n$  coordonnées par :

$$y = XQa = Xu = \sum_{j=1}^p x^j u_j \in \mathbb{R}^n \text{ avec } u = Qa$$

- Nous avons donc créé une **nouvelle variable  $y$**  avec :
  - un **axe** de vecteur unitaire  $a$ ,
  - un vecteur  $y$  de l'espace des variables, **composante principale**
  - une forme linéaire  $u$  qui est appelée **facteur**.

103

### 3. Résumés multidimensionnels

#### Transformation linéaire d'un tableau de données

- Si nous faisons  $q$  **transformations de l'ensemble du tableau** définies par la matrice  $A = [a^1 \mid \dots \mid a^q]$  tq :

$$Y_{n \times q} = X_{n \times p} Q_{p \times p} A_{p \times q}$$

- Supposons ici que  $Q = I_p$ 
  - On peut montrer que la **matrice de dispersion  $S_Y$**  de  $Y$  se déduit de celle de  $X$  par :

$$S_{Y, q \times q} = A_{q \times p}^T S_{X, p \times p} A_{p \times q}$$

104

### 3. Résumés multidimensionnels

#### Transformation linéaire d'un tableau de données

- Il est souvent utile d'imposer des conditions à la matrice de transformation  $A$  :
  - les  $q$  vecteurs sont **unitaires** :  $\|a^k\|^2 = 1; k = 1, \dots, q$
  - les  $q$  vecteurs sont **orthogonaux 2 à 2** :  $(a^k)^T a^l = 0; k \neq l = 1, \dots, q$
- La matrice  $A$  est une **matrice orthogonale** si elle est carrée (c'est-à-dire ici si  $q = p$ ) et  $AA^T = A^T A = I_q$ .
- quand  $q = p$ , la transformation est une simple **rotation** du repère de départ (avec  $Q = I_p$ ) :
  - $x_i \in \mathbb{R}^p$  transformée en  $y_i = A^T x_i; \forall i = 1, \dots, n$
  - Chaque composante du nouveau vecteur est la **projection** de  $x_i$  sur les  $p$  vecteurs  $a^j; j = 1, \dots, p$

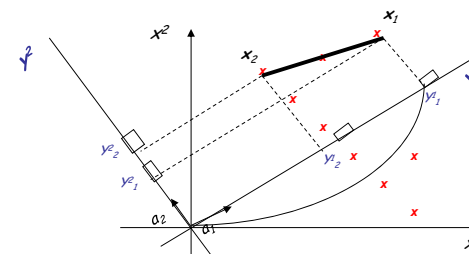
105

### 3. Résumés multidimensionnels

#### Transformation linéaire d'un tableau de données

- **Propriété intéressante et importante** :  
La **distance** entre 2 obs est **invariante** ds une **transformation linéaire** :

$$(y_r - y_s)^T (y_r - y_s) = (x_r - x_s)^T A^T A (x_r - x_s) = (x_r - x_s)^T (x_r - x_s)$$



106

### 3. Résumés multidimensionnels

#### Décomposition d'une matrice de données X

La plupart des méthodes d'Analyse Factorielle des Données peuvent être présentées dans un cadre commun :

celui de l'extension du théorème de la DVS au cadre d'espaces Euclidiens généraux.

Elle correspond à la DVS du triplet  $(X, Q, D)$

107

### 3. Résumés multidimensionnels

#### Décomposition d'une matrice de données X

- **Rappel : Décomposition en valeurs singulières (DVS)**

**Version « maigre »** : Toute matrice  $X_{n \times p}$  peut être décomposée de la façon suivante :

$$X_{n \times p} = U_{n \times r} \Theta_{r \times r} (V_{p \times r})^T$$

Où

- $\Theta_{r \times r} = \Lambda^{1/2}$  matrice **diagonale** contenant les  $r$  **valeurs singulières** de la matrice  $X$  ; i.e. les racines carrées des v.p. **positives et non nulles**  $\lambda_i$  de  $X^T X$  ou de  $XX^T$  ;  $\lambda_1 \geq \lambda_2 \geq \dots \lambda_r > 0$
- $r = \text{rg}(X) \leq \min(n, p)$ . Si  $r = p$ , la matrice  $X$  est dite de **plein rang**.
- $U_{n \times r} = [u^1 \dots u^r]$  unitaire  $n \times r$  et tq  $u^i$  est vecteur propre de  $XX^T$  associé à la valeur propre non nulle  $\lambda_i$   
 $\Rightarrow$  Vecteurs principaux de l'espace des colonnes
- $V_{p \times r} = [v^1 \dots v^r]$  unitaire  $p \times r$  et tq  $v^i$  est vecteur propre de  $X^T X$  associé à la valeur propre non nulle  $\lambda_i$   
 $\Rightarrow$  Vecteurs principaux de l'espace des lignes
- **U et V sont des matrices colonnes-orthonormales** (i.e. matrices contenant des vecteurs colonnes unitaires et orthogonaux 2 à 2).

108

### 3. Résumés multidimensionnels

#### Décomposition d'une matrice données X

##### Schéma de dualité

- Q est la matrice d'un produit scalaire de  $\mathbb{R}^p$  (*L'espace des individus*), soit une matrice carrée symétrique def. pos. qui définit la fonction :

$$(x, y) = \left( \begin{bmatrix} x_1 \\ \vdots \\ x_p \end{bmatrix}, \begin{bmatrix} y_1 \\ \vdots \\ y_p \end{bmatrix} \right) \in \mathbb{R}^p \times \mathbb{R}^p \mapsto x^T Q y = \langle x, y \rangle_Q \in \mathbb{R}$$

- D est la matrice d'un produit scalaire de  $\mathbb{R}^n$  (*L'espace des variables*), soit une matrice carrée symétrique def. pos. qui définit la fonction :

$$(x, y) = \left( \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix}, \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix} \right) \in \mathbb{R}^n \times \mathbb{R}^n \mapsto x^T D y = \langle x, y \rangle_D \in \mathbb{R}$$

- L'essentiel est que ces 4 matrices X,  $X^T$ , Q et D s'assemblent car les produits  $XQ$ ,  $QX^T$ ,  $X^T D$  et  $DX$  ont un sens. On les dispose dans un schéma dit *schéma de dualité* :

109

### 4. Dualité en analyse de données

#### Décomposition d'une matrice données

##### Schéma de dualité

Cf pour plus de détails <http://pbil.univ-lyon1.fr/R/cours/stage3.pdf>

- Idée d'origine de P. Cazes, il est appelé analyse générale par Greenacre (1984, Annexe A2 p. 346) ; introduit en écologie par Escoufier (1987).
- Un schéma est *constitué de trois éléments* donnant un triplet (X, Q, D).
- X est une *matrice de données* en gl issue d'une matrice de données brutes  $\tilde{X}$  à l'aide d'une transformation préalable. X a n lignes et p colonnes.
  - Les n lignes de X sont des vecteurs de  $\mathbb{R}^p$  tandis que
  - les p colonnes de X sont des vecteurs de  $\mathbb{R}^n$ .

110

### 4. Dualité en analyse de données

#### Décomposition d'une matrice données

##### Schéma de dualité

- On les dispose dans un schéma dit *schéma de dualité* :

$$\begin{array}{ccc} \mathbb{R}^p & \xrightarrow{Q} & \mathbb{R}^{p*} \\ X^T \uparrow & & \downarrow X \\ \mathbb{R}^{n*} & \xleftarrow{D} & \mathbb{R}^n \end{array}$$

- Rigoureusement,  $\mathbb{R}^{p*}$  est le dual de  $\mathbb{R}^p$  (ensemble des applications linéaires de  $\mathbb{R}^p$  dans  $\mathbb{R}$ ),  $\mathbb{R}^{n*}$  est le dual de  $\mathbb{R}^n$  (ensemble des applications linéaires de  $\mathbb{R}^n$  dans  $\mathbb{R}$ ),
- Q est vue comme la matrice d'une application linéaire définie par  $(Q(x))(y) = \langle x, y \rangle_Q$ ,
- D est vue comme la matrice d'une application linéaire définie par  $(D(x))(y) = \langle x, y \rangle_D$ .

111

### 4. Dualité en analyse de données

#### Décomposition d'une matrice données

##### Schéma de dualité

- Pour l'utilisateur, la simplification :

$$(X, Q, D) \Leftrightarrow \begin{array}{ccc} [p] & \xrightarrow{Q} & [p] \\ (X, Q, D) \Leftrightarrow X^T \uparrow & & \downarrow X \\ [n] & \xleftarrow{D} & [n] \end{array}$$

suffit pour se souvenir que les produits de matrices :  
 $X^T D X Q$ ,  $D X Q X^T$ ,  $X Q X^T D$  et  $Q X^T D X$   
 ont un sens.

112

## 4. Dualité en analyse de données

### Décomposition d'une matrice données

#### Schéma de dualité - DVS du triplet (X, Q, D)

- On définit deux matrices  $S = X^T D X$  et  $W = X Q X^T$  qui s'insèrent par :

$SQ = S_I$  et  $WD$  sont les opérateurs d'Escoufier.

- Le schéma a des **propriétés théoriques** très gales qui prennent des significations propres lorsqu'on utilise un ensemble de paramètres originaux.
- En particulier les **ACP classiques**, l'analyse des correspondances (**AFC**), les analyses de correspondances non symétriques (**ANSC**), les analyses discriminantes (**AFD**), l'analyse canonique (**AC**), les analyses de co-inertie (**ACO**), les analyses sur variables instrumentales (**ACPVI**), l'analyse des correspondances multiples (**ACM**), l'analyse factorielle multiple (**AFM**) et diverses extensions sont des **conséquences directes des propriétés générales**.

113

## 4. Dualité en analyse de données

### Décomposition d'une matrice de données X

#### DVS de (X, Q, D)

Soient les matrices réelles  $X_{n \times p}$  de rang  $r$  et les métriques  $Q_{p \times p}$  et  $D_{n \times n}$  de  $\mathcal{R}^p$  et de  $\mathcal{R}^n$ . Il existe :

- Une matrice  $U_{p \times r} = [u_1 \dots u_r]$   $Q$ -orthonormée  
les colonnes  $u_i$  sont les vecteurs propres de  $X^T D X Q = S Q = S_I$  associés aux valeurs propres  $\lambda_1 \geq \lambda_2 \geq \dots \lambda_r > 0$ .  
 $\Rightarrow$  Vecteurs principaux de l'espace des indiv.
- Une matrice  $V_{n \times r} = [v_1 \dots v_r]$   $D$ -orthonormée  
les colonnes st les vecteurs propres de  $X Q X^T D = W D$  associés aux valeurs propres  $\lambda_1 \geq \lambda_2 \geq \dots \lambda_r > 0$ .  
 $\Rightarrow$  Vecteurs principaux de l'espace des variables.
- Une matrice diagonale  $\Theta_{r \times r} = \Lambda^{1/2}$  contenant les  $r$  **valeurs singulières** du triplet (X, Q, D) tq X se décompose en :

$$X_{n \times p} = V_{n \times r} \Theta_{r \times r} (U_{p \times r})^T = \sum_{i=1}^r \sqrt{\lambda_i} v_i (u_i)^T$$

114

## Ce qu'il faut retenir

- Les questions préalables à toute analyse :**
  - Les unités de mesure. Comment les données ont été acquises ?
  - La structure du tableau de données, le type des variables.
- Résumé d'une variable :**
  - graphiques,
  - paramètres numériques : position, dispersion, forme.
- Résumé d'un tableau de variables :**
  - centre de gravité,
  - matrice de dispersion,
  - les deux espaces de représentation.
- DVS d'une matrice de données.**

115

## Exercices Td/TP ch I: calcul var, cor, DVS

### Données supplémentaires

Données	Description
Eaux1	Corpus 20 eaux minérales décrites par 7 variables
Eaux2	Corpus données supplémentaires
hemo_cir_quali.txt	Corpus 136 sujets décrits par 24 variables
Poumon.txt	Corpus 72 patients décrits par 7 variables
ChaZeb-a	Corpus de 23 bovins (Charolais & Zébus) décrits par 6 variables pondérales.
Kangourou	Corpus de 151 kangourous des 2 sexes appartenant à 3 espèces décrits par 18 variables quantitatives.
Loup	Description de 43 crânes Chien/Loup par 6 variables. Identification d'un crane d'origine inconnue.

Travail à rendre par écrit en binôme .... hemo\_cir\_quali.txt

Données : hémachromatose

=> A rendre lors du dernier TP semaine 43 ou 44

116



## Références

- L. Bellanger, R. Tomassone, *Exploration de données et méthodes statistiques : Data analysis & Data mining avec R. Collection Références Sciences*, Editions Ellipses, Paris, 2014.
  - J.-M. Bouroche & G. Saporta, *L'analyse des données*. Presses Universitaires de France : Que sais-je ? 85, Paris, 1992.
  - Site: [www.math.univ-montp2.fr/~durand](http://www.math.univ-montp2.fr/~durand)
- Support intitulé « Elts de Calcul matriciel et d'Analyse Factorielle de Données »*
- L. Lebart, A. Morineau, M. Piron, *Statistique exploratoire multidimensionnelle*. Dunod, Paris, 2006.
  - J.-P. Nakache, J. Confais, *Approche pragmatique de la Classification*. Editions Technip, Paris, 2005.
  - G. Saporta, *Probabilités, Analyse des données*. Editions Technip, Paris, 2006.
  - R Foundation, <http://www.r-project.org>
  - R. Tomassone, C. Dervin & J.-P. Masson. *Biométrie : modélisation de phénomènes biologiques*. Masson, Paris, 1993.