

---

# TRAITEMENT DE L'INFORMATION

---

# Contents

<b>1</b>	<b>Introduction</b>	<b>2</b>
1.1	Constitution du tableau de donnée . . . . .	2
1.2	Les variables et les individus . . . . .	2
1.2.1	les individus . . . . .	2
1.2.2	les variables . . . . .	2
1.3	Les types de tableaux . . . . .	3
1.3.1	Tableau de données quantitatives . . . . .	3
1.3.2	tableau de contingence . . . . .	3
1.3.3	tableaux binaires (logique d'incidence) . . . . .	3
1.3.4	Tableaux de préférences . . . . .	3
1.3.5	Tableaux de modalités . . . . .	3
1.3.6	tableaux de proximité . . . . .	3
1.4	Changement de variables . . . . .	4
1.5	Elements descriptifs d'un tableau de données . . . . .	4
1.6	Choix d'une mesure de ressemblance . . . . .	5
<b>2</b>	<b>Analyse en composantes principales : <i>ACP</i></b>	<b>6</b>
2.1	Principe de la méthode . . . . .	6
2.2	Formalisation du problème . . . . .	6
2.3	Resolution du problème . . . . .	6
2.3.1	Choix de $a$ . . . . .	6
2.3.2	choix de $[u_1, \dots, u_q]$ . . . . .	6

# Chapter 1

## Introduction

### 1.1 Constitution du tableau de donnée

Classe	Survivant	Nom	Sexe	Age	Parent présents	Ref billet	Prix billet	Port	
1	1	Allen Elisabeth	F	29	0	null	211	B5	S
3	0	Dyker Adolf	N	23	1	null	7	null	C
3	1	Dyker Anna	F	22	1	null	null	null	C
3	0	Emir Farell	N	null	0	null	null	null	null

Toute les variables doivent être du même type ( $\rightarrow$  quels sont les types de var ?)

- Quid des données manquantes ? ( $\rightarrow$  les remplacer ? par la moyenne ?)
- Quid des donnée aberrantes ?
- Lignes redondantes ? ( $\rightarrow$  à supprimer ?)
- Colonnes redondantes ?

### 1.2 Les variables et les individus

#### 1.2.1 les individus

Éléments de la population étudiée

#### 1.2.2 les variables

Application définie par

$$\omega(\text{population}) \longrightarrow O(\text{espace d'observation}) \text{ struct de cet espace}$$

1<sup>ère</sup> typologie

card O	struct de O	O ensemble continu de $\mathbb{R}$	O fini ou dénombrable	.
sans struct =, $\neq$		/	C8P lieu de residence	nominale
struct ordi. $\leq$		Age imperative	rang note	ordinal
corps ordonnée		salaire	/	mesurable
.		quantifiable	qualifiable	.

2<sup>ème</sup> typologie :  $\rightarrow$  variable d'incidence

$\rightarrow$  variable relationnelle

- variable d'incidence
  - $\rightarrow$  attribut descriptif, réponse oui = 1 OU non = 0
  - $\rightarrow$  var. numérique espace d'observation =  $\mathbb{R}$

- var. relationnelles  
 → var rang → échelle de notes suffisamment fine pour que 2 individu n'aient pas (pas de suite c'est effacé)  
 → var. présentant des modalités non ordonnées (bac, csp, ...)  
 → var. mesures sur  $\mathbb{R} \times \mathbb{R}$

## 1.3 Les types de tableaux

CN Tableaux doivent être homogènes 1 seul types de variable

### 1.3.1 Tableau de données quantitatives

(pourcentage de minerai par sondage / profondeur)

$$x_{i \leftarrow \text{ligne}}^{j \leftarrow \text{colonne}}$$

traitement

### 1.3.2 tableau de contingence

croisement de 2 var qualitatives

Traitement : Analyse fonctionnelle des correspondances

### 1.3.3 tableaux binaires (logique d'incidence)

tableaux d'attributs descriptif de 1,0

Traitement : classifications hiérarchiques

### 1.3.4 Tableaux de préférences

Notes données à des marques de parfum

Traitement : Analyse fonctionnelle des correspondances multiples

### 1.3.5 Tableaux de modalités

lisez vous tel journal ?

- 5 : tout le temps
- 4 : régulièrement
- 3 : parfois
- 2 : rarement
- 1 : jamais

### 1.3.6 tableaux de proximité

mesure sur  $\mathbb{R} \times \mathbb{R}$

	r1	r2	r3
r1			
r2			
r3			

## 1.4 Changement de variables

→ pour rendre le tableau homogène

	Sexe	couleur yeux
pere	M	marron
mere	F	bleu
enfant	M	vert

codage disjonctif complet

	sexe M	sexe F	yeux marr	yeux bleu	yeux vert
pere	1	0	1	0	0
mere	0	1	0	1	0
enfant	1	0	0	0	1

→ Tableau de Burt

${}^tD \times D$

	sexe M	sexe F	Y m	Y b	Y v
sexe M	2	0	1	0	1
sexe F	0	1	0	1	0
Y m	1	0	1	0	0
Y b	0	1	0	1	0
Y v	1	0	0	0	1

Tableaux d'effectifs → Analyse fonctionnelle des correspondances

autre exemple : Tableau de Burt du Titanic

	classe 1	classe 2	classe 3	enf	adul	F	M	surv 0	surv 1
classe 1	325	0	0	6	319	147	180	122	203
classe 2	0	285	0	24	261	106	179	167	118
classe 3	0	0	706	79	627	null	null	null	null
etc	⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮

## 1.5 Elements descriptifs d'un tableau de données

n : individus/lignes

p : variable/colonnes

individu  $i = x_i = (x_i^1 \quad \dots \quad x_i^p) \in \mathbb{R}^p$  la variable  $j : x_j = \begin{pmatrix} x_i^1 \\ \vdots \\ x_i^p \end{pmatrix} \in \mathbb{R}^n$

On va étudier le nuage  $N_I = \{x_i, i = 1, \dots, n, \text{poids } p_i = \frac{1}{n}\}$

$N_J = \{x_j, j = 1, \dots, p, \text{poids } p_j = \frac{1}{p}\}$

Soit le nuage des variables

A chaque variable, on peut associer: sa moyenne

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

sa variance  $Var \underline{x}^j = \sum_{i=1}^n p_i (x_i^j - \bar{x}^j)^2$

TROU

Pour le nuage  $N_I$  on peut calc le centre de gravité :

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

$$\underline{\bar{x}} = \begin{pmatrix} \bar{x}^1 \\ \vdots \\ \bar{x}^p \end{pmatrix}$$

Sur le nuage  $N_j$  on peut calculer le centre de gravité

$$\underline{\bar{x}} = \begin{pmatrix} \bar{x}_1 \\ \vdots \\ \bar{x}_n \end{pmatrix}$$

L'inertie du nuage

des individus  $I_{N_I} = \sum_{i=1}^n p_i d^2(x_i, \underline{\bar{x}}) \longrightarrow \underline{\text{distance}}$

## 1.6 Choix d'une mesure de ressemblance

définition 1 : indice de similarité

$$s\Omega \longrightarrow \mathbb{R}^+$$

s est un indice de similarité si :

1. s est symétrique,  $s(x, y) = s(y, x)$
2.  $s(x, x) = s(y, y) \geq s(x, y)$

définition 2 : indice de dissimilarité

$$s\Omega \longrightarrow \mathbb{R}^+$$

s est un indice de dissimilarité si :

1. symétrique
2.  $s(x, x) = 0$

définition 3 : une distance est un indice de dissimilarité qui vérifie en plus

1.  $s(x, y) = 0 \Leftrightarrow x = y$
  2.  $s(x, z) \leq s(x, y) + s(y, z)$
- distance euclidienne (A en composantes principales)
  - distance dite du  $\chi^2$  (Analyse factorielle des correspondances)

$$d^2(\underline{x}_i, \underline{x}_{i'}) = \sum_{j=1}^p \frac{1}{x_i^j} \left( \frac{x_i^j x_i}{x_{i'}^j} - x_{i'}^j \right)^2$$

- Tableau d'incidence
- TROU

Distance entre groupes

(indice d'appréparation)

distance du lien max

distance du lien min

distance des centres de gravité

trou mais fin de chap

distance fréquemment utilisé  
 $\delta(A, B) = I \quad \delta(A, B) = d(\underline{\bar{x}}, \underline{\bar{y}})$

## Chapter 2

# Analyse en composantes principales : *ACP*

Traite les tableaux de données quantitativement positives

### 2.1 Principe de la méthode

Le pb qui se pose dans la mesure où 1 partie des variables sont liées, c'est de passer d'un tableau X de dimension  $n \times p$  à un tableau Y de dim  $n \times q$   $q < p$  en réduisant le nbr de variable descriptive tout en perdant le moins possible d'info

la methode utilisé en ACP pour passer de p var à q var q<p consiste à projeter le nuage des individus sur un sous espace W de dim q en deformant le moins possible le nuage lors de sa projection

Les composantes principales sont les nouvelles variables. elles vont s'interpreter comme des" 'synthèses' des variables initiales  $\underline{x}^j$ .

TROU flemme 3 lignes + fin du cours

Le principe de l'ACP, c'est de réduire le nbr de variables décrivant les individus en perdant le moins possible d'information

### 2.2 Formalisation du problème

Rechercher l'espace de projection  $W$  ( $a$  : origine,  $u_1, \dots, u_q$  : système générateur de  $W$ ) de telle façon que la perte d'info fut minimale

$$\forall i \hat{x}_i = \underline{a} + \sum_{k=1}^q y_i^k \underline{u}_k$$

$\hat{x}_i$  est la projection orthogonale de  $x_i$  sur  $W$

On cherche le referentiel  $W$  tq  $\sum_W = \sum_{i=1}^n p_i d^2(x_i, \hat{x}_i)$

### 2.3 Resoplution du problème

#### 2.3.1 Choix de a

#### 2.3.2 choix de $[\underline{u}_1, \dots, \underline{u}_q]$