

# Classification Automatique.

Avant propos.

Le cours est adressé aux étudiants de la 2<sup>ème</sup> année STID.

Trois notions essentielles seront données dans ce cours

- Partition
- Hierarchie
- Discrimination.

Nous donnons des méthodes simples de construction d'une partition, une hiérarchie et fonction discriminante.

Pour ce but, le premier chapitre présente les notions de base et définitions nécessaires.



# Classification automatique

## I Introduction

La classification automatique est un axe des statistiques qui regroupe un ensemble de méthodes permettant de rassembler ds une même classe les elts (objets ou individus) qui se ressemblent.

Citons : classification automatique par :

- partition



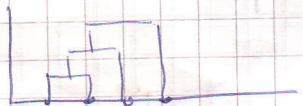
Partition = class disj dont  $U = E$

- Arbre



grand-père — fils — petits fils —

- Hierarchie



departement — faculté — université — Académie

- Classes empiétantes : les classes peuvent se couper.

Nous allons nous intéresser à deux entités : partition basée sur la ressemblance : si deux inds se ressemblent ils seront ds la même classe sinon ils seront ds des classes  $\neq$ . Et la hiérarchisation.

Les données sont souvent resumées ds un tabl. ind  $\times$  variables :  $n$  inds caractérisés par  $p$  variables.

On resume les étapes comme suit :

étape 1 : Choix des données : la population à étudier

étape 2 : Construction d'un tabl de ressemblance, en calculant la ressemblance ou dissimilarité entre tous les inds  $i$  à  $j$ .

étape 3 : Choisir un algorithme de classification : phase de regroupement d'inds.

étape 4 : Interprétation des résultats.

## II Définition de la distance et dissimilarité

def 1 : Soit  $E$  une population de  $n$  inds ou objets.

On appelle indice de dissimilarité  $d$ , toute application

$$d: E \times E \longrightarrow \mathbb{R}_+ \\ (x, y) \longmapsto d(x, y) \quad \text{tp}$$



•  $\forall (x, y) \in E \times E$   $d(x, y) = d(y, x)$  symétrie.

•  $\forall x \in E$   $d(x, x) = 0$ .

si de plus  $d$  vérifie l'inégalité triangulaire

$$\forall x, y \text{ et } z \in E \quad d(x, z) \leq d(x, y) + d(y, z)$$

alors  $d$  est une distance.

Rp: Une distance est donc un indice de dissimilarité, l'inverse est faux.

def 2: On appelle indice de similitude  $s$ , toute application

$$s: E \times E \longrightarrow \mathbb{R}_+ \\ (x, y) \longmapsto s(x, y) \quad \text{tp}$$

•  $d(x, y) = d(y, x)$  symétrie

•  $d(x, x) \geq d(x, y)$ .

## II.2) Types de distances et dissimilarité

Soit un tab  $X$  regroupant  $n$  inds opérés par  $k$  variables quantitatives

(a) Distance (euclidienne).

A chaque ind  $x \in E$ , on associe un  $k$ -uplet  $(x_1, \dots, x_k) \in \mathbb{R}^k$  où  $x_j$  représente la mes par la  $j^{\text{ème}}$  variable.

- La distance euclidienne entre deux inds  $x, y$  est donnée par

$$d(x, y) = \sqrt{\sum_{j=1}^k (x_j - y_j)^2}$$

On note par  $\sigma_j = \sqrt{\text{Var}(x_j)}$  écart type de la  $j^{\text{ème}}$  var.

- La distance euclidienne normalisée est donnée par.

$$d_N(x, y) = \sqrt{\sum_{j=1}^k \frac{1}{\sigma_j^2} (x_j - y_j)^2}$$

- La distance  $L_1$ :

$$d_{L_1}(x, y) = \sum_{j=1}^k |x_j - y_j|$$

- La distance Minkowski  $L_q$

$$d_{L_q}(x, y) = \left( \sum_{j=1}^k (x_j - y_j)^q \right)^{1/q} \quad q=2 \quad L_2 \text{ distance euclidienne}$$



- La distance de Mahalanobis.

$$d(x, y) = \sqrt{(x - y)^T V^{-1} (x - y)}$$

ou  $V$  est la matrice de var-covar de  $k$  vars.  $V$  de dim  $k \times k$

- La distance de khi-2.

$$d(x, y) = \sqrt{\sum_{j=1}^k \frac{n}{n_{.j}} \left( \frac{n_{xj}}{n_{x.}} - \frac{n_{yj}}{n_{y.}} \right)^2}$$

utilise ds le cas des tabs de contingence, ou

$n_{xj}$  est la fréquence de  $(x, j)$ .

$$n_{x.} = \sum_{j=1}^k n_{xj}$$

b) ~~similitude~~ similarité

sont une population de  $n$  inds à classifier et décrits par ~~une~~ variables qualitatives ou  $k$  modalités (caractéristiques)

On note a 1 si l'ind possède la caractéristique et par 0 sinon.

Une dissimilarité est obtenue en combinant les similarités présentée par les nombre  $a, b, c$  et  $d$  suivants:

$a =$  nombre de caractéristiques commune à deux inds  $i_1$  et  $i_2$ .

$b =$  " " possède par  $i_1$  et non par  $i_2$ .

$c =$  " " " par  $i_2$  et non par  $i_1$ .

$d =$  " " qui ne sont possédés ni par  $i_1$  ni par  $i_2$ .

Plusieurs de similarité ont été proposées, utilisant tous les  $n^{br}$   $a, b, c, d$

- La similarité de Jaccard :

$$s(i_1, i_2) = \frac{a}{a + b + c}$$

- La similarité de Dice ou de Brakanowski

$$s(i_1, i_2) = \frac{2a}{2a + b + c}$$

- La similarité d'Ochiai

$$s(i_1, i_2) = \frac{a}{\sqrt{(a+b)(a+c)}}$$

Yull

$$\frac{ad - bc}{ad + bc}$$

Pearson  $\frac{|ad - bc|}{[(a+b)(c+d)(a+c)(b+)]}$



- la similitude de Rogers et Tanimoto

$$s(x_1, x_2) = \frac{a+d}{a+d+2(b+c)}$$

- la similitude de Russel et Rao

$$s(x_1, x_2) = \frac{a}{a+b+c+d} = \frac{a}{b}$$

- Sokal et Sneath (Sneath)

$$\frac{a}{a+2(b+c)}$$

- Sokal et Michener

$$\frac{a+b}{a+b+c+d}$$

- Kulzinsky  $\frac{a}{a+b} + \frac{a}{a+c}$

Rp: toutes ces similitudes sont comprises entre 0 et 1.

En A chaque indice de similitude correspond un indice de dissimilarité, défini par  $d = 1 - s \geq 0$ .

Exemples

1/  $E = \{0, 1, 2, 3\}$ .

$d: E \times E \rightarrow \mathbb{R}^+$

$d(1, 0) = d(0, 1) = 1$        $d(0, 2) = d(2, 0) = 2$ .

$d(1, 3) = d(3, 1) = 4$        $d(0, 3) = d(3, 0) = 3$ .

$d(1, 2) = d(2, 1) = 3$

$d(2, 3) = d(3, 2) = 6$

$d(0, 0) = d(1, 1) = d(2, 2) = d(3, 3) = 0$

$d$  est un indice de dissimilarité mais pas une distance.

$d(2, 3) = 6$

$d(2, 0) + d(0, 3) = 2 + 3 = 5 < 6$

l'inégalité triangulaire n'est pas vérifiée.

On définit le tableau de dissimilarité suivant:

E	0	1	2	3
0	0	1	2	3
1	1	0	3	4
2	2	3	0	6
3	3	4	6	0

symétrique de diagonale nulle.



ex 2/  $E = \{0, 1, 2, 3\}$ .

•  $d(0,0) = d(1,1) = d(2,2) = d(3,3) = 0$ .

•  $d(0,1) = 1 = d(1,0)$     •  $d(1,2) = d(2,1) = 1$     •  $d(2,3) = d(3,2) = 1$

•  $d(0,2) = d(2,0) = 2$     •  $d(1,3) = d(3,1) = 2$

•  $d(0,3) = d(3,0) = 3$

Le tableau de distances est

	0	1	2	3
0	0	1	2	3
1	1	0	1	2
2	2	1	0	1
3	3	2	1	0

= distance entre 1 et 2.

ex 3/ Soit 4 inds à qui on a posé 4 questions A, B, C et D auxquelles ils doivent répondre par oui ou non.

On note par 1 si la réponse est oui et par 0 si la réponse est non et soit le tabl de données suivant:

Quests \ inds	Q <sub>1</sub>	Q <sub>2</sub>	Q <sub>3</sub>	Q <sub>4</sub>
i <sub>1</sub>	non	oui	oui	oui
i <sub>2</sub>	0	0	1	1
i <sub>3</sub>	1	0	0	1
i <sub>4</sub>	1	1	0	1

Calculons les quantités a, b, c et d entre i<sub>1</sub> et i<sub>2</sub>.

$a = 2$     2 réponse oui communes.

$b = 1$     i<sub>1</sub>: Q<sub>1</sub> = 1 et Q<sub>2</sub> = 0.

$c = 0$

$d = 1$      $Q_3^{i_1} = Q_3^{i_2} = 1$

l'indices de similarité de entre i<sub>1</sub> et i<sub>2</sub>.

• Jaccard     $s(i_1, i_2) = \frac{2}{2+1+0} = \frac{2}{3}$

• Dice ou Cra-Kanowski     $s(i_1, i_2) = \frac{4}{4+1+0} = \frac{4}{5}$

• Ochiaia     $s(i_1, i_2) = \frac{2}{\sqrt{(2+1)(2+0)}} = \frac{2}{\sqrt{6}}$

• Rogers et Tanimoto     $s(i_1, i_2) = \frac{2+1}{2+1+2(1+0)} = \frac{3}{5}$