

TD d'analyse de données. 3

1. Pour un tableau N de contingence comportant m_1 lignes et m_2 colonnes, formé de valeurs entières : $N = [n^{(1)} \ n^{(2)} \ \dots \ n^{(m_2)}] = {}^t [n_1 \ n_2 \ \dots \ n_{m_1}]$, on note $n_{i\cdot}$ la somme des éléments de la i -ème ligne n_i , et $n_{\cdot j}$ la somme des éléments de la j -ème colonne $n^{(j)}$ du tableau. On note n la somme de tous les éléments de N , et D_1 désigne la matrice diagonale d'éléments $n_{1\cdot}, n_{2\cdot}, \dots, n_{m_1\cdot}$. Indiquez la forme du profil $D_1^{-1} N$, et du centre de gravité $g = \frac{1}{n} {}^t (D_1^{-1} N) D_1 {}^t 1_{m_1}$. On définit la distance d du Khi deux entre deux lignes i et i' de N par :

$$d^2(i, i') = \sum_{j=1}^{m_2} \frac{n}{n_{i\cdot} n_{i'\cdot}} \left(\frac{n_{ij}}{n_{i\cdot}} - \frac{n_{i'j}}{n_{i'\cdot}} \right)^2$$

- i) Vérifiez que d est bien une distance sur N^{m_1}

- ii) Vérifiez que dans l'écriture de d^2 on peut remplacer les valeurs n_c par les fréquences $f_c = \frac{n_c}{n}$

- iii) Qu'obtient-on si on échange les rôles des lignes et des colonnes de N ?

2. Si d est une application de $E \times E$ dans \mathbb{R} qui vérifie $d(x, y) = d(y, x) \geq 0$, $\forall (x, y) \in E^2$, d est appelée *dissimilarité* si $\forall x \in E$, $d(x, x) = 0$, et *similarité* si $\forall (x, y) \in E^2$, $d(x, x) \geq d(x, y)$. Une similarité peut-elle être une distance ? Donnez des exemples de similarités et de dissimilarités (autres que des distances).

3. Indiquez la forme du produit (est-ce un produit scalaire ?) en dimension 2, associé à la matrice $M = \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix}$. Quelle est la condition pour que la 'distance' entre deux points soit nulle ?

4. Comment peut s'écrire l'inertie d'un nuage de M points ? Comment définir une inertie intra classe et une inertie inter classes ? Quelles sont les relations qui relient ces valeurs ?

5. On souhaite comparer le fonctionnement d'une classification hiérarchique et celui d'une méthode k -means. Observez le comportement de ces deux algorithmes en dimension 1 – on pourra prendre des valeurs entières, comme les n premiers entiers, ou leurs carrés. Considérez des exemples en dimension 2, en proposant des critères d'agrégation des classes, pour la distance euclidienne dans le plan.

6. Soit un nuage de 10 points : $(0, 0); (2, 0); (4, 0); (8, 0); (0, 1); (0, 6); (8, 2); (3, 3); (6, 5); (2, 8)$. En utilisant la distance euclidienne du plan, écrivez une Classification Ascendante Hiérarchique jusqu'à obtenir 3 classes exactement. Ecrivez pour l'ensemble des étapes (agrégation de classes) la valeur de l'inertie intra classe totale. Comparez avec une approche par moyennes mobiles pour 3 classes, en prenant des centres initiaux dans des classes distinctes de la CAH.

7. Soit $A = (a_{ij})$ une matrice d'appartenance, ie une matrice binaire, pour laquelle $a_{ij} = 1$ signifie que le i -ème individu (ligne i) appartient à la classe no j . Comment peuvent être écrits les effectifs des différentes classes ? Quelle est la condition pour qu'un individu appartienne à au plus une classe (respectivement exactement une classe) ? Peut-on représenter cette situation par un graphe ? Peut-on utiliser une matrice d'appartenance non déterministe, dont les éléments sont des nombres compris entre 0 et 1 ?

Analyse de données: TDS

Exo 2

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

- i) $d(x, y) = d(y, x)$
ii) $d(x, y) = 0 \Rightarrow x = y \rightarrow$ dissim
iii) $d(x, y) \leq d(x, x) \rightarrow$ sim
iv) $d(x, x) = 0 \rightarrow$ dissim
v) $d(x, z) \leq d(x, y) + d(y, z) \rightarrow$ nécessaire pour une distance

Si on prend à la fois
sim et dissim on a que
 $d(x, x) = 0$ et $d(x, x)$
majore \rightarrow fonction nulle

distance = cas particulier de dissim

Exemple de sim

$$\begin{cases} d(x, x) = 1 \\ d(x, y) = 0 \text{ si } x \neq y \end{cases}$$

$$d(x, y) = \min(|x|, |y|)$$

$$C_{x,y} = \frac{\text{cov}(x,y)}{\sqrt{V(x)V(y)}}, \quad \text{cov}(x,y) = E([X-E(X)] \cdot [Y-E(Y)])$$

$$\hookrightarrow -1 \leq C_{x,y} \leq 1 \Rightarrow d(x,y) = |C_{x,y}|$$

$$C_{x,x} = 1 \text{ car } \text{cov}(x,x) = V(x)$$

Exemple dissim

$$d(x,y) = (x-y)^2 \rightarrow \text{inégalité plus vraie}$$

• Faire une composition avec une distance ne donne pas
une distance mais une similarité en général

Exo 3

produit scalaire associé à M :

$$\langle u, u' \rangle = {}^t M u u'$$

$$M = \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix}, u = \begin{bmatrix} x \\ y \end{bmatrix}, u' = \begin{bmatrix} x' \\ y' \end{bmatrix}$$

$$\langle u, u' \rangle = \begin{bmatrix} x & y \end{bmatrix} \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix} \begin{bmatrix} x' \\ y' \end{bmatrix} = xx' + xy' + yx' + yy'$$

$$\begin{aligned} \langle u, u \rangle &= x^2 + 2xy + y^2 \\ &= (x+y)^2 \end{aligned}$$

$$\langle u, u \rangle = 0 \Leftrightarrow x+y=0$$

$$\Leftrightarrow x = -y \Leftrightarrow x=y=0$$

Ce n'est donc pas un produit scalaire (car $\langle u, u \rangle = 0 \nRightarrow u = 0$ en fait)

Donc $\|u\| = \sqrt{\langle u, u \rangle}$ ne définit pas une norme.

et donc pas réellement une distance

$$d(u, u') = \|u - u'\| = 0$$

$$\text{Si } \langle u - u', u - u' \rangle = 0$$

$$u - u' = \begin{bmatrix} x - x' \\ y - y' \end{bmatrix}$$

$$(x - x') + (y - y') = 0 \Rightarrow x + y = x' + y'$$

vérifie que $d(x, x) = 0$ et $d(x, y) = d(y, x) \rightarrow \underline{\text{ok}}$
 \hookrightarrow dissim

Et les inégalités triangulaires? Apparemment ça vérifie

Exo 5 : Classifications hiérarchique

Pour une classif hiérarchique:

$d(x, y) \leftarrow$ * définition d'une fonction d'écart (dissim en général)

$\delta(c, c') \leftarrow$ * définition d'une fonction d'écart entre classes

* Parmi points on définit m classes $C_k = \{x_k\}$, $k=1, \dots, m$
 base de classe fixé \rightarrow ou un absc $K \geq 1$ pour être plus pertinent

* while $(k > 1)$ on fusionne les classes les \oplus proches
 ($\delta(c, c')$ minimal), donc $k--$

Dans le cas où on repz par des binaires

$n = 5$

$E = \{x_1, x_2, x_3, x_4, x_5\}$

$C_3 = \{x_3, x_4\} \rightarrow 00110$

→ Au début la matrice C des classes repz
 par un binaire est I_n (quand on fait n classes
 de 1 pr au début) puis ça évolue

$C = \begin{bmatrix} \text{class 1} \\ \text{class 2} \\ \vdots \end{bmatrix} \rightarrow$ voir "longer" ci

Fusionner les classes $c_i, c_j \rightarrow c_i + c_j$ avec
 si la suite qui repz
 sa classe

Moyenne des éléments de C

3x arcs (m, d)
 dans mathlab

$$S = [11 \ 0 \dots 0 \ 0] = C_i$$

nbre elem de la
 classe

$$nC = S \times \begin{bmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{bmatrix}$$

\rightarrow car $S =$ présence des élem dans la classe, x_i à x_n
 transfo en sa valeur pour multiplier par
 un vecteur obtenu full 1 le fait.

centre des
 classes

$$\bar{C} = \frac{1}{nC} (u_1 + u_2 + \dots + u_n) \rightarrow \text{SOMME}$$

$$E = \{u_1, u_2, \dots, u_n\}, u_i = \begin{pmatrix} x_i \\ y_i \end{pmatrix}$$

$$\bar{C} = \begin{bmatrix} \frac{1}{nC} \cdot S \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix} \\ \frac{1}{nC} \cdot S \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix} \end{bmatrix}$$

ex (n) \rightarrow crée I_m sous mathlab

on peut s'arrêter sur autre chose que le nombre de classes:

→ inertie intra classe

$$C: \sum_{u,v \in C} d(u,v)^2 = i_C$$

La condition :

$$\left(\frac{\sum i_C}{C} \right) \geq 1$$

distancé entre tous les points d'une même classe
somme des inerties de classes
divise
inertie totale (élément
distancé entre tous les points)

$\in [0,1]$ ← avec que des classes taille 1 = 0
avec 1 classe = 1
et croissant

• Les k-means

→ nombre de classes fixés K

→ choix de K points → $C = \{x_i\}$ classe à 1 elem

$$F = \{u_1, u_2, \dots, u_K\} \text{ et } C_k = \{u_k\}$$

→ tant que $E \setminus F \neq \emptyset$

choisir $u \in E \setminus F$

associer u à une des K classes

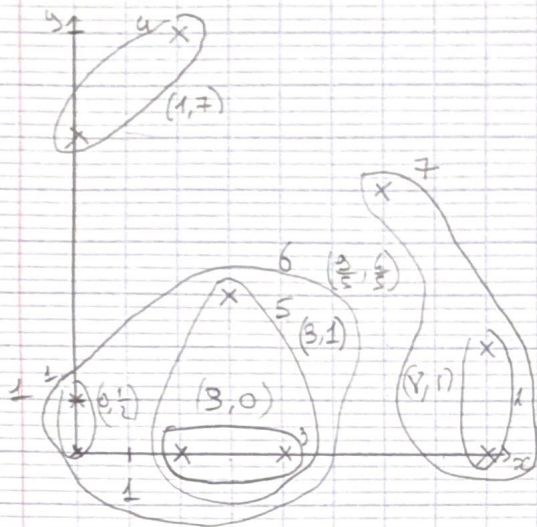
$$F = F \cup \{u\}$$

minimisation des distances

On peut obtenir un meilleur résultat en prenant le résultat d'une classif hiérarchique de condition d'arrêt K . On utilise les moyennes des classes obtenues. (on "remplace" des pts par des moyennes)
Mieux que de prendre des points au pif.

Analyse de données : TD 3

Exo 6



On a 10 classes de 1 point.
On va regrouper par deux

$(0,0)$ et $(0,1) \rightarrow \{(0,0), (0,1)\}$

$$\rightarrow \bar{m} = (0, \frac{1}{2})$$

$(8,0)$ et $(8,2) \rightarrow \{(8,0), (8,2)\}$

$$\rightarrow \bar{m} = (8, 1)$$

$(2,0)$ et $(4,0) \rightarrow \{(2,0), (4,0)\}$

$$\rightarrow \bar{m} = (3, 0)$$

$(0,6)$ et $(2,8) \rightarrow \{(0,6), (2,8)\}$

$$\rightarrow \bar{m} = (1, 7)$$

$$5: \bar{m} = \frac{(2,0) + (4,0) + (3,3)}{3} = (3, 1)$$

$$6: \bar{m} = \frac{(2,0) + (4,0) + (3,3) + (0,0) + (0,1)}{5} = (\frac{9}{5}, \frac{4}{5})$$

↳ Reste 4 classe à cette étape

$$7: \bar{m} = \frac{(8,0) + (8,2) + (6,5)}{3} = (\frac{22}{3}, \frac{7}{3})$$