

Méthodes de classifications

Dhafer Malouche

ESSAI-U2S-ENIT

<http://essai.academia.edu/DhaferMalouche>

dhafer.malouche@me.com

Juin-2013

ANR—Do Well Be, Saint Nectaire

Plan

Introduction

Classification hiérarchique

Partitionnement

Indices de validité

Méthodes probabilistes pour la classification

Introduction

Classification hiérarchique

Partitionnement

Indices de validité

Méthodes probabilistes pour la classification

Définition

- ▶ Algorithmes qui permettent de classer des objets observés dans des classes (appelées *clusters*).
- ▶ Les objets d'une même classe doivent être "similaires" et les objets de deux classes différentes doivent être "distincts".
- ▶ Différents algorithmes de classification
 - ▶ Méthodes Hiérarchiques (*classification hiérarchiques...*)
 - ▶ Partitionnement (*k-means..*)
 - ▶ Méthodes Probabilistes (*EM-algorithme...*) et encore d'autres...

Les objets

- ▶ $E = \{w_1, \dots, w_n\}$ un ensemble de n objets (individus) pour lesquels on a observé d "variables" (qualitatives et/ou quantitatives) x^1, \dots, x^d où

$$x_i^j = x^j(w_i)$$

est la $i^{\text{ème}}$ observation de w_i selon x^j .

- ▶ On considère $w_i = (x_{i,t})_{t \in T_i}$ $i = 1, \dots, n$ des signaux ou séries temporelles (ST) .
- ▶ Une matrice $D = (d_{ij'})$, $n \times n$ (n est le nombre d'objets à classer), où $d_{ij'}$ est une mesure de **similarité** entre deux objets w_i et $w_{j'}$,

Etape 1 : calcul des similarités

Cas de variables quantitatives.

- ▶ Distance *Euclidienne*

$$d(w_i, w_{i'}) = \sqrt{\sum_{j=1}^d (x_i^j - x_{i'}^j)^2}.$$

- ▶ ou la distance de *Minkowski* pour tout $p > 0$

$$d(w_i, w_{i'}) = \left(\sum_{j=1}^d (x_i^j - x_{i'}^j)^p \right)^{1/p}.$$

- ▶ encore d'autres exemples de distances : *Canberra*, *Manhattan*...
- ▶ Sous R : `dist`, il faut spécifier
 - ▶ une matrice de données quantitative
 - ▶ le type de distance

Distance entre signaux ou ST

- ▶ Distance Euclidienne ou l'une des distances déjà définies (Séries sont toutes de même "longueurs").
- ▶ *Dynamic time wrapping* ou Déformation temporelle dynamique (par Keogh & Pazzani, 2001¹)
 - ▶ Permet de comparer deux séquences de ST pas forcément de même longueur.
 - ▶ Dilatation ou compression des séquence pour obtenir le meilleur alignement possible
 - ▶ Sous R, package dtw.

¹Keogh, E. J. and Pazzani, M. J. (2001). Derivative dynamic time warping. In the 1st SIAM Int. Conf. on Data Mining (SDM-2001), Chicago, IL, USA.

Introduction

Classification hiérarchique

Partitionnement

Indices de validité

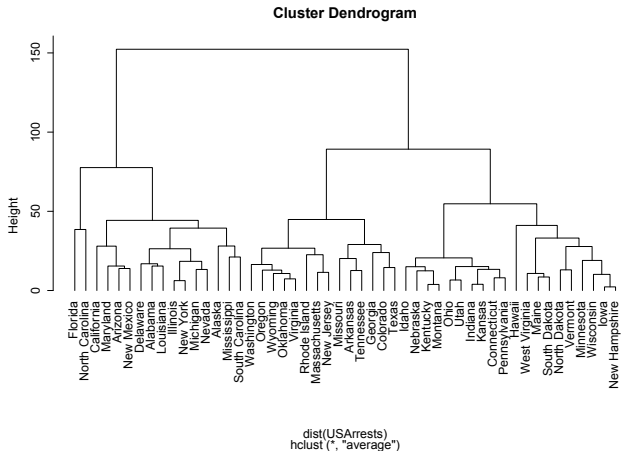
Méthodes probabilistes pour la classification

Algorithme ascendant (agrégatif) hiérarchique

1. $k = n$ le nombre de classes
2. Chaque classes contient une observation : n classes C^1, \dots, C^n .
3. Calculer la matrice de distances **entre classes** (phase agrégative de l'algorithme).
4. Considérer les deux classes C^i et C^j telles que

$$d(C^i, C^j) = \min_{(i', j')} d(C^{i'}, C^{j'})$$

- .
5. Agréger C^i et C^j dans une seule classe C^{ij} . Donc $k = k - 1$.
 6. $k = 1$? Si oui **STOP**, sinon, **aller à 3**.



Méthodes d'agrégations

$$d(C^l, C^{ij}) = \alpha_i d(C^l, C^i) + \alpha_j d(C^l, C^j) + \beta d(C^i, C^j) + \gamma |d(C^l, C^i) - d(C^l, C^j)|$$

Algorithme de classification	α_i	α_j	β	γ
<i>Single Linkage</i>	1/2	1/2	0	-1/2
<i>Complete Linkage</i>	1/2	1/2	0	1/2
<i>Centroïd Linkage</i>	$\frac{n_j}{n_i+n_j}$	$\frac{n_i}{n_i+n_j}$	$-\frac{n_i n_j}{n_i+n_j}$	0
Méthode de Ward	$\frac{n_j+n_l}{n_i+n_j+n_l}$	$\frac{n_i+n_l}{n_i+n_j+n_l}$	$-\frac{n_l}{n_i+n_j+n_l}$	0

Méthodes d'agrégations

- ▶ *Single Linkage* : $d(C^l, C^{ij}) = \min(d(C^l, C^i), d(C^l, C^j))$
- ▶ *Complete Linkage* : $d(C^l, C^{ij}) = \max(d(C^l, C^i), d(C^l, C^j))$
- ▶ *Centroid Linkage* : $d(C^l, C^{ij}) = d(m_l, m_{ij})$ où m_l et m_{ij} sont resp. les centres de gravité de C^l et C^{ij} .
- ▶ *Méthode de Ward* (après une ACP²) : elle a pour objectif de minimiser l'accroissement dans la variance intra-classe :

$$E = \sum_{l=1}^k \sum_{x_i \in C^l} d(x_i, m_l)^2.$$

²Analyse en composantes principales

Algorithme descendant hiérarchique (DIANA).

1. $k = 1$, Toutes sont dans la même classe C_1
2. Considérer C_i la classe de diamètre maximale ($d(C_i) = \max_{w, w' \in C_i} d(w, w')$).
3. Introduire $C_j = \emptyset$.
 - 3.1 Pour chaque $w \in C_i$, calculer

$$d(w, C_i \setminus \{w\}) = \frac{1}{|C_i| - 1} \sum_{w' \in C_i \setminus \{w\}} d(w, w').$$

- 3.2 Calculer la différence

$$\Delta(w, C_i, C_j) = d(w, C_i \setminus \{w\}) - d(w, C_j)$$

4. Si $w_m = \operatorname{argmax}_{w \in C_i} d(w, C_i \setminus \{w\})$ et $\Delta(w_m, C_i, C_j) > 0$ alors $w_m \ni C_j$
5. si $k < n$ alors $k = k + 1$ revient à 1/ sinon STOP.

Introduction

Classification hiérarchique

Partionnement

Indices de validité

Méthodes probabilistes pour la classification

- ▶ **Objectif** A partir d'une fonction critère J donnée, partitionner w^1, \dots, w_n en K classes (K est fixé d'avance) en minimisant ou maximisant la fonction J
- ▶ Méthode exhaustive : calculer pour toute les partitions possible C , le critère $J(C)$ et la solution est alors

$$C_{sol} = \operatorname{argmax} J(C).$$

- ▶ Solution impossible si n devient grand. En effet le nombre de partitions possible est égal à

$$P(n, K) = \frac{1}{K!} \sum_{m=1}^K (-1)^{K-m} C_K^m m^n.$$

Pour $n = 30$, et $K = 3$, $P(n, K)$ est de l'ordre de 2×10^{14} .

K -means

- ▶ Critère à minimiser : la somme des carrées des erreurs (SSE) :

$$J(C) = \sum_{c \in C} \sum_{w \in c} d(w, m(c))$$

où $m(c)$ est le centre de gravité de c .

- ▶ Algorithme :

1. Initialise K -centres m_1, \dots, m_K et
2. Affecter chaque individu vers la classe la plus proche

$$C = \{c_1, \dots, c_K\} :$$

$$w \in c_l \text{ si } d(w, m_l) = \min_{l'=1, \dots, K} d(w, m_{l'}).$$

3. Re-calculer les centres de gravité des nouvelles classes,
4. Répéter 2 et 3 jusqu'à convergence.

Introduction

Classification hiérarchique

Partionnement

Indices de validité

Méthodes probabilistes pour la classification

Mesure de la validité d'une classification

Deux types de mesures :

- ▶ Indices de validité **interne** : à partir d'un résultat de classification et de l'information intrinsèque dans la base de données on mesure la qualité du résultat :
 - ▶ indice de connectivité
 - ▶ indice de silhouette
 - ▶ indice de *Dunn*

- ▶ Indice de validité **externe** : comparer la classification obtenue avec une autre classification (soit une partition déjà connue ou une autre classification obtenue par une autre méthode)

Notations

Soit

- ▶ n le nombre des observations,
- ▶ d le nombre des variables,
- ▶ x_1, \dots, x_n les vecteurs d'observations à classer,
- ▶ $\mathcal{C} = \{C_1, \dots, C_k\}$ une classification à tester,
- ▶ $C(i)$ est la classe contenant x_i , i.e., $C(i) \ni x_i$.

Pour tout x_i , et $j \in \{1, \dots, n\} \setminus \{i\}$: $x_i^{(j)}$ est le $j^{\text{ième}}$ voisin le plus proche de x_i :

$$d(x_i, x_{(1)}) \leq d(x_i, x_{(2)}) \leq \dots d(x_i, x_{(j)}) \leq \dots \leq d(x_i, x_{(n-1)})$$

On pose pour tout $i, j \in \{1, \dots, n\}$

$$\delta(x_i, x_i^{(j)}) = \begin{cases} 0 & \text{si } x_i, x_i^{(j)} \in C(i) \\ 1/j & \text{sinon} \end{cases}$$

Indices de Connectivité

Pour une classification \mathcal{C} ,

$$\text{Conn}(\mathcal{C}) = \sum_{i,j=1,\dots,n} \delta(x_i, x_i^{(j)})$$

$\text{Conn}(\mathcal{C})$ prend ses valeurs dans $[0, +\infty[$ et c'est un coefficient qui doit être minimisée :

$\text{Conn}(\mathcal{C})$ proche de zéro \iff une bonne connectivité à l'intérieure des classes.

R : clValid package, connectivity

Largeur de la silhouette

- ▶ C'est la moyenne de toutes les valeurs de la silhouette.
- ▶ Valeur de la silhouette pour une observation x_i :

$$S(i) = \frac{b_i - a_i}{\max(b_i, a_i)}$$

où

- ▶ a_i est la distance moyenne entre x_i et tous les autres $x_j \in C(i)$.

$$a_i = \frac{1}{n(C(i))} \sum_{x \in C(i)} d(x_i, x).$$

- ▶ b_i est la distance moyenne entre x_i et les observations dans la classe la plus proche de x_i :

$$b_i = \min_{C \in \mathcal{C} \setminus C(i)} \frac{1}{n(C)} \sum_{x \in C} d(x_i, x).$$

Largeur de la silhouette

- ▶ Définie par

$$S(C) = \frac{1}{n} \sum_i S(i)$$

- ▶ La largeur de la silhouette est un nombre qui est toujours compris entre -1 et $+1$.
- ▶ S doit être maximisée.
- ▶ une valeur de $S(i)$ négative indique que cet individu n'est pas dans "la bonne classe" et il pourrait être déplacé vers la classe la plus proche.
- ▶ R, package cluster, silhouette

Indices externes

- ▶ \mathcal{P} une partition existante de l'ensemble des individus E ,
 $\mathcal{P} = \{G_1, \dots, G_r\}$.
- ▶ $\mathcal{C} = \{C_1, \dots, C_k\}$ une classification obtenue à l'aide d'un algorithme de classification (*hiérarchique, kmeans...*)
- ▶ Objectif : comparer la classification \mathcal{C} et la partition \mathcal{P}
- ▶ Définition de quelques indices de validité externes : package `clv` sur R

Calcul des indices externes

Soient deux individus x_i et $x_{i'}$. Quatre cas sont possibles

Cas 1 $\exists C \in \mathcal{C}$ et $\exists G \in \mathcal{P}$ tel que

$$x_i, x_{i'} \in C \text{ et } x_i, x_{i'} \in G$$

Cas 2 $\exists C \in \mathcal{C}$ et $\exists G \neq G' \in \mathcal{P}$ tel que

$$x_i, x_{i'} \in C \text{ et } x_i \in G, x_{i'} \in G'$$

Cas 3 $\exists C \neq C' \in \mathcal{C}$ et $\exists G \in \mathcal{P}$ tel que

$$x_i \in C, x_{i'} \in C' \text{ et } x_i, x_{i'} \in G$$

Cas 4 $\exists C \neq C' \in \mathcal{C}$ et $\exists G \neq G' \in \mathcal{P}$ tel que

$$x_i \in C, x_{i'} \in C' \text{ et } x_i \in G, x_{i'} \in G'$$

Calcul des indices externes

Notons par \mathcal{A}_l l'ensemble des paires $(x_i, x_{i'})$ vérifiant le cas l et par $a_l = |\mathcal{A}_l|$. Trois indices :

- ▶ Indice de Rand

$$\text{Rand} = \frac{a_1 + a_4}{a_1 + a_2 + a_3 + a_4}$$

- ▶ Indice de Jaccard

$$\text{Jaccard} = \frac{a_1}{a_1 + a_2 + a_3}$$

- ▶ Indice de Folkes et Mallows

$$\text{FM} = \sqrt{\frac{a_1}{a_1 + a_3} \frac{a_1}{a_1 + a_2}}$$

Introduction

Classification hiérarchique

Partionnement

Indices de validité

Méthodes probabilistes pour la classification

l'EM-algorithme

- ▶ Soit $\underline{X} = (X_1, \dots, X_n) \sim g(x | \theta)$ un n -échantillon de vecteurs aléatoire dans \mathbb{R}^d . On suppose

$$g(x | \theta) = \int f(x, z | \theta) dz$$

- ▶ Objectif : calculer $\hat{\theta} = \operatorname{argmax}_{\theta} L(\underline{x} | \theta) = \operatorname{argmax}_{\theta} \prod_{i=1}^n g(x_i | \theta)$
où $\underline{x} = (x_1, \dots, x_n)$ est une réalisation de \underline{X} .
- ▶ On pose $(X, Z) \sim f(x, z | \theta)$, et la densité de $Z | \theta, X = x$ est

$$k(z | \theta, x) = \frac{f(x, z | \theta)}{g(x | \theta)}.$$

l'EM-algorithme

- ▶ Donc

$$\log g(x | \theta) = \log f(x, z | \theta) - \log k(z | \theta, x)$$

- ▶ Pour un θ_0 fixé,

$$\begin{aligned} \log g(x | \theta)k(z | \theta_0, x) &= \log f(x, z | \theta)k(z | \theta_0, x) \\ &\quad - \log k(z | \theta, x)k(z | \theta_0, x) \quad (1) \end{aligned}$$

- ▶ En intégrant (1) par rapport à z , on obtient

$$\log g(x | \theta) = \mathbb{E}_{\theta_0}(\log f(x, Z | \theta)) - \mathbb{E}_{\theta_0}(\log k(Z | x, \theta))$$

l'EM-algorithme

- ▶ La fonction vraisemblance s'écrit alors

$$\begin{aligned}\log L(\underline{x} \mid \theta) &= \sum_{i=1}^n \log g(x_i \mid \theta) \\ &= \underbrace{\mathbb{E}_{\theta_0}(\log L^c(\underline{x}, Z \mid \theta))}_{Q(\theta \mid \underline{x}, \theta_0)} - \sum_{i=1}^n \mathbb{E}_{\theta_0}(\log k(Z \mid x_i, \theta))\end{aligned}$$

où L^c est appelée la vraisemblance *complète*.

- ▶ Ainsi, dans la recherche du maximum de $\log L(\underline{x} \mid \theta)$, dans l'EM-algorithme on maximize, d'une façon itérative, $Q(\theta \mid \underline{x}, \theta_0)$.

l'EM-algorithme

1. Considérons une valeur initiale $\hat{\theta}_{(0)}$
2. Calculer l'espérance (*E-step*)

$$Q(\theta \mid \underline{x}, \hat{\theta}_{(m)}) = \mathbb{E}_{\hat{\theta}_{(m)}}(\log L^c(\underline{x}, Z \mid \theta))$$

où la moyenne a été calculée par rapport à $k(z \mid \hat{\theta}_{(m)}, x)$ et $m = 0$

3. Maximiser la fonction $Q(\theta \mid \underline{x}, \hat{\theta}_{(m)})$ en θ (*M-step*) :

$$\hat{\theta}_{(m+1)} = \operatorname{argmax}_{\theta} Q(\theta \mid \underline{x}, \hat{\theta}_{(m)})$$

et $m = m + 1$.

4. Répéter les étapes 2-3 jusqu'à le point fix soit atteint :
 $\hat{\theta}_{(m+1)} = \hat{\theta}_{(m)}$

l'EM-algorithme, dans la classification

- ▶ On suppose que les données $\underline{x} = (x_1, \dots, x_n) \in (\mathbb{R}^d)^n$ est un n -échantillon qu'on suppose issu d'une loi de probabilité

$$g(\underline{x} \mid \theta) = \sum_{l=1}^k \alpha_l g_l(\underline{x} \mid \theta_l)$$

où $\alpha_l \in]0, 1[$, $l = 1, \dots, k$ et $\sum_{l=1}^k \alpha_l = 1$ et $g_l(\underline{x} \mid \theta_l)$ est la densité d'une loi Gaussienne multivariée :

$$f_l(\underline{x} \mid \theta_l) = \frac{1}{(2\pi)^{d/2} |\Sigma_l|^{1/2}} \exp\left(-\frac{1}{2}(\underline{x} - \mu_l)' \Sigma_l^{-1} (\underline{x} - \mu_l)\right)$$

- ▶ Donc $X \sim g(\underline{x} \mid \theta)$ et Z est la variable à valeurs dans $\{1, \dots, k\}$, où k est le nombre de classes et

$$\mathbb{P}(Z = l) = \alpha_l, \quad \forall l = 1, \dots, k.$$

Mise en œuvre sous R : le package `mclust`

- ▶ Description détaillé du package dans
 - ▶ C. Fraley and A. E. Raftery (2006). MCLUST Version 3 for R : Normal Mixture Modeling and Model-Based Clustering, Technical Report no. 504, Department of Statistics, University of Washington
 - ▶ C. Fraley and A. E. Raftery (2002). Model-based clustering, discriminant analysis, and density estimation. Journal of the American Statistical Association 97 :611-631
- ▶ Utiliser la commande `Mclust`,
- ▶ `Mclust` applique l'EM-algorithme (avec des différentes paramétrisations), et utilise le BIC comme critère de sélection du modèle : chercher le modèle M_k correspondant au BIC maximal.

$$\text{BIC} = 2 \log \text{Vraisemblance} - \text{Nbre paramètres} \times \log(n)$$