

Classification de données longitudinales

Ndèye Niang
Conservatoire National des Arts
et Métiers

STA112 Niang

Plan

- Introduction
- Rappels sur la classification
- Distances en analyse de données évolutives
- Typologies de trajectoires

(Référence : analyse de données évolutives méthodes et applications Dazy, Le Barzic Saporta, Lavallard Technip 1996)

STA112 Niang

Introduction

- Données longitudinales **multidimensionnelles**
- **Plusieurs individus** - plusieurs variables
- Méthodes « factorielles » permettent visualisation des trajectoires des individus sur les plans factoriels mais interprétation difficile si grand nombre de trajectoires
- Objectif classification = **réduction du nombre de trajectoires**

STA112 Niang

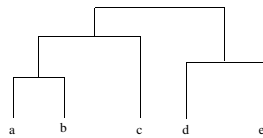
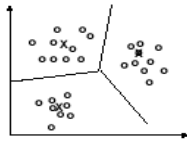
Introduction

- Démarche:
 - Application d'une méthode factorielle (STATIS,...)
 - Coordonnées factorielles des individus sur les q axes retenus
 - Application d'une méthode de classification en ayant au préalable défini une distance
 - Les trajectoires des individus de chaque classe sont représentées par la trajectoire du centre de gravité de la classe (réduction)

STA112 Niang

Rappels sur la classification

- **Méthodes de partitionnement :**
 - **une partition en un nombre fixe de classes**
- **Méthodes hiérarchiques :**
 - **suite de partitions emboîtées**



STA112 Niang

Rappels sur la classification

- **L'objectif des méthodes de classification automatique est la construction d'une partition ou d'une suite de partitions emboîtées d'un ensemble d'objets.**
- **Les classes formées doivent être le plus homogènes possible d'où la nécessité de définir un critère à optimiser.**

STA112 Niang

Rappels sur la classification

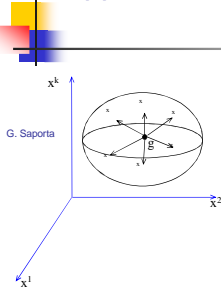
- Inertie = « mesure » de la dispersion du nuage de points de \mathbb{R}^p espace euclidien

I_g par rapport à g centre de gravité

- Plus I est faible plus la classe est homogène

STA112 Niang

Rappels sur la classification



$I =$ moyenne des carrés des distances à g

$$\sum p_i d^2(i, g)$$

Inertie = variance généralisée = $\sum Var(x_j)$

STA112 Niang

$$I_W = \sum P_i I_i \quad \text{inertie intra classe}$$

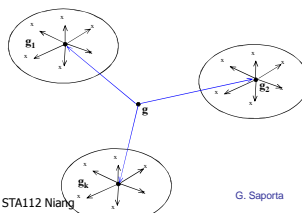
$$I_B = \sum P_i d^2(g_i; g) \quad \text{inertie inter classe}$$

Relation de Huyghens : $I = I_W + I_B$

partition en k classes

g_1, \dots, g_k les centres de gravité des classes

I_1, \dots, I_k les inerties des k classes



STA112 Niang

G. Saporta

Rappels sur la classification

- Critère usuel de classification = chercher la partition qui minimise I_W cad qui maximise I_B
- I_W ne permet pas de comparer 2 partitions avec des nombres de classes différents.

STA112 Niang

■ Méthode des centres mobiles (Forgy) K-means

Etape 1

- a) configuration initiale : $C_1^{(0)}, C_2^{(0)}, \dots, C_k^{(0)}$
- b) chaque individu i est affecté à une classe et une seule $E_i^{(0)}$ de centre $C_i^{(0)}$ telle que : soit minimum en parcourant tous les centres $C_1^{(0)}, C_2^{(0)}, \dots, C_k^{(0)}$
à la fin de cette étape on a k classes $E_1^{(0)}, E_2^{(0)}, \dots, E_k^{(0)}$

Etape 2

- a) centres de gravité des classes précédentes : $C_1^{(1)}, C_2^{(1)}, \dots, C_k^{(1)}$
- b) chaque individu i est affecté à une classe et une seule $E_i^{(1)}$ de centre $C_i^{(1)}$ telle que : soit minimum en parcourant tous les centres $C_1^{(1)}, C_2^{(1)}, \dots, C_k^{(1)}$
à la fin de cette étape on a k classes $E_1^{(1)}, E_2^{(1)}, \dots, E_k^{(1)}$

arrêt de la procédure :

- 2 étapes successives ne changent pas les classes
- le nombre d'itérations fixé est atteint
- la valeur du critère reste inchangée

STA112 Niang

Rappels sur la classification

- Remarque : formes fortes
- Problème : la partition finale dépend du nombre de classes et du choix des centres initiaux
- Solution : appliquer l'algorithme sur s tirages différents, croiser les s partitions pour obtenir une partition dite en formes fortes ou regroupements stables
- (formes fortes = ensembles d'éléments ayant toujours été regroupés dans la partition finale pour les s passages de l'algorithme).

STA112 Niang

Rappels sur la classification

Méthodes hiérarchiques

- Elles consistent à fournir un ensemble de partitions de E en classes de moins en moins fines par regroupements successifs de parties.
- On obtient une hiérarchie représentée par un arbre de classification ou dendrogramme.
- On associe au système de classes résultant une échelle de niveau : à chaque partition on associe une valeur numérique représentant le niveau auquel ont lieu les regroupements
- Différentes méthodes selon la stratégie de regroupement

STA112 Niang

Rappels sur la classification

Méthodes hiérarchiques

Comment mesurer la distance entre une partie et un élément, entre deux parties?

Stratégie d'agrégation:

- ✓ saut minimal (single linkage): $d(A,B) = \inf(a,b)$
- ✓ diamètre (complete linkage): $d(A,B) = \sup(a,b)$
- ✓ Moyenne (average linkage)
- ✓ Ward $d_w(A,B) = (p_A p_B / (p_A + p_B)) d^2(g_A, g_B)$

STA112 Niang

Rappels sur la classification

Classification mixte

Les algorithmes classiques sont plus ou moins adaptés à la gestion d'un nombre importants d'objets à classer :

- partitionnement : ensemble volumineux à faible coût mais la partition dépend des centres initiaux et du nombre de classes.
- hiérarchique : non adaptée aux vastes ensembles

D'où les algorithmes mixtes

- * centres mobiles
- * classification hiérarchique des groupes obtenus
- * réaffectation par centres mobiles

STA112 Niang

Distances en analyse de données évolutives

- Trois types de distances
 - Distances euclidiennes basées sur les positions
 - Distances euclidiennes basées sur les évolutions des individus
 - Distances compromis positions-évolutions

STA112 Niang

Distances en analyse de données évolutives

- Données et notations
- Rappel : les coordonnées des individus sont obtenus à l'étape intrastructure par projection des T tableaux sur les axes du compromis
- Tableau de données de la classification = tableau des coordonnées factorielles: q axes retenus * T instants
- Pour chaque individu q*T coordonnées

STA112 Niang

Distances en analyse de données évolutives

- Voir au tableau

STA112 Niang

Distances euclidiennes basées sur les positions

- Distance 1: euclidienne classique

$$d_1(i, i')^2 = \sum_{t=1}^T \sum_{j=1}^q [(c_i^j)^t - (c_{i'}^j)^t]^2$$

- $d_1(i, i') = 0$ Ssi égalité de toutes les $c(i, j, t)$
- On aura une classification selon les positions
- Pas de prise en compte de l'aspect temps

STA112 Niang

Distances euclidiennes basées sur les évolutions des individus

- On définit l'évolution des individus par:

$$(exd_i^j)^t = (c_i^j)^t - (c_i^j)^{t-1}$$

- Et la distance évolution est :

$$d_2(i, i')^2 = \sum_{t=2}^T \sum_{j=1}^q [(exd_i^j)^t - (exd_{i'}^j)^t]^2$$

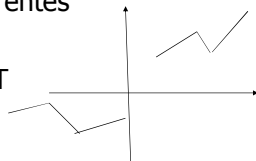
- $d_2(i, i') = 0$ Ssi les vecteurs de coordonnées sont liés par une translation

STA112 Niang

Distances euclidiennes basées sur les évolutions des individus

- Plus adaptée aux données temporelles, longitudinales.

- Quelques inconvénients cependant:
- Classe identique car même évolution mais caractéristiques différentes
- Classification de T-1 trajectoires au lieu de T



STA112 Niang

Distances compromis positions-évolutions

Distance compromis:

- A. Carlier propose

$$d_3(i, i')^2 = \beta_1 d_1(i, i')^2 + \beta_2 d_2(i, i')^2$$

- Pondération pour « homogénéiser » d1 et d2, choix non explicité
- 2T coordonnées alors que T suffisent

STA112 Niang

Distances compromis positions-évolutions

Dazy et al proposent:

$$d_4(i, i')^2 = \sum_{j=1}^q ((c_i^j)^1 - (c_{i'}^j)^1)^2 + d_2(i, i')^2$$

- Premier terme de d1 + d2
 - Point de départ = position à t=1, et T-1 évolutions
- Déséquilibre trop de terme dans d2
- Autre approche:

STA112 Niang

Distances compromis positions-évolutions

- Principe de l'approche proposée par Dazy et al
- Déterminer les instants t parmi T pour lesquels on utilise la coordonnée-position et ceux pour lesquels on utilise la coordonnée-évolution
- On se ramène à un problème de minimisation d'une fonction différence entre inertie calculée avec les coordonnées-position et inertie calculée avec les coordonnées-évolution

STA112 Niang

Application aux données de croissance des jeunes filles

- Après STATIS, on récupère les coordonnées intrastructure des points individus
- On utilise SPAD pour faire la classification basée sur les distances position
- Voir listing document joint.

STA112 Niang

Remarques, autres cas

- Cas général de la classification de trajectoires: trajectoires de longueur différentes, non synchronisées, ...
- Mesures de distances spécifiques: par ex:
 - LCSS (Longest Common Subsequence)
 - DTW (Dynamic Time Warping)
 - Références Vlachos 2003, Aris Anagnostopoulos et Vlachos 2006

STA112 Niang
