

# ANR CLASEL

**Mohamed Nadif<sup>1</sup>, Gérard Govaert<sup>2</sup>**

<sup>1</sup>LIPADE, Université Paris Descartes, France

<sup>2</sup>Heudiasyc, Université de Technologie de Compiègne, France

## Introduction

- Types d'algorithmes de classification croisée
- Intérêts
- Cas des données continues
- Approche modèle de mélange

## Livrables

- Modèles pour la classification croisée
  - Livrable 1.1 : Etat de l'art
  - Livrable 1.2 : Modèles pour les données continues
- Sélection de modèle
  - Livrable 2.3 : Aspects asymptotiques

## Travaux en cours et programme

## Données

- matrix  $\mathbf{x} = (x_{ij})$
- $i \in I$  ensemble de  $n$  lignes
- $j \in J$  ensemble de  $d$  colonnes

## Partition $\mathbf{z}$ de $I$ en $g$ classes

- $\mathbf{z} = (z_{ik})$  : une matrice de classification de taille  $(n \times g)$
- $z_{ik} = 1$  si  $i \in k$ ème classe and  $z_{ik} = 0$  sinon

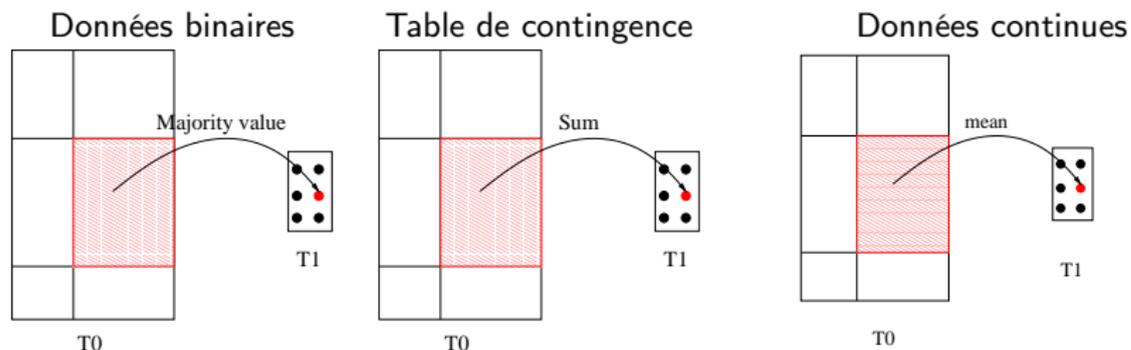
$$\begin{array}{c|ccc} 3 & 0 & 0 & 1 \\ 2 & 0 & 1 & 0 \\ 3 & 0 & 0 & 1 \\ 2 & 0 & 1 & 0 \\ 1 & 1 & 0 & 0 \end{array}$$

## Partition $\mathbf{w}$ de $J$ en $m$ classes

- $\mathbf{w} = (w_{j\ell})$  : une matrice de classification de taille  $(d \times m)$
- $w_{j\ell} = 1$  si  $j \in \ell$ ème classe et  $w_{j\ell} = 0$  sinon

Bloc  $k\ell$  défini par  $i, j$  tel que  $z_{ik}w_{j\ell} = 1$

## Principe général



## Critères

| Données              | $\mu_{kl}$ | Critère $W(\mathbf{z}, \mathbf{w}, \mu)$  |
|----------------------|------------|---|
| Binaires             | Mode       | $\sum_{i,j,k,l} z_{ik} w_{jl}  x_{ij} - \mu_{kl} $  |
| Table de contingence | Somme      | $\chi^2(\mathbf{z}, \mathbf{w}) = K \sum_{k,l} \frac{(f_{kl} - f_k \cdot f_l)^2}{f_k \cdot f_l}$    |
| Continues            | Moyenne    | $\sum_{i,j,k,l} z_{ik} w_{jl} (x_{ij} - \mu_{kl})^2 = \ \mathbf{x} - \mathbf{z}\mu\mathbf{w}^T\ ^2$ |

## Complémentaires aux méthodes factorielles

- ACP, AC, etc.

## Réduction de la taille des données

- Offre un résumé (de même nature) de la table initiale

## Méthodes capables de gérer des données de grande taille

- Moins de calcul par l'utilisation de matrices de taille réduite
- Remédie à la grande dimension et la *sparsité*

## Applications

- *Classification de documents* : classification des documents et des mots simultanément
- *Données des biopuces* : classification des gènes et des tissus simultanément
- *Systèmes de recommandation*

Minimisation du critère  $W(\mathbf{z}, \mathbf{w}, \boldsymbol{\mu}) = \|\mathbf{x} - \mathbf{z}\boldsymbol{\mu}\mathbf{w}^T\|^2$

## Double $k$ -means

- (a) Recherche de  $\boldsymbol{\mu} = (\mu_{kl})$  et  $\mathbf{z}$  en minimisant  $W(\mathbf{z}, \boldsymbol{\mu}|\mathbf{w}) = \sum_k \sum_{j,\ell} w_{j\ell} (x_{ij} - \mu_{kl})^2$   
 (b) Recherche de  $\boldsymbol{\mu} = (\mu_{kl})$  et  $\mathbf{w}$  en minimisant  $W(\mathbf{w}, \boldsymbol{\mu}|\mathbf{z}) = \sum_\ell \sum_{i,k} z_{ik} (x_{ij} - \mu_{kl})^2$

## Algorithme *Croec* : Govaert (1983)

Working on intermediate  $(n \times m)$  matrix  $\mathbf{u} = (u_{il})$  et  $(g \times d)$  matrix  $\mathbf{v} = (v_{kj})$

- (a) minimisation de  $W(\mathbf{z}, \boldsymbol{\mu}|\mathbf{w}, \mathbf{u})$  où  $u_{il} = \sum_j w_{jl} x_{ij} / w_\ell$  ( $w_\ell = \sum_j w_{j\ell}$ )  
 (a.1)  $k$ -means sur  $\mathbf{u}$  et on obtient  $\mathbf{z}$  et  $\boldsymbol{\mu}$   
 (b) minimisation de  $W(\mathbf{w}, \boldsymbol{\mu}|\mathbf{z}, \mathbf{v})$  où  $v_{kj} = \sum_i z_{ik} x_{ij} / z_k$  ( $z_k = \sum_i z_{ik}$ )  
 (b.1)  $k$ -means sur  $\mathbf{v}$  et on obtient  $\mathbf{w}$  et  $\boldsymbol{\mu}$

## *Croec* vs Double $k$ -means et autres variantes

- Résultats équivalents
- *Croec* est plus rapide

Plusieurs problèmes avec ce critère (prop. varia.) : approche *Latent Block Model*

## Latent block model, Govaert and Nadif (2003)

### Définition du modèle de mélange sur $I \times J$

$$f(\mathbf{x}, \theta) = \sum_{\mathbf{u} \in U} P(\mathbf{u}) f(\mathbf{x} | \mathbf{u}; \alpha)$$

où  $U$  est l'ensemble de toutes les partitions de  $I \times J$

### Hypothèses supplémentaires pour considérer la classification croisée

- $\mathbf{u} = \mathbf{z} \times \mathbf{w}$
- $f(\mathbf{x} | \mathbf{z}, \mathbf{w}; \alpha) = \prod_{i,j} \varphi(x_{ij}; \alpha_{z_i, w_j})$  où  $\varphi(\cdot, \alpha)$  est pdf sur  $\mathbb{R}$

### Latent block model

$$f(\mathbf{x}; \theta) = \sum_{(\mathbf{z}, \mathbf{w}) \in \mathcal{Z} \times \mathcal{W}} \underbrace{\prod_{i,k} \pi_k^{z_{ik}}}_{prop.} \underbrace{\prod_{j,\ell} \rho_\ell^{w_{j\ell}}}_{prop.} \underbrace{\prod_{i,j,k,\ell} \varphi(x_{ij}; \alpha_{k\ell})^{z_{ik} w_{j\ell}}}_{dens.}$$

où  $\theta = (\pi_1, \dots, \pi_g, \rho_1, \dots, \rho_m, \alpha_{11}, \dots, \alpha_{gm})$

# Types de données : différentes approches, ML, CML et floue

*Bernoulli latent block model*, Govaert and Nadif (2003, 2005a, 2008)

- Données binaires
- $\varphi$  distribution de Bernoulli  $\mathcal{B}(\alpha_{kl})$

*Poisson latent block model*, Govaert and Nadif (2005b)

- Table de contingence
- $\varphi$  distribution de Poisson  $\mathcal{P}(\mu_i \nu_j \alpha_{kl})$ 
  - $\mu_i$  et  $\nu_j$  effets de  $i$  et de  $j$
  - $\alpha_{kl}$  effet du bloc  $kl$ .

## Avantages

- Modèles plus parcimonieux qu'un modèle de mélange classique sur  $I$  et sur  $J$
- Critères plus riches et algorithmes de type *Block EM* plus efficaces pour la classification croisée

## A faire

- Extension au cas continu possible
- Sélection de modèle

## Historique

- Travaux depuis 1970

## Difficultés

- Approches et objectifs différents

## Domaines

- Statistique
- Analyse de données textuelles
- Filtrage collaboratif
- Factorisation de matrice non-négative
- Bioinformatique
- Approche modèle de mélange

## Résultat

- Un rapport avec 124 références
- Une soumission dans une revue prévue avant fin décembre



## Caractéristiques du modèle

- Nombre de paramètres égal  $g + (g \times m) + (g \times m)$  au lieu de  $g + (g \times d) + (g \times d)$ .
- Parcimonie très utile lorsque  $n < d$

## Algorithme GEM

- E-step : Calcul de  $s_{ik}^{(c)} = P(z_{ik} = 1 | \mathbf{x}, \boldsymbol{\theta}^{(c)}) \propto \pi_k^{(c)} \varphi_k(\mathbf{x}_i; \mathbf{w}^{(c)}, \boldsymbol{\alpha}^{(c)})$
- M-step : Maximisation de  $Q(\boldsymbol{\theta}, \boldsymbol{\theta}^{(c)}) = \sum_{i,k} s_{ik}^{(c)} \log(\pi_k \varphi_k(\mathbf{x}_i; \mathbf{w}, \boldsymbol{\alpha}))$

## Algorithme CGEM

- Lorsque toutes les variances sont supposées égales 1) la maximisation de la vraisemblance classifiante est équivalente à la minimisation de  $\|\mathbf{x} - \mathbf{z}\boldsymbol{\mu}\mathbf{w}^T\|^2$  2) La version CGEM est équivalente à *Croeu*c

## Résultats

- Une publication dans CAP'2009
- Une publication dans ICMLA'2010
- Préparation d'une version étendue

## Sélection de modèle

- Extension du critère BIC :  $BIC(g, m) = L - d \log(m) - \frac{\nu}{2} \log(nd)$
- CAP'2009

## Thèse, Heudiasyc

- Expression du critère de BIC au modèle *Latent block model* par une approximation variationnelle étendue
- Sujet de thèse de Aurore Lomet (financée par le ministère)

## Post-doc, Heudiasyc

- Principe du bootstrap pour corriger l'erreur de l'estimateur de la vraisemblance ou la vraisemblance classifiante
- Modèles de mélange Gaussiens simples
- Plusieurs critères bootstrap expérimentés

# Travaux dans les 6 mois à venir

## Heudiasyc & LIPADE

- Plusieurs travaux en cours de finalisation
- Etude des liens entre le modèle proposé et la tri-factorisation de matrice non négative (post-doc, LIPADE, 01/09/10)
- Visualisation des données par GTM (en cours de soumission)

## Heudiasyc

- Version Stochastique de *Block EM*
- Sélection de modèle

## LIPADE

- Données manquantes : Algorithmes qui fonctionnent
- Bagging dans la classification croisée (ECML/PKDD 2010)