

MIXMOD

A software for model-based classification for
quantitative and qualitative data

C.Biernacki¹ G.Celeux² G.Govaert³ F.Langrognet⁴

¹UMR824 - Univ. Lille & CNRS - Lille (France)

²INRIA Orsay (France)

³UMR6599 - UTC & CNRS - Compiègne (France)

⁴UMR6623- Univ. Franche-Comté & CNRS - Besançon (France)

June 2008

MIXMOD

A software for model-based classification for
quantitative and qualitative data

- 1 MIXMOD functionalities
- 2 Software overview
- 3 Illustration : supervised classification with qualitative data
- 4 Perspectives

MIXMOD

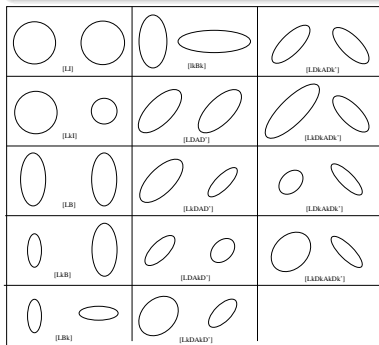
- 1 MIXMOD functionalities
- 2 Software overview
- 3 Illustration : supervised classification with qualitative data
- 4 Perspectives

Parsimonious models

Gaussian models

14 Gaussian models

based on eigenvalue decomposition
of variance matrices
(see Celeux and Govaert 1995)



High Dimensional Data

8 specific models for High Dimensional Data
(see Bouveyron 2007)

Multinomial models

5 Multinomial models

based on a parameterization of Bernoulli
distributions
(see Celeux and Govaert 1991)

Other functionalities

Algorithms

To maximize Likelihood or Completed Likelihood

- **EM** (Expectation Maximisation)
- **SEM** (Stochastic EM)
- **CEM** (Classification EM)

Criteria

- **BIC** (Bayesian Information Criterion)
- **ICL** (Integrated Completed Likelihood)
- **NEC** (Normalized Entropy Criterion)
- **CV** (Cross Validation)

Initializations and Strategies

- **6 initializations**
Ex : 'random', 'short runs of EM'
- **Chained algorithms**
Ex : 100 iterations of **SEM** and 50 iterations of **EM**

Others

- **Partial labeling** of individuals
- **Weighted** individuals

MIXMOD

- 1 MIXMOD functionalities
- 2 Software overview**
- 3 Illustration : supervised classification with qualitative data
- 4 Perspectives

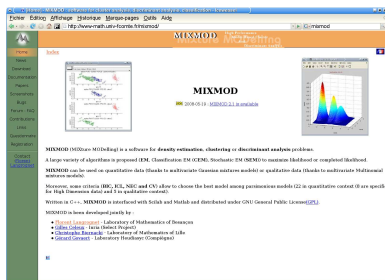
Software overview

Historical background

- Sep. 2000 : Mixmod 1.0
- Feb. 2007 : Mixmod 2.0 (with qualitative data processing)
- May 2008 : **Mixmod 2.1** (with models for High Dimension data)

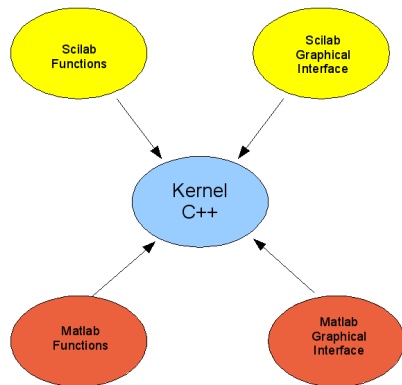
Distribution

- <http://www-math.univ-fcomte.fr/mixmod/>
800 visits and 300 downloads per month
- GNU **GPL** License
- **Windows and Linux** packages (source and binary)
- **Documentations** (statistical, userguide, quickStart...)



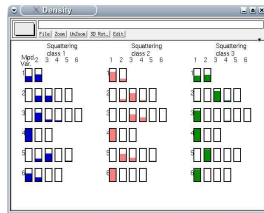
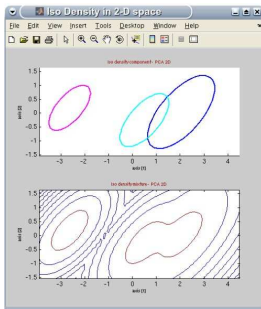
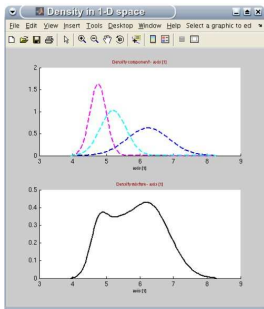
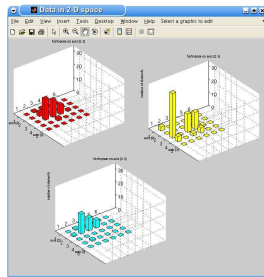
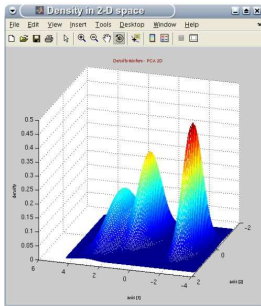
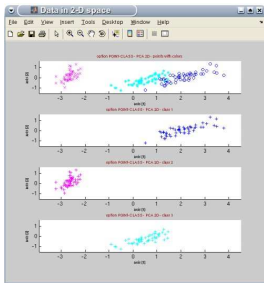
The screenshot shows the homepage of the Mixmod software project. The page features a navigation menu on the left, a central banner with the Mixmod logo and version information (2008-05-19 - MIXMOD 2.1 is available), and a main content area with a 3D plot and descriptive text. The text describes Mixmod as a software for density estimation, clustering, and discriminant analysis, and lists the authors: Olivier Lacombe, Sébastien Ballester, Christophe Berruyer, and Sébastien Lacroix.

Software Architecture



- **Mixmod Kernel** : C++
- **Two levels of use** :
 - graphical interfaces (Scilab, Matlab)
 - functions for Scilab and Matlab
- **Mixmod library**

Screenshots



MIXMOD

- 1 MIXMOD functionalities
- 2 Software overview
- 3 Illustration : supervised classification with qualitative data**
- 4 Perspectives

Supervised classification with qualitative data

Puffins



Data

- Number of **individuals** : $n = 69$
- Number of **species** (or clusters) : $K = 2$
- Number of **variables** : $d = 6$
- **Individual i** :
 $(x_i, z_i) = ((x_i^j)_{j=1, \dots, d}, z_i)$

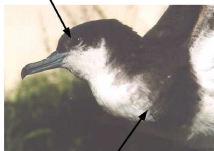
n°	z_i	x_i^1	x_i^2	x_i^3	x_i^4	x_i^5	x_i^6
1	1	1	2	2	1	2	2
2	1	2	1	3	1	3	1
\vdots	\vdots			\vdots			
68	2	1	4	1	1	2	1
69	2	1	3	1	1	2	1

Observations

Variables description

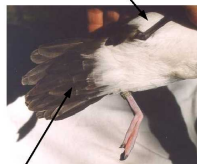
variable	number of response levels	values
sex	2	2 values (male, female)
eyebrows	5	5 values (none -> heavy)
collar	6	6 values (none -> unless)
strips	3	3 values (none, ...)
subcaudal	5	5 values (black, white, ...)
liseret	4	4 values (...)

sourcil



collier (sur la poitrine blanche)

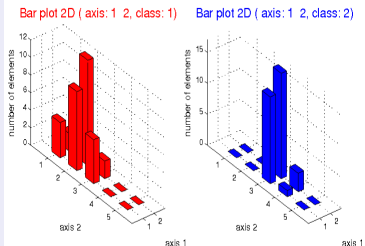
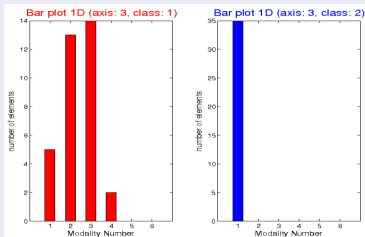
zébrures (sur les flancs)



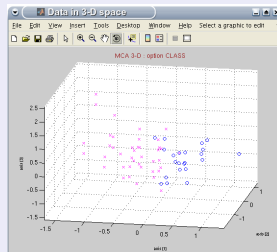
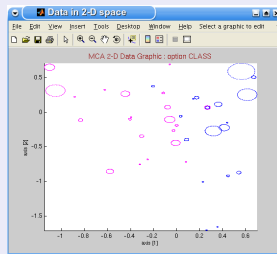
sous-caudales (plumes sous la queue)

Data visualisation

barplots



scatterplots

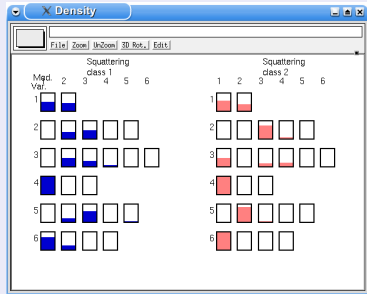


Discriminant Analysis with MIXMOD (step 1)

Classification rule

Estimated parameters

$$\hat{\rho}_1 = 0.49, \quad \hat{\rho}_2 = 0.51$$
$$\hat{\alpha}_1 = \begin{bmatrix} \hat{\alpha}_1^1 = 0.53 & 0.47 & & & & & \\ \hat{\alpha}_2^1 = 0.00 & 0.44 & 0.56 & 0.00 & 0.00 & & \\ \hat{\alpha}_3^1 = 0.00 & 0.62 & 0.30 & 0.08 & 0.00 & 0.00 & \\ \hat{\alpha}_4^1 = 1.00 & 0.00 & 0.00 & & & & \\ \hat{\alpha}_5^1 = 0.00 & 0.20 & 0.72 & 0.00 & 0.08 & & \\ \hat{\alpha}_6^1 = 0.81 & 0.19 & 0.00 & 0.00 & & & \end{bmatrix}$$
$$\hat{\alpha}_2 = \begin{bmatrix} \hat{\alpha}_1^2 = 0.59 & 0.41 & & & & & \\ \hat{\alpha}_2^2 = 0.00 & 0.00 & 0.89 & 0.11 & 0.00 & & \\ \hat{\alpha}_3^2 = 0.69 & 0.00 & 0.15 & 0.16 & 0.00 & 0.00 & \\ \hat{\alpha}_4^2 = 1.00 & 0.00 & 0.00 & & & & \\ \hat{\alpha}_5^2 = 0.00 & 0.97 & 0.03 & 0.00 & 0.00 & & \\ \hat{\alpha}_6^2 = 1.0 & 0.00 & 0.00 & 0.00 & 0.00 & & \end{bmatrix}$$



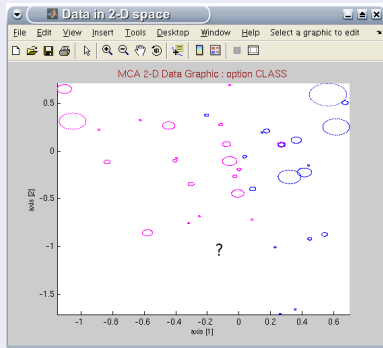
Quality of the classification rule

- Cross Validation rate (leave-one-out) : 94.8%
- Reclassification by MAP method : 96.7%

Discriminant Analysis with MIXMOD (step 2)

Classification of a new observation

n°	z_j	x_j^1	x_j^2	x_j^3	x_j^4	x_j^5	x_j^6
70	?	1	3	1	1	5	2



Results

$$z_{70} = 2$$

$$p(z_{70} = 1) = 0.06 \text{ and } p(z_{70} = 2) = 0.94$$

MIXMOD

- 1 MIXMOD functionalities
- 2 Software overview
- 3 Illustration : supervised classification with qualitative data
- 4 Perspectives

Perspectives

New functionalities

- **Semi-supervised** classification
- **Heterogeneous data** (quantitative and qualitative data)
- ...

Software evolutions

- **New graphical user interface** (without Scilab and Matlab)
- **Double license** distribution (GPL and ...)
- ...

Thank you for your attention

Web : <http://www-math.univ-fcomte.fr/mixmod/>

Authors :

- christophe.biernacki@math.univ-lille1.fr
- gilles.celeux@math.u-psud.fr
- gerard.govaert@hds.utc.fr
- florent.langrognet@univ-fcomte.fr

ANNEX

Finite mixture models Latent Class Model

Finite mixture models

- $\mathbf{x} = (\mathbf{x}_1, \dots, \mathbf{x}_n)$: n independent vectors in \mathbf{R}^d .
Each \mathbf{x}_i arises from a probability distribution with density : $f(\mathbf{x}) = \sum_{k=1}^K p_k f_k(\mathbf{x})$
 - ▶ K : number of groups
 - ▶ p_k : mixing proportions
 - ▶ $f_k(\cdot)$: densities of components
 - ★ Gaussian densities for quantitative data
 - ★ Multinomial densities for qualitative data
 - $\mathbf{z} = \{\mathbf{z}_1, \dots, \mathbf{z}_n\}$: a sample of indicator vectors or labels
-
- Unsupervised classification : \mathbf{z} is unknown
 - Supervised classification : \mathbf{z} is known
Estimate \mathbf{z}_{n+1} , the label of a new observation \mathbf{x}_{n+1}

Qualitative data : Latent Class Model (LCM)

Data

$$\mathbf{x}_i = (x_i^{jh}; j = 1, \dots, d; h = 1, \dots, m_j)$$

- m_j : the **number of response levels** of the variable j

- $$\begin{cases} x_i^{jh} = 1 & \text{if } i \text{ has level } h \text{ for variable } j \\ x_i^{jh} = 0 & \text{otherwise.} \end{cases}$$

Example : $X_i \begin{cases} x_i^1 = 10 \\ x_i^2 = 0010 \\ x_i^3 = 010 \end{cases} \quad (d = 3, m_1 = 2, m_2 = 4, m_3 = 3)$

Multivariate multinomial distribution

$$f(\mathbf{x}_i; \theta) = \sum_k p_k m_k(\mathbf{x}_i; \alpha_k) = \sum_k p_k \prod_{j,h} (\alpha_k^{jh})^{x_i^{jh}}$$

$\theta = (p_1, \dots, p_K, \alpha_1^{11}, \dots, \alpha_K^{dm_d})$: the parameter of the latent class model

- α_k^{jh} : probability that variable j has level h in cluster k
- p_k : mixing proportions

LMC : Reparameterization

$$\forall k, j : (\alpha_k^{j1}, \dots, \alpha_k^{jm_j}) \longrightarrow (\mathbf{a}_k^{j1}, \dots, \mathbf{a}_k^{jm_j}, \varepsilon_k^{j1}, \dots, \varepsilon_k^{jm_j})$$

$$\bullet \mathbf{a}_k^{jh} = \begin{cases} 1 & \text{if } h = \arg \max_h \alpha_k^{jh} \\ 0 & \text{otherwise} \end{cases} \quad \textit{center}$$

$$\bullet \varepsilon_k^{jh} = \begin{cases} 1 - \alpha_k^{jh} & \text{if } \mathbf{a}_k^{jh} = 1 \\ \alpha_k^{jh} & \text{if } \mathbf{a}_k^{jh} = 0. \end{cases} \quad \textit{scattering}$$

Example : (0.2, 0.7, 0.1) \longrightarrow (0, 1, 0, 0.2, 0.3, 0.1).

Five latent class models

Different **constraints** to the **scattering** parameters ε_k^{jh} .

model	the scattering is depending on...
$[\varepsilon_k^{jh}]$	clusters, variables and levels
$[\varepsilon_k^j]$	clusters and variables
$[\varepsilon_k]$	clusters
$[\varepsilon^j]$	variables
$[\varepsilon]$	none