



HAL
open science

Simultaneous dimension reduction and multi-objective clustering using probabilistic factorial discriminant analysis

Vincent Vandewalle

► **To cite this version:**

Vincent Vandewalle. Simultaneous dimension reduction and multi-objective clustering using probabilistic factorial discriminant analysis. CMStatistics 2016, Dec 2016, Sevilla, Spain. hal-01424965

HAL Id: hal-01424965

<https://inria.hal.science/hal-01424965>

Submitted on 3 Jan 2017

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

Simultaneous dimension reduction and multi-objective clustering using probabilistic factorial discriminant analysis

V. Vandewalle

University Lille 2, EA 2694 & Inria Lille, Modal team

ERCIM 2016
Sevilla
December 11th, 2016

Dimension reduction and clustering

Summarizing the data

- **Dimension reduction:** find some principal components explaining the major *variability* in the data
- **Clustering:** find some clusters explaining the major *heterogeneity* of the data

Link between dimension reduction and clustering

Often when visualizing the data one is interested in *visualizing clusters* on the *visualization space*.

In practice

Dimension reduction and clustering are performed separately or sequentially but rarely simultaneously.

Combining dimension reduction and clustering through probabilistic models

Generative models allow to combine visualization and clustering

- Tibshirani & Hastie (1996): Reduced rank discriminant analysis
- Kumar & Andreou (1998): Heteroscedastic discriminant analysis
- Bouveyron & Brunet (2012): Model-based clustering and visualization in the Fisher discriminative subspace

In a rigorous probabilistic way

- A unique homogeneous criterion to optimize: the likelihood
- Simultaneous selection of the number clusters and of the number of components: BIC
- Missing data naturally taken into account: EM algorithm

Multi-objective clustering

Motivation: clustering part

- Usually clustering summarizes the data information by only **one** latent variable, the clustering variable.
- But we would like to allow for **several** views of the data with potentially more than one clustering variable.
- Often performances of clustering methods measured based on classification of reference, but many possible references in many settings (sex, species, status, ...).

Motivation: visualisation part

- View each clustering variable on a **clustering component**
- Heterogeneity-based visualisation rather than inertia-based visualisation

Illustration (1/2)

Data in the initial space

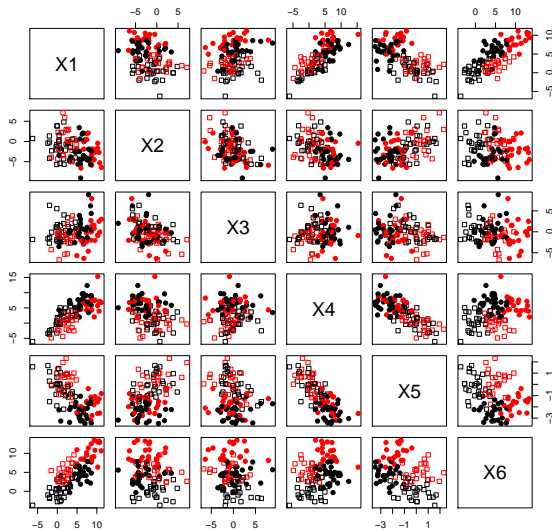
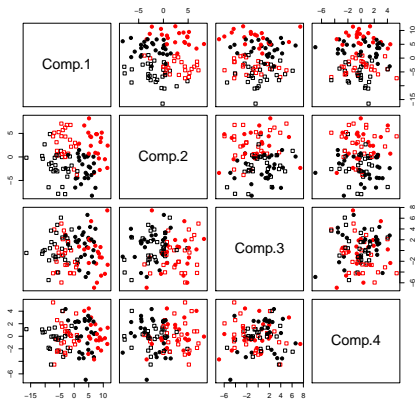
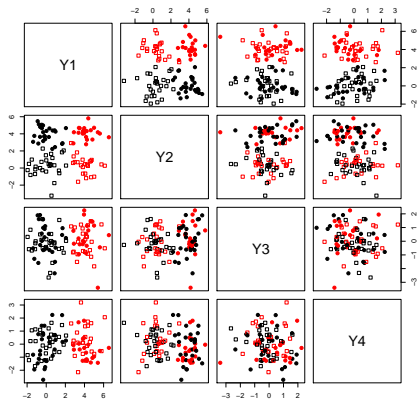


Illustration (2/2)

Principal Components Analysis



Principal Cluster Components



Outline

1. Probabilistic interpretation of the factorial discriminant analysis
2. Multi-objective mixture model
3. Estimation of the parameters and model choice
4. Experiments on simulated and real datasets

Linear discriminant analysis (LDA)

Data

- $\mathbf{x} \in \mathbb{R}^d$: the observed continuous variables
- $\mathbf{z} \in \{1, \dots, K\}$: the cluster (latter considered as unknown)

Assumptions

n independent replicates of (\mathbf{x}, \mathbf{z})

- $\mathbf{z} \sim \text{Mult}(\pi_1, \dots, \pi_K)$ with $\pi_k = p(z_k = 1)$
- $\forall k \in \{1, \dots, K\}, \mathbf{x} | z_k = 1 \sim \mathcal{N}_d(\boldsymbol{\mu}_k, \boldsymbol{\Sigma})$

Parameters estimation

$\hat{\pi}_k = \frac{n_k}{n}$, $\hat{\boldsymbol{\mu}}_k = \bar{\mathbf{x}}_k$ and $\hat{\boldsymbol{\Sigma}} = \mathbf{W}$ (empirical within-class covariance matrix)

Classification

$$p(z_k = 1 | \mathbf{x}, \hat{\theta}) = \frac{\hat{\pi}_k \phi_d(\mathbf{x}; \hat{\boldsymbol{\mu}}_k, \hat{\boldsymbol{\Sigma}})}{\sum_{k'=1}^K \hat{\pi}_{k'} \phi_d(\mathbf{x}; \hat{\boldsymbol{\mu}}_{k'}, \hat{\boldsymbol{\Sigma}})} \Rightarrow \hat{\mathbf{z}} = \operatorname{argmax}_{k \in \{1, \dots, K\}} \hat{\boldsymbol{\mu}}_k^\top \hat{\boldsymbol{\Sigma}}^{-1} \mathbf{x} - \frac{1}{2} \boldsymbol{\mu}_k^\top \hat{\boldsymbol{\Sigma}}^{-1} \boldsymbol{\mu}_k - \log(\hat{\pi}_k)$$

\Rightarrow linear classification boundary.

Factorial discriminant analysis (FDA)

FDA principle

- The first component of FDA finds $\mathbf{v}_1 \in \mathbb{R}^d$ maximizing the variance explained by the cluster:

$$\mathbf{v}_1 = \arg \max_{\mathbf{v} \in \mathbb{R}^d} \frac{\mathbf{v}^T \mathbf{B} \mathbf{v}}{\mathbf{v}^T (\mathbf{W} + \mathbf{B}) \mathbf{v}} = \arg \max_{\mathbf{v} \in \mathbb{R}^d} \frac{\mathbf{v}^T \mathbf{B} \mathbf{v}}{\mathbf{v}^T \mathbf{W} \mathbf{v}}$$

with \mathbf{B} the empirical between class covariance matrix.

- \mathbf{v}_1 is given by the eigen vector associated with the highest eigen value of $\mathbf{W}^{-1} \mathbf{B}$
- The remaining discriminant components are obtained through the remaining eigen vectors of $\mathbf{W}^{-1} \mathbf{B}$

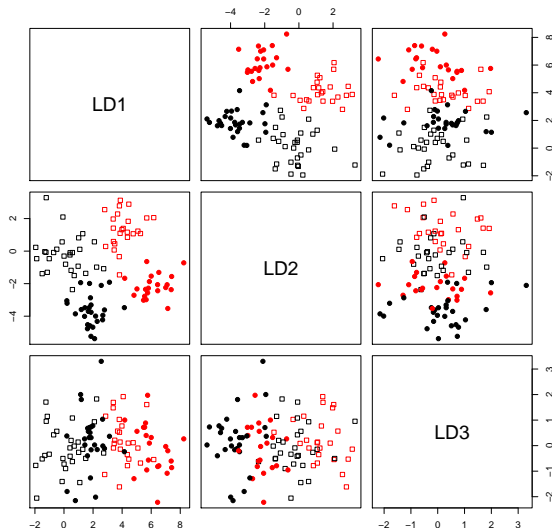
Remark: The number of non null eigen-values of $\mathbf{W}^{-1} \mathbf{B}$ is at most $\min(d, K - 1)$.

Classification rule

If components within class variance is scaled to one, classification can be performed by nearest class center with the Euclidean distance on the projected data.

Illustration of FDA

Factorial Discriminant analysis



Equivalence between LDA and FDA

Campbell (1984)

FDA \Leftrightarrow LDA s.c. $\text{rank}(\{\boldsymbol{\mu}_1, \dots, \boldsymbol{\mu}_K\}) = p$, where $p \leq K - 1$.

Reparametrization of LDA under constrains

$$\mathbf{x} | z_k = 1 \sim \mathcal{N}_d \left(\mathbf{A} \begin{pmatrix} \boldsymbol{\nu}_k \\ \boldsymbol{\gamma} \end{pmatrix}, \mathbf{A} \mathbf{A}^\top \right)$$

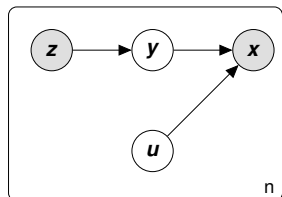
where $\boldsymbol{\nu}_k \in \mathbb{R}^p$, $\boldsymbol{\gamma} \in \mathbb{R}^{d-p}$ and $\mathbf{A} \in \mathcal{M}_{d,d}(\mathbb{R})$.

Generative interpretation

Let $\mathbf{y} \in \mathbb{R}^p$ and $\mathbf{u} \in \mathbb{R}^{d-p}$. The generative model is:

- Draw $\mathbf{z} : \mathbf{z} \sim \text{Mult}(1; \pi_1, \dots, \pi_K)$
- Draw $\mathbf{y} | \mathbf{z} : \mathbf{y} | z_k = 1 \sim \mathcal{N}_p(\boldsymbol{\nu}_k, \mathbf{I}_p)$
- Draw $\mathbf{u} : \mathbf{u} \sim \mathcal{N}_{d-p}(\boldsymbol{\gamma}, \mathbf{I}_{d-p})$
- Compute \mathbf{x} based on \mathbf{y} and \mathbf{u} :

$$\mathbf{x} = \mathbf{A} \begin{pmatrix} \mathbf{y} \\ \mathbf{u} \end{pmatrix}$$



Posterior class membership probabilities

Posterior probabilities

Let $\mathbf{V} \in \mathcal{M}_{p,d}(\mathbb{R})$ and $\mathbf{R} \in \mathcal{M}_{d-p,d}(\mathbb{R})$ the matrices such that $\mathbf{y} = \mathbf{V}\mathbf{x}$, $\mathbf{u} = \mathbf{R}\mathbf{x}$. Obviously $\begin{pmatrix} \mathbf{V} \\ \mathbf{R} \end{pmatrix} = \mathbf{A}^{-1}$.

$$p(z_k = 1 | \mathbf{x}) = p(z_k = 1 | \mathbf{y}, \mathbf{u}) = p(z_k = 1 | \mathbf{y}) = p(z_k = 1 | \mathbf{V}\mathbf{x}).$$

$\mathbf{V}\mathbf{x}$ conveys all the clustering information.

$$p(z_k = 1 | \mathbf{x}) = \frac{\pi_k \phi_p(\mathbf{V}\mathbf{x}; \boldsymbol{\nu}_k, \mathbf{I}_p)}{\sum_{k'=1}^K \pi_{k'} \phi_p(\mathbf{V}\mathbf{x}; \boldsymbol{\nu}_{k'}, \mathbf{I}_p)}.$$

Parameters estimation

$$\begin{aligned}
 \ell(\boldsymbol{\theta}; \mathbf{x}, \mathbf{z}) = & n \log \left| \det \begin{pmatrix} \mathbf{V} \\ \mathbf{R} \end{pmatrix} \right| - \overbrace{\frac{1}{2} \sum_{i=1}^n \sum_{k=1}^K z_{ik} \|\mathbf{V} \mathbf{x}_i - \boldsymbol{\nu}_k\|^2}^{\text{classifying}} \\
 & - \underbrace{\frac{1}{2} \sum_{i=1}^n \|\mathbf{R} \mathbf{x}_i - \boldsymbol{\gamma}\|^2}_{\text{not classifying}} - \frac{n}{2} \log(2\pi) + \sum_{i=1}^n \sum_{k=1}^K z_{ik} \log(\pi_k).
 \end{aligned}$$

Maximum likelihood estimation

- $\hat{\mathbf{V}}$ and $\hat{\mathbf{R}}$ are obtained through eigen value decomposition of $\mathbf{W}^{-1} \mathbf{B}$.
- $\hat{\nu}_1, \dots, \hat{\nu}_K$ and $\hat{\gamma}$ explicit given $\hat{\mathbf{V}}$ and $\hat{\mathbf{R}}$.

Application in the clustering setting (\mathbf{z} unknown)

EM algorithm

Until convergence iterate the following steps:

- E step: Compute $t_{ik}^{(r+1)} = p(\mathbf{z}_i | \mathbf{x}_i; \boldsymbol{\theta}^{(r)})$
- M step:
 - Compute $t_{ik}^{(r+1)}$ -weighted versions of \mathbf{W} , \mathbf{B} and $\bar{\mathbf{x}}_k$
 - Deduce $\mathbf{V}^{(r+1)}$ and $\mathbf{R}^{(r+1)}$, $\nu_1^{(r+1)}$, \dots , $\nu_K^{(r+1)}$ and $\gamma^{(r+1)}$

Remark

When $p < \min(K - 1, d)$, not equivalent to perform a standard EM algorithm and then perform FDA at the end of the EM algorithm.

Presentation of the model (supervised setting)

Model description

- H independent class variables
 - $\mathbf{z}^1, \dots, \mathbf{z}^H$ with K_1, \dots, K_H modalities
 - $p(z_k^h = 1) = \pi_k^h$
- H classifying latent variables
 - $\mathbf{y}^h \in \mathbb{R}^{p_h}$ with $p_h \leq K_h - 1$
 - $\mathbf{y}^h | z_k^h = 1 \sim \mathcal{N}_{p_h}(\boldsymbol{\nu}_k^h, I_{p_h})$
- Not classifying latent variables
 - $\mathbf{u} \in \mathbb{R}^{d-p_\bullet}$ with $p_\bullet = \sum_{h=1}^H p_h$
 - $\mathbf{u} \sim \mathcal{N}_{d-p_\bullet}(\boldsymbol{\gamma}, I_{d-p_\bullet})$

- \mathbf{x} defined by $\mathbf{x} = \begin{pmatrix} \mathbf{V}_1 \\ \vdots \\ \mathbf{V}_H \\ \mathbf{R} \end{pmatrix}^{-1} \begin{pmatrix} \mathbf{y}^1 \\ \vdots \\ \mathbf{y}^H \\ \mathbf{u} \end{pmatrix}$.

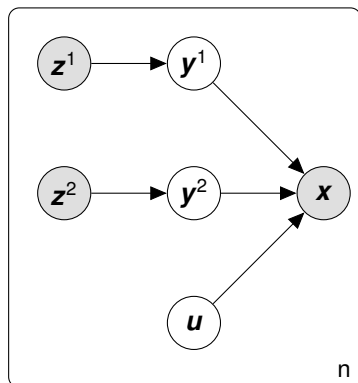


Figure: Illustration for $H = 2$

Possible models and number of parameters

Remarks

- Need to choose H, K_1, \dots, K_H and p_1, \dots, p_H : possible selection by model choice.
- $\mathbf{V}_1, \dots, \mathbf{V}_H$ and \mathbf{R} estimated up to isometric transformations.

Number of parameters

$$\dim(\Theta) = \sum_{h=1}^H (K_h - 1) + \sum_{h=1}^H \frac{p_h(2K_h + 2d - p_h + 1)}{2} + \frac{(d - p_\bullet)(d + p_\bullet + 3)}{2}$$

Estimation in the supervised setting

- The likelihood cannot be maximized directly \Rightarrow alternate optimisation.
- Constrain $\mathbf{V}_h^{(r+1)}$ and $\mathbf{R}^{(r+1)}$ to be linear combinations of $\mathbf{V}_h^{(r)}$ and $\mathbf{R}^{(r)}$, the parameters of other clustering dimension being fixed:

$$\begin{pmatrix} \mathbf{V}_h^{(r+1)} \\ \mathbf{R}^{(r+1)} \end{pmatrix} = \mathbf{M} \begin{pmatrix} \mathbf{V}_h^{(r)} \\ \mathbf{R}^{(r)} \end{pmatrix} = \begin{pmatrix} \mathbf{M}_1 \\ \mathbf{M}_2 \end{pmatrix} \begin{pmatrix} \mathbf{V}_h^{(r)} \\ \mathbf{R}^{(r)} \end{pmatrix},$$

- by denoting $\mathbf{y}^{h(r)} = \mathbf{V}_h^{(r)} \mathbf{x}$ and $\mathbf{u}^{(r)} = \mathbf{R}^{(r)} \mathbf{x}$,
- finding \mathbf{M} and other parameters is equivalent to perform probabilistic FDA on the data $\begin{pmatrix} \mathbf{y}^{h(r)} \\ \mathbf{u}^{(r)} \end{pmatrix}$

Estimation in the unsupervised setting

Similar to the supervised setting except that the data at each iteration are now weighted by $t_{ik}^{h(r+1)}$ instead of z_{ik}^h

EM algorithm

Until convergence, for $h \in \{1, \dots, H\}$ iterate the following steps:

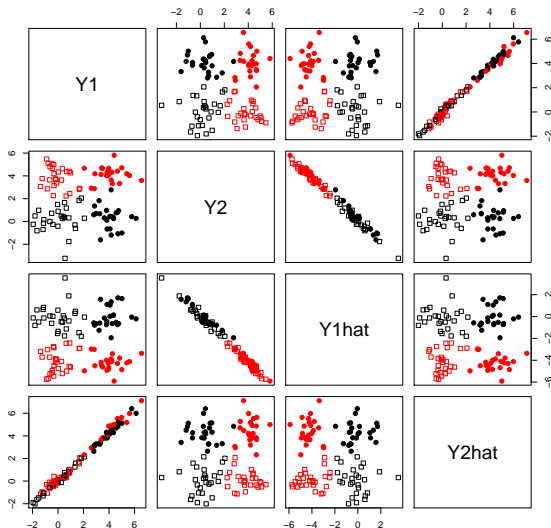
- E step: compute

$$t_{ik}^{h(r+1)} = p(\mathbf{z}_i^h | \mathbf{x}_i; \theta^{(r)}) = \frac{\pi_k \phi_{p_h}(\mathbf{y}_i^{h(r)}; \boldsymbol{\nu}_k^{h(r)}, \mathbf{I}_{p_h})}{\sum_{k'=1}^K \pi_{k'} \phi_{p_h}(\mathbf{y}_i^{h(r)}; \boldsymbol{\nu}_{k'}^{h(r)}, \mathbf{I}_{p_h})}$$

- M step: update \mathbf{V}_h , \mathbf{R} , $\boldsymbol{\nu}_1^h, \dots, \boldsymbol{\nu}_{K_h}^h$, $\pi_1^h, \dots, \pi_{K_h}^h$ and γ

Simulated data set: illustration

Comparison of true and estimated projections



Simulated data: model choice

- Performed using the BIC criterion

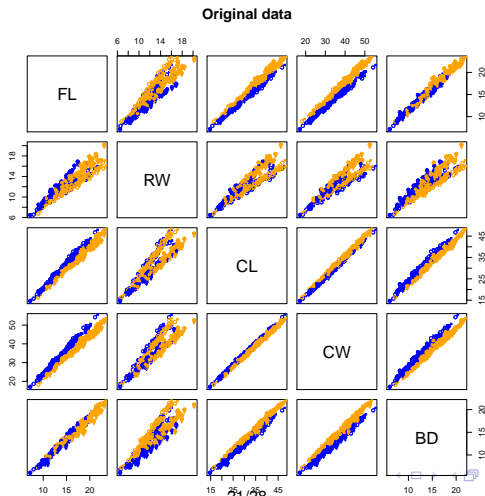
$$BIC = \ell(\mathbf{x}; \hat{\theta}) - \frac{\dim(\Theta)}{2} \log n$$

- Model collection: $H = 2$, $K_1 \in \{1, \dots, 5\}$, $K_2 \in \{1, \dots, 5\}$,
 $p_1 = p_2 = 1$

$K_1 \setminus K_2$	1	2	3	4	5
1	-1414.69	-1395.53	-1397.55	-1395.65	-1396.88
2		-1383.76	-1384.64	-1397.17	-1393.86
3			-1384.05	-1396.36	-1395.36
4				-1392.30	-1394.01
5					-1389.14

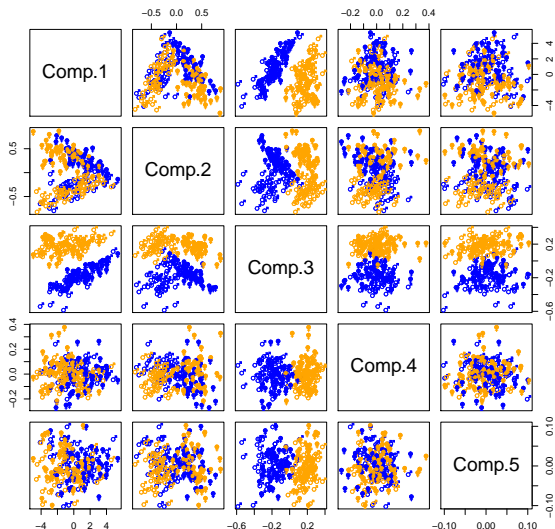
Crabs data set

200 crabs morphological data, 5 morphological variables, 50 males orange, 50 males blue, 50 females orange, 50 females blue.



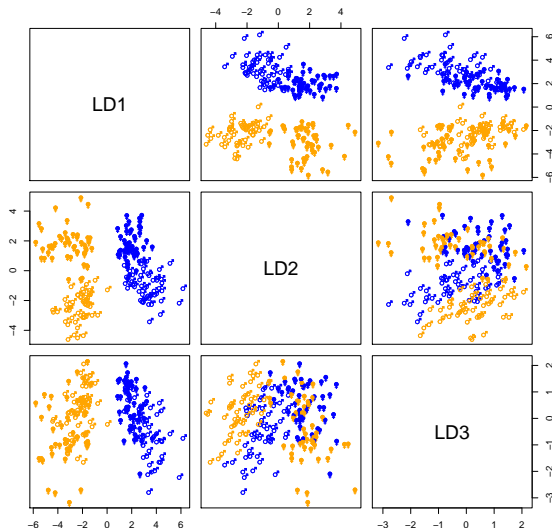
PCA on the crab data set

Data after PCA



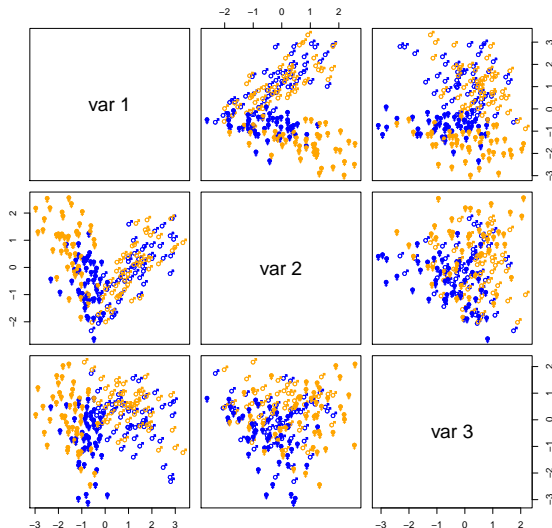
FDA on the crab data set: supervised setting

Data after Factorial Discriminant analysis



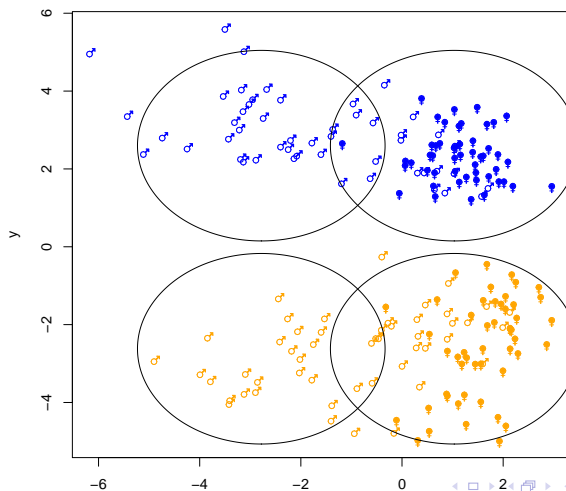
FDA on the crab data set: unsupervised setting

Data after Factorial Discriminant Analysis



Principal clustering components: $H = 2$, $K_1 = K_2 = 2$, $p_1 = p_2 = 1$

Data on the principal clustering components



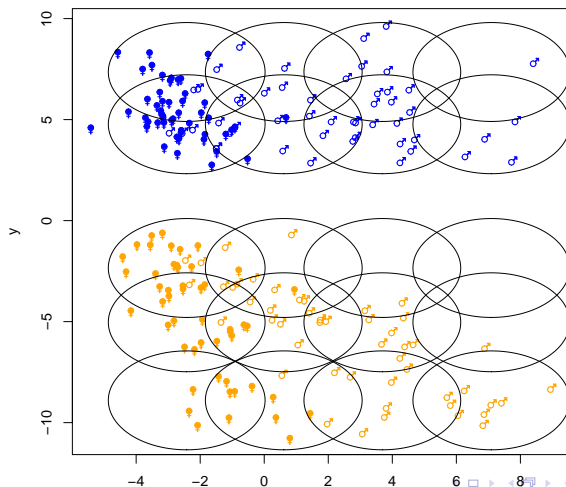
Principal clustering components: model choice

- Performed using the BIC criterion
- Model collection: $H = 2$, $K_1 \in \{1, \dots, 5\}$, $K_2 \in \{1, \dots, 5\}$, $p_1 = p_2 = 1$

$K_1 \setminus K_2$	1	2	3	4	5
1	-313.02	-252.35	-244.20	-250.60	-249.40
2		-229.73	-254.76	-254.82	-233.91
3			-248.82	-255.49	-233.17
4				-238.88	-231.19
5					-269.75

Principal clustering components: $H = 2$, $K_1 = 4$, $K_2 = 5$, $p_1 = p_2 = 1$

Data on the principal clustering components



Conclusion and perspectives

Conclusion

- model combining visualization and clustering with many clustering view points
- possibility to perform model choice

Perspectives

- consider sparse estimates of the \mathbf{R} matrix for the high dimensional setting
- extension to the heteroscedastic setting