

AG DANK/BCS Meeting 2013 in London  
University College London, 8/9 November 2013

# MODELS FOR SIMULTANEOUS CLASSIFICATION AND REDUCTION OF THREE-WAY DATA

**Roberto Rocci**  
University "Tor Vergata", Rome



# A general classification model: Gaussian Mixtures

Let  $\mathbf{x}=[x_1, x_2, \dots, x_J]'$  be a random vector of  $J$  variables. We assume

$$f(\mathbf{x}) = \sum_{g=1}^G p_g \phi_g(\mathbf{x}), \quad p_g \geq 0, \sum_{g=1}^G p_g = 1 \quad \text{mixture model}$$

where each component represents an underlying group, in our case

$$\phi_g(\mathbf{x}) = (2\pi)^{-\frac{J}{2}} |\boldsymbol{\Sigma}_g|^{-\frac{1}{2}} \exp\left\{-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_g)' \boldsymbol{\Sigma}_g^{-1} (\mathbf{x} - \boldsymbol{\mu}_g)\right\} \quad \text{Gaussian}$$

and each observation is assigned to a group by computing

$$p(g | \mathbf{x}) = \frac{p_g \phi_g(\mathbf{x})}{\sum_h p_h \phi_h(\mathbf{x})} \quad \text{posterior probabilities}$$

Given a sample of  $N$  i.i.d. observations, the parameters are estimated by maximizing

$$L(\boldsymbol{\vartheta}) = \sum_n \log\left(\sum_g p_g \phi_g(\mathbf{x}_n)\right) \quad \text{log-likelihood}$$

## Problems

- a very large number of parameters;
- difficult to understand which are the “discriminant” variables, i.e. the variables that describe the clustering structure.

## Idea

The mixture model induces the following covariance structure

$$\begin{aligned} \text{Var}(\mathbf{x}) &= \overbrace{\sum_{g=1}^G p_g (\boldsymbol{\mu}_g - \boldsymbol{\mu})(\boldsymbol{\mu}_g - \boldsymbol{\mu})'}^{\text{Between}} + \overbrace{\sum_{g=1}^G p_g \boldsymbol{\Sigma}_g}^{\text{Within}} && \text{variance decomposition} \\ &= \boldsymbol{\Sigma}_B + \boldsymbol{\Sigma}_W \end{aligned}$$

## Model the Between covariance matrix to:

- reduce the number of parameters;
- find the components (linear combinations of the variables) explaining the “largest information” about the classification.

# Reduction Model

The model is a “component analysis” of the centroid matrix.

## Scalar

$$\mu_{jg} = \mu_j + \sum_{q=1}^Q b_{jq} \eta_{qg}, \quad \sum_{g=1}^G p_g \eta_{qg} = 0$$

where:

$\mu_{jg}$  is the mean of variable  $j$  in component  $g$ ;

$\eta_{qg}$  is the mean of *prototype variable*  $q$  in component  $g$ ;

$b_{jq}$  is the loading of variable  $j$  on *prototype variable*  $q$ .

## Vector

$$\boldsymbol{\mu}_g = \boldsymbol{\mu} + \mathbf{B}\boldsymbol{\eta}_g, \quad \sum_{g=1}^G p_g \boldsymbol{\eta}_g = \mathbf{0}$$

## Matrix

$$\mathbf{M} = \mathbf{N}\mathbf{B}', \quad \mathbf{1}'\mathbf{N} = \mathbf{0}$$

where:

- $\mathbf{M} = [\boldsymbol{\mu}_1, \boldsymbol{\mu}_2, \dots, \boldsymbol{\mu}_G]' - \mathbf{1}\boldsymbol{\mu}'$ , (centred) centroid matrix;
- $\mathbf{N} = [\boldsymbol{\eta}_1, \boldsymbol{\eta}_2, \dots, \boldsymbol{\eta}_G]'$ , centroid matrix on the reduced space.

The component model is not identified. In fact

$$\mathbf{B}\boldsymbol{\eta}_g = \mathbf{B}\mathbf{F}^{-1}\mathbf{F}\boldsymbol{\eta}_g = \tilde{\mathbf{B}}\tilde{\boldsymbol{\eta}}_g.$$

We exploit such rotational freedom by requiring that

$$\mathbf{B}'\boldsymbol{\Sigma}^{-1}\mathbf{B} = \mathbf{I}_Q.$$

# ML Estimation (homoscedastic case): EM algorithm

Maximization of the loglikelihood

$$L(\vartheta) = \sum_{n=1}^N \log \left( \sum_{g=1}^G p_g \phi_g(\mathbf{x}_n) \right) \quad \text{objective}$$

is equivalent to the maximization of the “fuzzy” function (Hathaway, 1986)

$$l(\vartheta) = \sum_{ng} u_{ng} \log(p_g \phi_g(\mathbf{x}_n)) - \sum_{ng} u_{ng} \log(u_{ng}) \quad \text{fuzzy objective}$$

where  $u_{ng} \geq 0$  and  $\sum_g u_{ng} = 1$ . This is so because  $l(\vartheta)$  reaches a maximum respect to  $\mathbf{U} = [u_{ng}]$  when

$$u_{ng} = \frac{p_g \phi_g(\mathbf{x}_n)}{\sum_h p_h \phi_h(\mathbf{x}_n)} \quad \text{posterior probabilities}$$

Substituting the previous in  $l(\vartheta)$  we obtain  $L(\vartheta)$ .

The algorithm is based on the conditional maximization of  $l(\mathcal{G})$  with respect to a subset of parameters given the others.

The fundamental steps are the following.

a) Update  $\mathbf{U}=[u_{ng}]$ :

$$u_{ng} = \frac{p_g \phi_g(\mathbf{x}_n)}{\sum_h p_h \phi_h(\mathbf{x}_n)}, n=1,2,\dots,N; g=1,2,\dots, G,$$

b) Update  $\mathbf{p}=[p_g]$ :

$$p_g = \frac{1}{N} \sum_n u_{ng}, g=1,2,\dots,G.$$

c) Update  $\Sigma$ :

$$\Sigma = \frac{1}{N} \sum_{ng} u_{ng} (\mathbf{x}_n - \boldsymbol{\mu}_g)(\mathbf{x}_n - \boldsymbol{\mu}_g)'$$

They are simply the steps of a ordinary EM algorithm.

d) Update  $\boldsymbol{\mu}$ :

We consider centered data,  $\boldsymbol{\mu} = \mathbf{0}$ .

e) Update  $\mathbf{N}$  and  $\mathbf{B}$ :

It can be shown that the objective function can be written as

$$l(\mathfrak{G}) = -\frac{1}{2} \text{tr}\{\mathbf{D}(\bar{\mathbf{X}} - \mathbf{NB}')\boldsymbol{\Sigma}^{-1}(\bar{\mathbf{X}} - \mathbf{NB}')'\} + c$$

where  $c$  is a constant term (independent of  $\mathbf{N}$  and  $\mathbf{B}$ ),  $\mathbf{D} = \text{diag}(u_{+1}, u_{+2}, \dots, u_{+G})$  and  $\bar{\mathbf{X}} = [\bar{\mathbf{x}}_1, \bar{\mathbf{x}}_2, \dots, \bar{\mathbf{x}}_G]'$  is the matrix of centroids,  $\bar{\mathbf{x}}_g = \frac{1}{\sum_n u_{ng}} \sum_n u_{ng} \mathbf{x}_n$ , computed on the centred variables.

This algorithm can be also seen as an ECM (Meng & Rubin, 1993).



# Use and interpretation of components

Step M of the EM algorithm shows that:

1) the within-standardized component loadings matrix  $\widehat{\mathbf{B}} = \boldsymbol{\Sigma}^{-\frac{1}{2}} \mathbf{B}$  derives from a PCA of the matrix of within-standardized centroids

$$\text{tr}\{\mathbf{D}(\bar{\mathbf{X}} - \mathbf{N}\mathbf{B}')\boldsymbol{\Sigma}^{-1}(\bar{\mathbf{X}} - \mathbf{N}\mathbf{B}')'\} = \left\| \mathbf{D}^{\frac{1}{2}}\bar{\mathbf{X}}\boldsymbol{\Sigma}^{-\frac{1}{2}} - \mathbf{D}^{\frac{1}{2}}\mathbf{N}\mathbf{B}'\boldsymbol{\Sigma}^{-\frac{1}{2}} \right\|^2 = \left\| \mathbf{D}^{\frac{1}{2}}\bar{\mathbf{Z}} - \mathbf{D}^{\frac{1}{2}}\mathbf{N}\widehat{\mathbf{B}}' \right\|^2 \rightarrow \min_{\mathbf{N}, \mathbf{B}}$$

2) the component scores

$$\mathbf{Y} = \mathbf{Z}\widehat{\mathbf{B}} = \mathbf{X}\boldsymbol{\Sigma}^{-\frac{1}{2}}\boldsymbol{\Sigma}^{-\frac{1}{2}}\mathbf{B} = \mathbf{X}\boldsymbol{\Sigma}^{-1}\mathbf{B} = \mathbf{X}\check{\mathbf{B}}$$

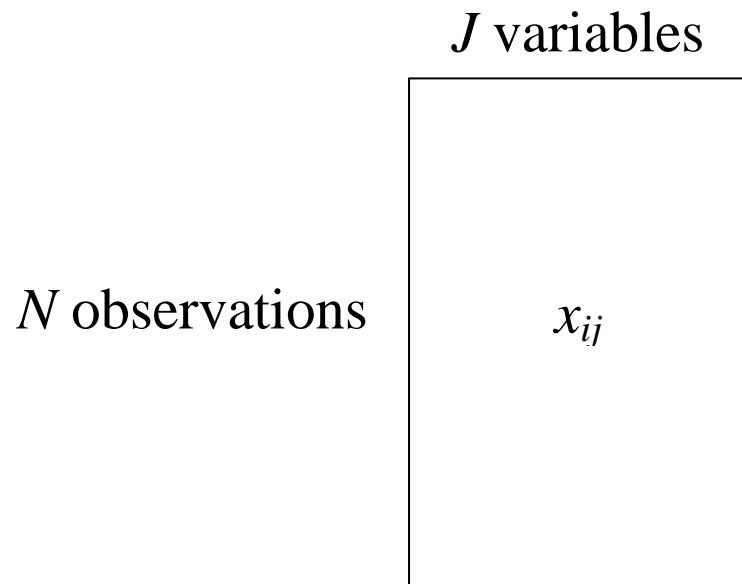
maximize the between variance subject to the constraint of unit within variance, i.e.

$$\begin{aligned} & \max \text{tr}(\check{\mathbf{B}}'\bar{\mathbf{X}}'\mathbf{D}\bar{\mathbf{X}}\check{\mathbf{B}}) \\ & \text{subject to } \check{\mathbf{B}}'\boldsymbol{\Sigma}\check{\mathbf{B}} = \mathbf{B}'\boldsymbol{\Sigma}^{-1}\boldsymbol{\Sigma}\boldsymbol{\Sigma}^{-1}\mathbf{B} = \mathbf{B}'\boldsymbol{\Sigma}^{-1}\mathbf{B} = \mathbf{I}_Q. \end{aligned}$$

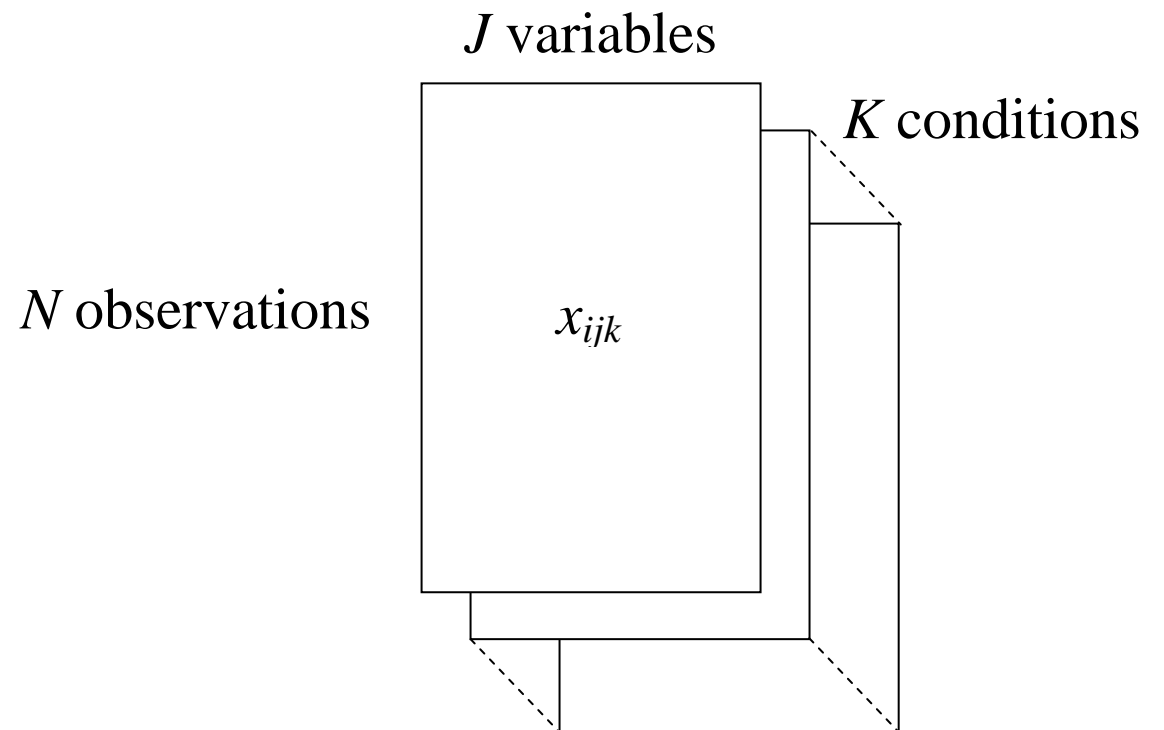
**Fisher linear discriminant analysis (LDA)**

# Three-way Extension

## Two-way sample



## Three-way sample



Let  $\mathbf{x} = [x_{11}, x_{21}, \dots, x_{J1}, \dots, x_{1K}, x_{2K}, \dots, x_{JK}]'$  be a random vector of  $J$  variables observed under  $K$  different conditions.

## General classification model

$$f(\mathbf{x}) = \sum_{g=1}^G p_g \phi_g(\mathbf{x})$$

**mixture model**

where

$$\phi_g(\mathbf{x}) = (2\pi)^{-\frac{JK}{2}} |\boldsymbol{\Sigma}_g|^{-\frac{1}{2}} \exp\left\{-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_g)' \boldsymbol{\Sigma}_g^{-1} (\mathbf{x} - \boldsymbol{\mu}_g)\right\}$$

**Gaussian components**

## Problems

- a very large number of parameters;
- difficult to understand which are the “discriminant” variables and/or occasions;
- difficult to distinguish the role of variables from that of occasions.

# Within Covariance Structure

## Model

Direct Product (Browne, 1984)

$$\Sigma_g = \Sigma_O \otimes \Sigma_V = \begin{bmatrix} \sigma_{11} \Sigma_V & \cdots & \sigma_{1K} \Sigma_V \\ \vdots & \ddots & \vdots \\ \sigma_{K1} \Sigma_V & \cdots & \sigma_{KK} \Sigma_V \end{bmatrix}$$

in scalar notation

$$\sigma_{jklm} = \sigma_{jl} \sigma_{km}$$

Basford & MacLachlan (1985) proposed

$$\Sigma_g = \mathbf{I}_K \otimes \Sigma_{V;g}$$

# Reduction Model

The model is a “Tucker 2 component analysis” of the centroid matrix.

## Scalar

$$\mu_{jkg} = \mu_{jk} + \sum_{q=1}^Q \sum_{r=1}^R b_{jq} c_{kr} \eta_{qrg}, \quad \sum_{g=1}^G p_g \eta_{qrg} = 0$$

where:

- $\mu_{jkg}$  is the mean of variable  $j$  under condition  $k$  in component  $g$ ;
- $\eta_{qrg}$  is the mean of *prototype variable*  $q$  under *prototype condition*  $r$  in component  $g$ ;
- $\sum_q b_{jq} \eta_{qrg}$  is the mean of variable  $j$  under *prototype condition*  $r$  in component  $g$ ;
- $b_{jq}$  is the loading of variable  $j$  on *prototype variable*  $q$ ;
- $\sum_r c_{kr} \eta_{qrg}$  is the mean of *prototype variable*  $q$  under condition  $k$  in component  $g$ ;
- $c_{kr}$  is the loading of occasion  $k$  on *prototype occasion*  $r$ .

Often used in Chemistry and Psychology, see <http://three-mode.leidenuniv.nl/>

## Vector

$$\boldsymbol{\mu}_g = \boldsymbol{\mu} + (\mathbf{C} \otimes \mathbf{B})\boldsymbol{\eta}_g, \quad \sum_{g=1}^G p_g \boldsymbol{\eta}_g = \mathbf{0}$$

## Matrix

$$\mathbf{M} = \mathbf{N}(\mathbf{C}' \otimes \mathbf{B}'), \quad \mathbf{1}'\mathbf{N} = \mathbf{0}$$

where:

- $\mathbf{M} = [\boldsymbol{\mu}_1, \boldsymbol{\mu}_2, \dots, \boldsymbol{\mu}_G]' - \mathbf{1}\boldsymbol{\mu}'$ , (centred) centroid matrix;
- $\mathbf{N} = [\boldsymbol{\eta}_1, \boldsymbol{\eta}_2, \dots, \boldsymbol{\eta}_G]'$ , centroid matrix on the reduced space.

The component model is not identified. In fact

$$(\mathbf{C} \otimes \mathbf{B})\boldsymbol{\eta}_g = (\mathbf{C} \otimes \mathbf{B})(\mathbf{D}^{-1} \otimes \mathbf{F}^{-1})(\mathbf{D} \otimes \mathbf{F})\boldsymbol{\eta}_g = (\mathbf{C}\mathbf{D}^{-1} \otimes \mathbf{B}\mathbf{F}^{-1})\tilde{\boldsymbol{\eta}}_g = (\tilde{\mathbf{C}} \otimes \tilde{\mathbf{B}})\tilde{\boldsymbol{\eta}}_g.$$

We exploit such rotational freedom by requiring that

$$\mathbf{B}'\boldsymbol{\Sigma}_V^{-1}\mathbf{B} = \mathbf{I}_Q, \quad \mathbf{C}'\boldsymbol{\Sigma}_O^{-1}\mathbf{C} = \mathbf{I}_R.$$

## ML Estimation (homoscedastic case): EM algorithm

An EM algorithm can be programmed following the analogous algorithm already seen for the two-way case.

About the update of  $\mathbf{N}$ ,  $\mathbf{B}$  and  $\mathbf{C}$ , it is interesting to note that the complete log-likelihood can be written as

$$l(\vartheta) = -\frac{1}{2} \text{tr} \left\{ \mathbf{D} [\bar{\mathbf{X}} - \mathbf{N}(\mathbf{C} \otimes \mathbf{B})'] (\boldsymbol{\Sigma}_o^{-1} \otimes \boldsymbol{\Sigma}_v^{-1}) [\bar{\mathbf{X}} - \mathbf{N}(\mathbf{C} \otimes \mathbf{B})']' \right\} + c$$

where  $c$  is a constant term and  $\bar{\mathbf{X}}$  is the matrix of centroids computed on the centred variables.

It follows that the parameters can be updated by computing a weighted least squares approximation of the centroid matrix.

# Use and interpretation of components

1) the within-standardized component loadings matrices  $\hat{\mathbf{B}} = \boldsymbol{\Sigma}_V^{-\frac{1}{2}} \mathbf{B}$  and  $\hat{\mathbf{C}} = \boldsymbol{\Sigma}_O^{-\frac{1}{2}} \mathbf{C}$  derive from a Tucker2 analysis of the matrix of within-standardized centroids

$$\left\| \mathbf{D}^{\frac{1}{2}} \bar{\mathbf{X}} (\boldsymbol{\Sigma}_O^{-\frac{1}{2}} \otimes \boldsymbol{\Sigma}_V^{-\frac{1}{2}}) - \mathbf{D}^{\frac{1}{2}} \mathbf{N} (\mathbf{C}' \otimes \mathbf{B}') (\boldsymbol{\Sigma}_O^{-\frac{1}{2}} \otimes \boldsymbol{\Sigma}_V^{-\frac{1}{2}}) \right\|^2 = \left\| \mathbf{D}^{\frac{1}{2}} \bar{\mathbf{Z}} - \mathbf{D}^{\frac{1}{2}} \mathbf{N} (\hat{\mathbf{C}}' \otimes \hat{\mathbf{B}}') \right\|^2 \rightarrow \min_{\mathbf{N}, \mathbf{B}, \mathbf{C}}$$

2) the component scores

$$\mathbf{Y} = \mathbf{Z} (\hat{\mathbf{C}} \otimes \hat{\mathbf{B}}) = \mathbf{X} (\boldsymbol{\Sigma}_O^{-1} \otimes \boldsymbol{\Sigma}_V^{-1}) (\mathbf{C} \otimes \mathbf{B}) = \mathbf{X} (\check{\mathbf{C}} \otimes \check{\mathbf{B}})$$

maximize the between variance subject to the constraint of unit within variance, i.e.

$$\begin{aligned} & \max \text{tr} [(\check{\mathbf{C}} \otimes \check{\mathbf{B}})' \bar{\mathbf{X}}' \mathbf{D} \bar{\mathbf{X}} (\check{\mathbf{C}} \otimes \check{\mathbf{B}})] \\ \text{subject to } & \check{\mathbf{C}}' \boldsymbol{\Sigma}_O \check{\mathbf{C}} = \mathbf{I}_R, \check{\mathbf{B}}' \boldsymbol{\Sigma}_V \check{\mathbf{B}} = \mathbf{I}_Q \Leftrightarrow (\check{\mathbf{C}} \otimes \check{\mathbf{B}})' (\boldsymbol{\Sigma}_O \otimes \boldsymbol{\Sigma}_V) (\check{\mathbf{C}} \otimes \check{\mathbf{B}}) = \mathbf{I}_R \otimes \mathbf{I}_Q \end{aligned}$$

**Bilinear discriminant analysis (BLDA)**



# BLDA: interpretation

## Constrained LDA

$$y_{qr} = \sum_{j=1}^J \sum_{k=1}^K x_{jk} w_{jkqr} \Leftrightarrow \mathbf{y} = (\check{\mathbf{c}} \otimes \check{\mathbf{b}})' \mathbf{x}$$

where

$$w_{jkqr} = \check{b}_{jq} \check{c}_{kr}$$

Dimensionality reduction of the variables

Dimensionality reduction of the occasions

## Hierarchical LDA

$$\left\{ \begin{array}{l} y_{qr} = \sum_{j=1}^J \check{b}_{jq} f_{jr} \Leftrightarrow f_{jr} = \sum_{k=1}^K \check{c}_{kr} x_{jk} \\ y_{qr} = \sum_{k=1}^K \check{c}_{kr} h_{qk} \Leftrightarrow h_{qk} = \sum_{j=1}^J \check{b}_{jq} x_{jk} \end{array} \right. \Rightarrow \left\{ \begin{array}{l} \check{b}_{jp} \text{ variable component weights} \\ \check{c}_{kq} \text{ occasion component weights} \end{array} \right.$$

# Application

## Data

58 units: soybeans;

8 conditions: 4 environments (Lawes, Brookstead, Nambour, Redland Bay)  
× 2 years (1970, 1971);

2 variables: yield Kg/Ha, protein.

## Model selection

Model considered:

$G = 2:7$ ,  $Q = 1:2$ ,  $R = 1:8$ ,  $\Sigma_o$  diagonal or with non null covariances only between the same locations.

Best model selected by BIC:

$$G = 7, Q = 2, R = 2 \text{ and } \Sigma_o \text{ diagonal.}$$

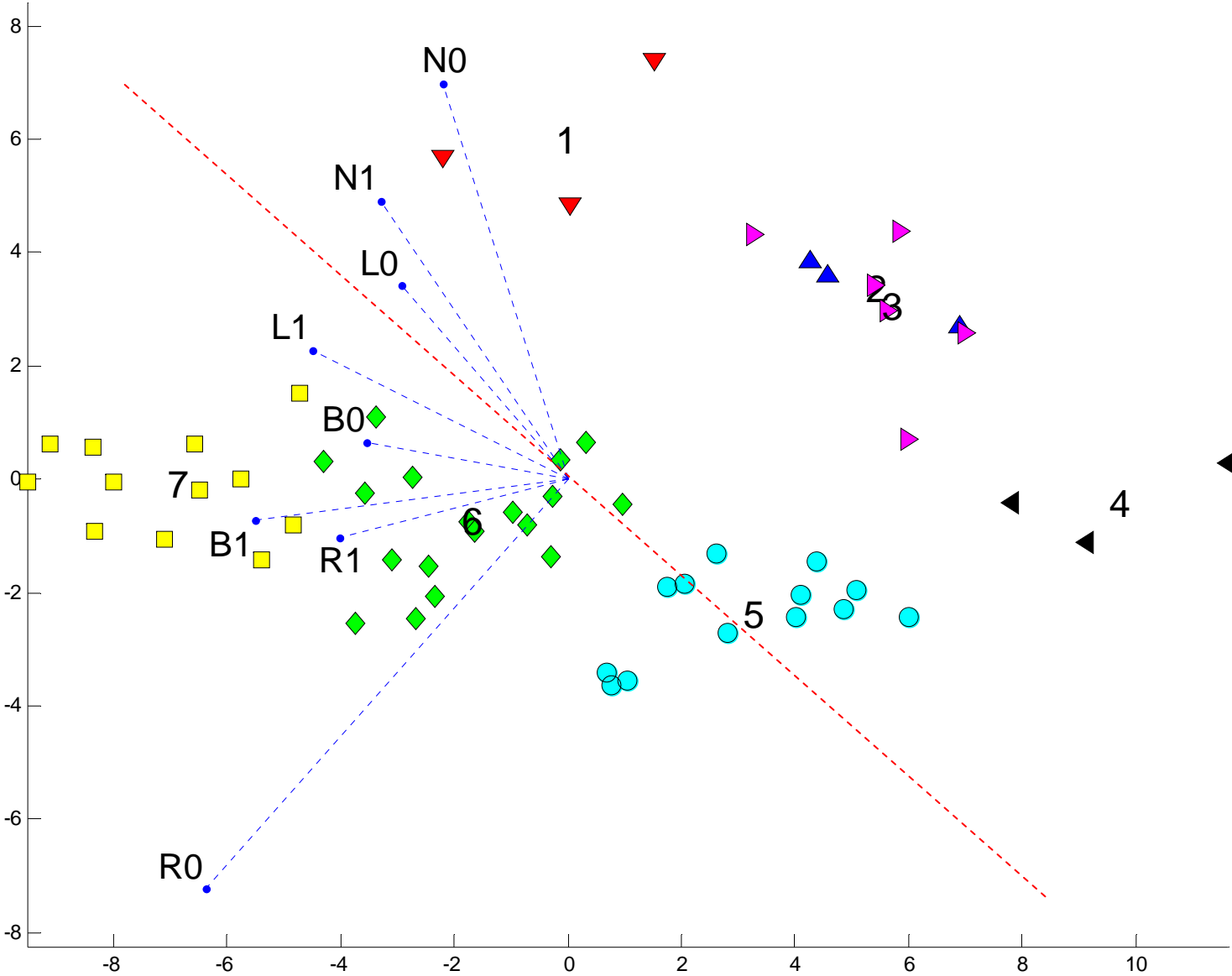
**Percentage of variation accounted for by the components  
on the within-standardized data**

<b>Variables</b>	<b>Occasions</b>		<b>Tot</b>
	<b>1</b>	<b>2</b>	
<b>1</b>	50.98	11.59	62.57
<b>2</b>	11.96	4.56	16.52
<b>Tot</b>	62.94	16.15	79.09

## Basford & McLachlan (B&M) and our (R) classification

B&M	R						
	1	2	3	4	5	6	7
1	3						
2		3					
3			6	3			
4					9		
5					3	6	
6					1	8	
7						4	12

# Biplot on the first latent variable at the two latent occasions



# Heteroscedastic case

## Reduction model

### Scalar

$$\mu_{jkg} = \mu_{jk} + \sum_{q=1}^J \sum_{r=1}^K b_{jq} c_{kr} \eta_{qrg}, \quad \sum_{g=1}^G p_g \eta_{qrg} = 0, \quad \eta_{qrg} = 0 \text{ if } q > Q \text{ and/or } r > R$$

### Vector

$$\boldsymbol{\mu}_g = \boldsymbol{\mu} + (\mathbf{C} \otimes \mathbf{B}) \boldsymbol{\eta}_g, \quad \sum_{g=1}^G p_g \boldsymbol{\eta}_g = \mathbf{0}$$

### Matrix

$$\mathbf{M} = \mathbf{N}(\mathbf{C}' \otimes \mathbf{B}'), \quad \mathbf{1}' \mathbf{N} = \mathbf{0}$$

where

- $\mathbf{C} = [\mathbf{C}_R, \mathbf{C}_{K-R}]$ , square,
- $\mathbf{B} = [\mathbf{B}_Q, \mathbf{B}_{J-Q}]$ , square.

## Within-covariance model

$$\mathbf{\Sigma}_g = (\mathbf{C} \otimes \mathbf{B}) \mathbf{\Omega}_g (\mathbf{C} \otimes \mathbf{B})'$$

where

$$\mathbf{\Omega}_g = \begin{bmatrix} \mathbf{\Omega}_{O,g} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix} \otimes \begin{bmatrix} \mathbf{\Omega}_{V,g} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix} + \mathbf{\Psi},$$

-  $\mathbf{\Psi}$  diagonal.

If  $K = 3$  and  $R = 2$ , we have

$$\mathbf{\Omega}_g = \left[ \begin{array}{cc|cc|cc} \omega_{11O,g} \mathbf{\Omega}_{V,g} & 0 & \omega_{12O,g} \mathbf{\Omega}_{V,g} & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ \hline \omega_{21O,g} \mathbf{\Omega}_{V,g} & 0 & \omega_{22O,g} \mathbf{\Omega}_{V,g} & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ \hline 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \end{array} \right] + \mathbf{\Psi}$$