

Model-based Statistical Learning



Pr. Charles BOUVEYRON

Professor of Statistics
Chair of the Institut 3IA Côte d'Azur
Université Côte d'Azur & Inria

charles.bouveyron@univ-cotedazur.fr
@cbouveyron

The mixture of PPCA

Tipping & Bishop proposed in 1996 the mixture of PPCA

$$Z \sim \mathcal{H}(1; \pi) \quad \pi = (\pi_1, \dots, \pi_K)$$

$$X_{|Z=h} = \gamma U_h + \epsilon_h \quad \text{where } U_h \text{ is a } p \times d \text{ matrix}$$

$$\epsilon_h \sim N(0, \Sigma_h \mathbb{I}_p)$$

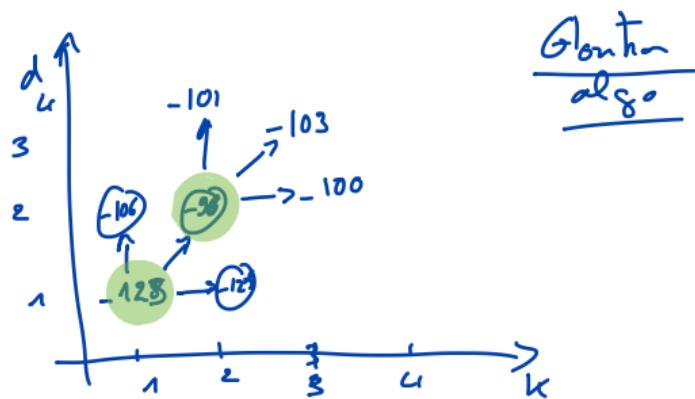
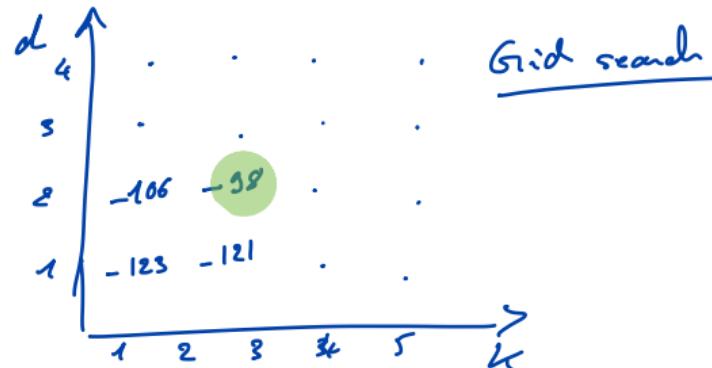
$$Y|Z=h \sim N(\mu_h, I_d)$$

$$\Rightarrow X \sim \sum_{h=1}^K \pi_h \phi(\mu_h U_h, U_h^t U_h + \Sigma_h^{-1})$$

The mixture of PPCA

The inference of MPPCA can be done using the EM algorithm.

The choice of both K and d can be done using model selection, with BIC for instance.



HDDC (High Dimensional Data Clustering, Bouveyron 06)

The idea of HDDC is to extend NPPCA on 2 directions:

- 1) allows the intrinsic dimensions of the subspaces to be different : $d \rightarrow d_k$.
- 2) allows sub-models of NPPCA to fit in different situations with fewer data.

HDDC

the model of HDDC, $Z \sim \text{cl}(1, \pi)$

$$X_{|Z=k} = YU_k + \varepsilon_k \quad \text{where } U_k \text{ is a } p \times d_k \text{ orthonormal matrix}$$

$$\varepsilon_k \sim N(0, \sigma_k^2 I_p)$$

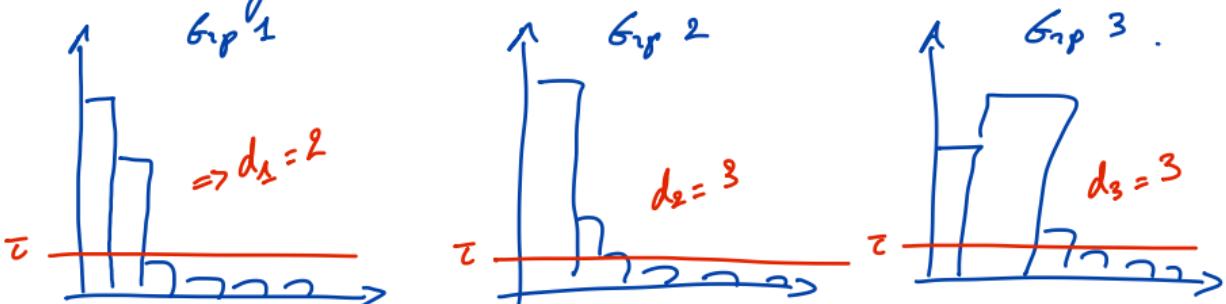
$$Y \sim N(\mu_k, \alpha_k I_{d_k}) \quad \text{where } \alpha_k \in \mathbb{R}^{d_k}$$

$$\Rightarrow X \sim \sum_{k=1}^n \pi_k N\left(\mu_k U_k, (\alpha_k U_k)(\alpha_k U_k)^T \sigma_k^2 I_p\right)$$

HDDC

Inferring this HDDC model is still possible with ER, but model selection is now very difficult because we have $(K+1)$ discrete hyper-parameters to choose.

A solution to avoid this combinatorial problem: the scree-test of Cattell



\Rightarrow Model selection here resumes to choosing K and τ .

In HDDC, it is also possible to get more parsimonious models by constraining some model parameters to be common between or within groups.

$$\forall h, d_h = d \Rightarrow \text{PPCA}$$

$$\forall h, \sigma_h^2 = \sigma^2 \Rightarrow [a_{hj} \ b_{hj} \ Q_h \ d_h]$$

$$a_{hj} = a_h \Rightarrow [a_h \ b_h \ Q_h \ d_h]$$

$$\sigma_h^2 = \sigma^2 \ \& \ a_{hj} = a_h \Rightarrow [a_h \ b \ Q_h \ d_h]$$

(...)

$$\Rightarrow [a \ b \ Q \ d]$$

\Rightarrow the R Library `HDDC` classif implements it.

\Rightarrow the Python package `HDDA` for Python github.com/mfauvel/HDDA

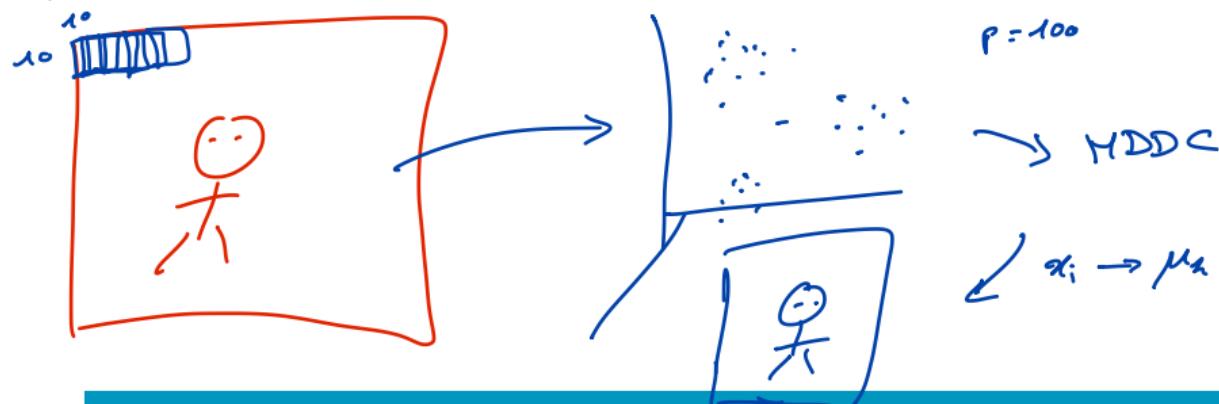
}

28 submodels.
which can be selected thanks to model selection.

In the supervised case: HDDA

In the supervised case, HDDA implements all the 28 models (including MPPCA) for supervised classification.

A remark: HDDC for image denoising turns out to be extremely performant: HDTI method implements it.



Conclusion on subspace clustering / classif.:

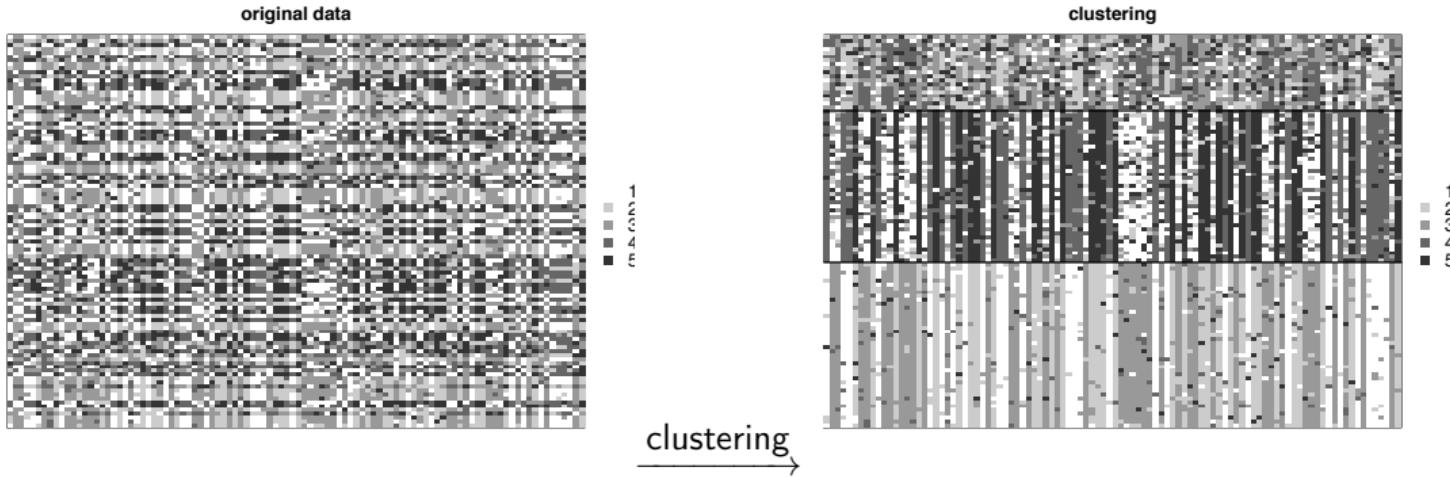
- ⊕ extremely performant models for HD data.
- ⊖ no easy way to visualize the data in the latent subspaces.
↳ the Fisher-EP algorithm aims to solve this visualization problem while keeping the idea of subspace clustering.

Outline

Co-Clustering

Co-clustering

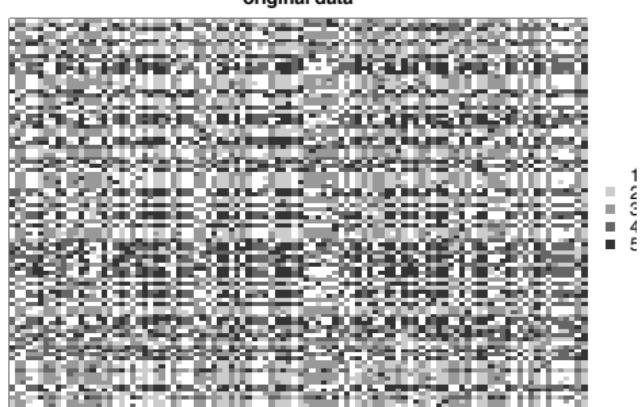
When the number of variables is high, the interpretation of the obtained clustering can remain difficult:



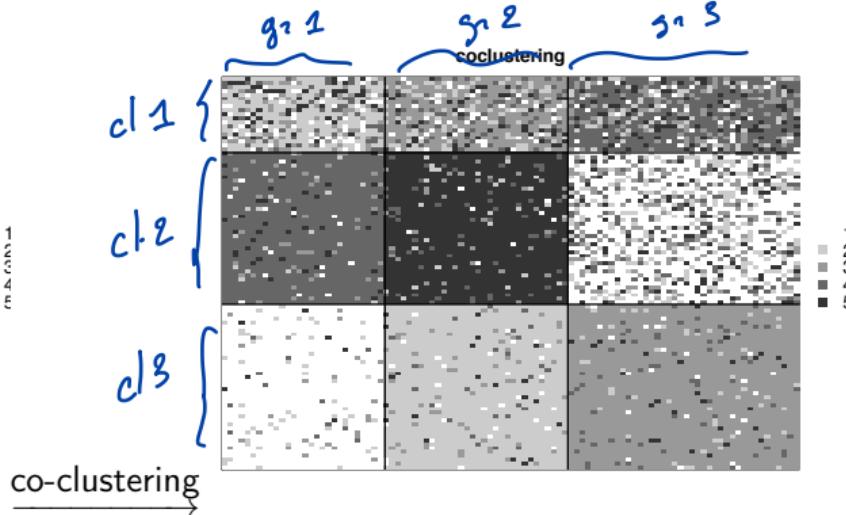
Co-clustering

In such a case, co-clustering may help by clustering simultaneously the rows and the columns of the data matrix:

cluster the rows



cluster the variables.



The latent block model (LBM)

The LBR model is designed for clustering both rows and columns of a data set:

Z : cluster variable for rows (into K groups)

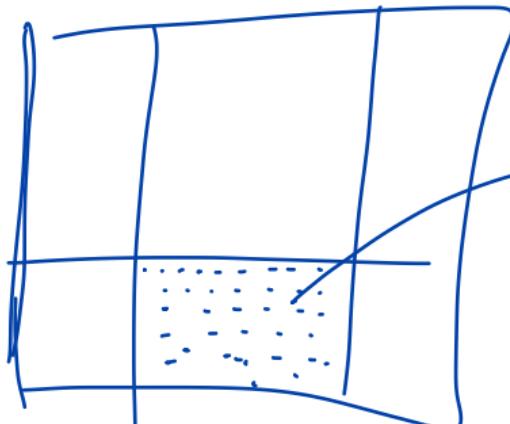
W : columns (into L groups)

$$Z \sim \text{dl}(1; \pi) \quad \text{and} \quad W \sim \text{dl}(1; \rho)$$

The latent block model (LBM)

Assumption: knowing Z and W , the individuals into a block defined by Z and W are all assumed to be independent of the same univariate pdf.

$$X_{ij} \mid Z_{ik}=1, W_{jl}=1 \sim f_{k\ell}(x_{ij}; \theta_{k\ell})$$



$\sim f_{k\ell}(x_{ij}; \theta_{k\ell})$
which is
a univariate
dist.

The latent block model (LBM)

As an illustration, the LBP model for binary data:

$$X_{ij} \in [0, 1]$$

$$Z \sim \text{dfl}(1; \pi)$$

$$W \sim \text{dfl}(1; \rho)$$

$$X_{ij} | Z_{ih}=1, W_{je}=1 \sim B(\rho_{he})$$

\Rightarrow the nb of parameters to estimate:

$$\begin{array}{lcl} \pi & \rightarrow & K-1 \\ \rho & \rightarrow & L-1 \\ \rho_{he} & \rightarrow & KL. \end{array} \quad \left. \right\} \quad KL + (K-1) + (L-1)$$

The latent block model (LBM)

- ⊕ This model is extremely parsimonious in term of models parameters and, therefore, has no other way to find very good clustering for rows and columns.
- ⊕ The LBN model can be adapted to various types of data by changing the univariate distribution.
 - Binary data \rightarrow Bernoulli.
 - count data \rightarrow Poisson
 - \rightarrow continuous data \rightarrow Gaussian.
 - categorical \rightarrow multinomial
 - compositional \rightarrow Dirichlet

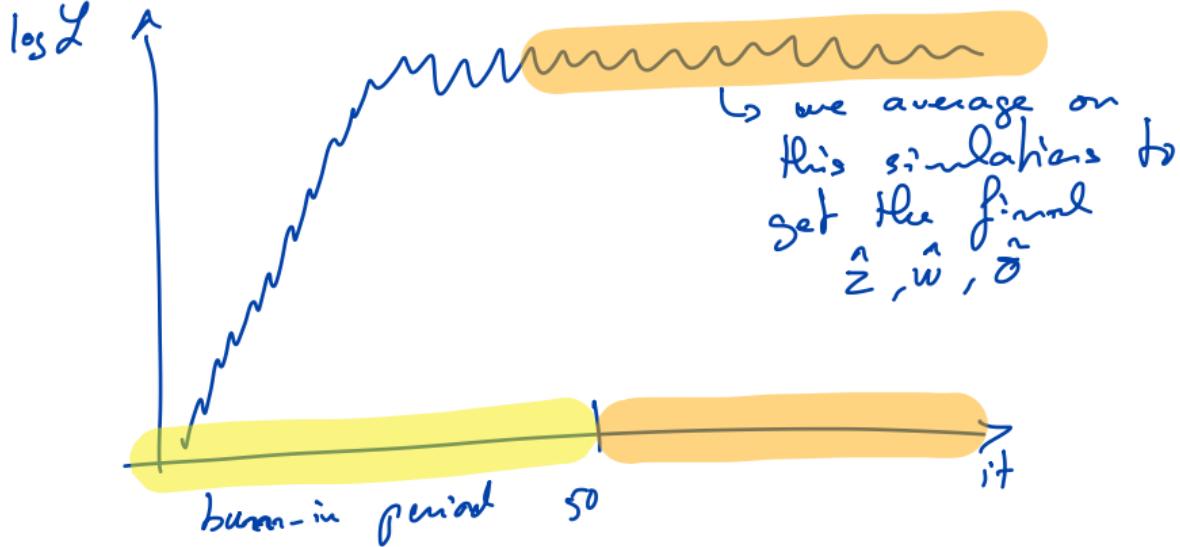
Inference of the LBM model

Unfortunately, it is not possible to directly use the EM algorithm because we have 2 latent variables.

One approach is to use a stochastic version of EM: the SEM-Gibbs algorithm.

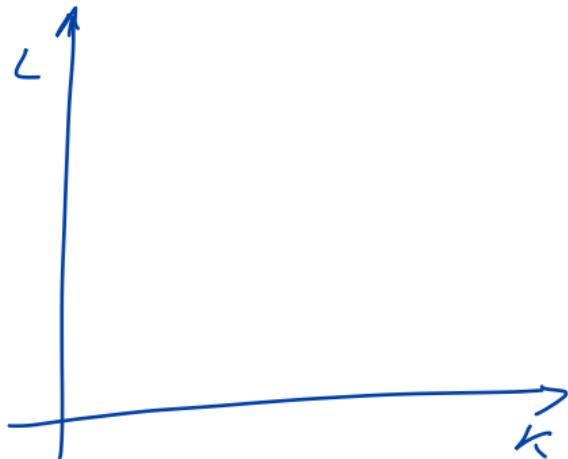
- SE step : we simulate Z and W according to the posterior conditional distribution (MCMC idea)
- T step : we find estimates of model parameters by maximizing the likelihood knowing some simulated values for Z and W .

This generates a sequence of posterior distributions which converges toward the actual posterior one.



Model selection

Model selection criteria (Bic, Aic, Icl) can be used to select K and L.



Grid search
or the gradient approach.

Application with R

The blockcluster package implements the LBR for several data types.

The Fisher-EM algorithm

$$X_{|z=h} = YU + \epsilon_t$$

but U is such that it discriminates best
the clusters

