

# Multi-components Data, Signal and Image Processing for Biological and Medical Applications

Ali Mohammad-Djafari

Laboratoire des Signaux et Systèmes  
UMR 8506 CNRS - CS - Univ Paris Sud  
CentraleSupélec, Gif-sur-Yvette.

[djafari@lss.supelec.fr](mailto:djafari@lss.supelec.fr)  
<http://djafari.free.fr>

January 6, 2017

## Summary 3: Data redundancy, Dimensionality Reduction, ...

- ▶ Redundancy and structure
- ▶ Dimensionality Reduction
- ▶ PCA and ICA
- ▶ PPCA and its extensions
- ▶ Stationarity / non-stationarity
- ▶ Discriminant Analysis (DA)
- ▶ Classification and Clustering
- ▶ Mixture Models
- ▶ Factor Analysis
- ▶ Blind Sources Separation

# Dimension reduction, PCA, Factor Analysis, ICA

- ▶  $M$  variables  $\mathbf{g}_1, \dots, \mathbf{g}_M$  are observed. They are redundant.  
Can we express them with  $N \leq M$  factors  $\mathbf{f}_1, \dots, \mathbf{f}_N$ ?  
**How many factors (Principal Components, Independent Components) can describe the observed data?**
- ▶ Each variable is observed  $T$  times. To index them, we use  $t$ , but this does not necessarily mean time. So, we have  $\{\mathbf{g}_1(t), \dots, \mathbf{g}_M(t), t = 1, \dots, T\}$ .
- ▶ We may represent these data either as a vector or a matrix:

$$\mathbf{g}(t) = \begin{bmatrix} \mathbf{g}_1(t) \\ \mathbf{g}_2(t) \\ \vdots \\ \mathbf{g}_M(t) \end{bmatrix}, t = 1, \dots, T \text{ or } \mathbf{G} = \begin{bmatrix} \mathbf{g}_1(1) & \mathbf{g}_1(1) & \cdots & \mathbf{g}_1(T) \\ \mathbf{g}_2(1) & \mathbf{g}_2(1) & \cdots & \mathbf{g}_2(T) \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{g}_M(1) & \mathbf{g}_M(1) & \cdots & \mathbf{g}_M(T) \end{bmatrix}$$

# Dimension reduction, PCA, Factor Analysis, ICA

- If we define the  $N < M$  factors as

$$\mathbf{f}(t) = \begin{bmatrix} \mathbf{f}_1(t) \\ \mathbf{f}_2(t) \\ \vdots \\ \mathbf{f}_N(t) \end{bmatrix}, t = 1, \dots, T \text{ or } \mathbf{F} = \begin{bmatrix} \mathbf{f}_1(1) & \mathbf{f}_1(1) & \cdots & \mathbf{f}_1(T) \\ \mathbf{f}_2(1) & \mathbf{g}_2(1) & \cdots & \mathbf{f}_1(T) \\ \vdots & & & \\ \mathbf{f}_M(1) & \mathbf{g}_M(1) & \cdots & \mathbf{f}_1(T) \end{bmatrix}$$

where we assume that each factor is a linear combination of data

$$\mathbf{f}_j(t) = \sum_{i=1}^M \mathbf{b}_{ji} \mathbf{g}_i(t) \quad \text{or inversely} \quad \mathbf{g}_i(t) = \sum_{j=1}^N \mathbf{a}_{ij} \mathbf{f}_j(t).$$

- Now, if we define the matrices  $\mathbf{B} = \{\mathbf{b}_{ji}\}$  or  $\mathbf{A} = \{\mathbf{a}_{ij}\}$  we can write

$$\mathbf{f}(t) = \mathbf{B}\mathbf{g}(t) \text{ or } \mathbf{g}(t) = \mathbf{A}\mathbf{f}(t)$$

- $\mathbf{B}$  is called Loading matrix and  $\mathbf{A}$  is called mixing matrix.
- Ideal case is then  $\mathbf{B} = \mathbf{A}^{-1}$ . But this is not interesting, because we want to have  $N < M$ .
- We may accept some errors:  $\mathbf{g}(t) = \mathbf{A}\mathbf{f}(t) + \epsilon(t)$

# Dimension reduction, PCA, Factor Analysis

- ▶  $M$  variables  $\mathbf{g}(t)$  are observed. They are redundant.  
Can we express them with  $N \leq M$  factors  $\mathbf{f}$ ?  
**How many factors (Principal Components, Independent Components) can describe the observed data?**

$$\mathbf{f}(t) = \mathbf{B}\mathbf{g}(t) \text{ or } \mathbf{g}(t) = \mathbf{A}\mathbf{f}(t)$$

or still

$$\mathbf{F} = \mathbf{B}\mathbf{G} \text{ or } \mathbf{G} = \mathbf{A}\mathbf{F}$$

- ▶ We assume all the variables to be centred.
- ▶ PCA uses the second order statistics

$$\text{cov}[\mathbf{g}] = \mathbf{A}\text{cov}[\mathbf{f}]\mathbf{A}^t$$

- ▶ We want principal Components (PC) to be non-correlated:  
 $\text{cov}[\mathbf{f}]$  be a diagonal matrix.

$$\text{cov}[\mathbf{f}] = \text{diag} [\nu_1, \dots, \nu_N]$$

- ▶ One solution: Estimate  $\text{cov}[\mathbf{g}] = \sum_{t=1}^T \mathbf{g}(t)\mathbf{g}'(t)$  and use it.

# Dimension reduction, PCA Algorithms

- ▶ Forward model

$$\mathbf{g}(t) = \mathbf{A}\mathbf{f}(t) \text{ or } \mathbf{G} = \mathbf{A}\mathbf{F}$$

- ▶ Classical PCA algorithm:

- ▶ Estimate  $\text{cov}[\mathbf{g}] = \sum_{t=1}^T \mathbf{g}(t)\mathbf{g}'(t)$
- ▶ Hoping that it is positive Definite, compute its SVD:  
 $\text{cov}[\mathbf{g}] = \mathbf{A}\mathbf{V}\mathbf{A}^t$
- ▶ Identify  $\tilde{\mathbf{A}} = \mathbf{A}\mathbf{V}^{1/2}$  and compute  $\hat{\mathbf{f}} = \mathbf{V}^{-1/2}\mathbf{A}^t\mathbf{g}$
- ▶ The number of factors is the number of non-zero singular values.
- ▶ The data can be retrieved exactly by

$$\hat{\mathbf{g}} = \hat{\mathbf{A}}\hat{\mathbf{f}} = \mathbf{A}\mathbf{V}^{1/2}\mathbf{V}^{-1/2}\mathbf{A}^t\mathbf{g} = \mathbf{A}\mathbf{A}^t\mathbf{g} = \mathbf{g}$$

- ▶ if we keep  $K$  SV, we make an error which can be used as criterion to determine the number of factors

$$E = \frac{\|\hat{\mathbf{g}} - \mathbf{g}\|^2}{\|\mathbf{g}\|^2}$$

# Dimension reduction: PPCA

- ▶ Forward model

$$\mathbf{g}(t) = \mathbf{Af}(t) + \epsilon(t) \text{ or } \mathbf{G} = \mathbf{AF} + \mathbf{E}$$

- ▶ Probabilistic PCA:

$$\text{cov}[\mathbf{g}] = \mathbf{A}\text{cov}[\mathbf{f}]\mathbf{A}^t + \text{cov}[\epsilon]$$

- ▶ Estimate  $\text{cov}[\mathbf{g}] = \sum_{t=1}^T \mathbf{g}(t)\mathbf{g}'(t)$
- ▶ Hoping that it is positive Definite, compute its SVD:  
 $\text{cov}[\mathbf{g}] = \mathbf{AV}\mathbf{A}^t$
- ▶ Keep the SV which are greater than  $v_\epsilon$  and Identify  
 $\widehat{\mathbf{A}} = \mathbf{AV}^{1/2}$  and compute  $\widehat{\mathbf{f}} = \mathbf{V}^{-1/2}\mathbf{A}^t\mathbf{g}$
- ▶ The number of factors is the number of singular values which are greater than  $v_\epsilon$ .
- ▶ The main difficulty is to estimate  $v_\epsilon$ .

# Application on a set of data

- ▶ Genes expressions Time series data in two organs:

## Colon:

Clock: Rev, Per2, Bmal1  
Metabolism: CE2, Top1, UGT, DBP  
CC: Wee1, Ccna2, Ccnb2  
Apoptose: Bcl2, Mdm2, Bax, P53

## Liver:

Clock: Rev, Per2, Bmal1  
Metabolism: CE2, Top1, UGT, DBP  
CC: Wee1, P21  
Apoptose: Bcl2, Mdm2, Bax, P53

## Physiology:

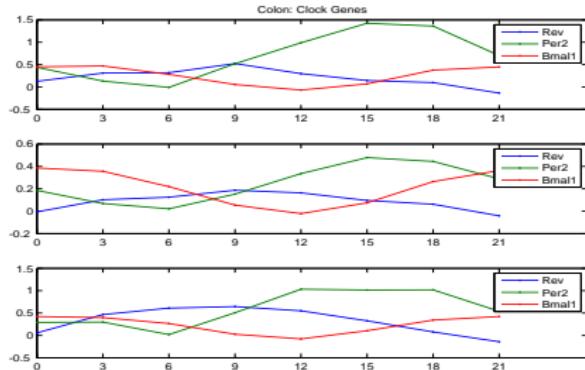
Temperature, Activity  
Cortico, Melato

- ▶ The data are obtained via the COSINOR Model

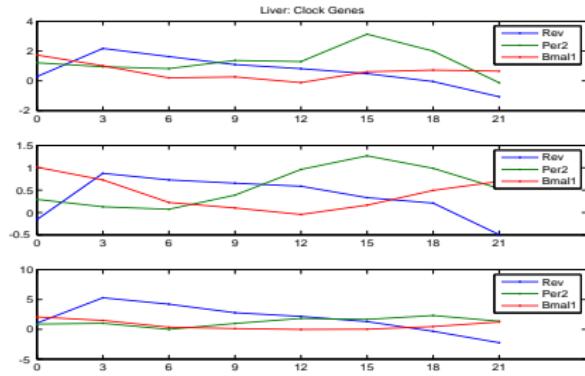
$$g(t) = M + \sum_{k=1}^3 A_k \cos(k\omega_0 t + \phi_k), \omega_0 = \frac{2\pi}{24}, t = 0, 3, 6, \dots, 21$$

# Genes Clock

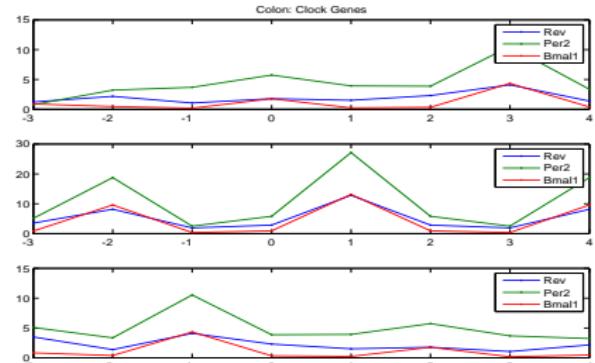
## Colon: Time Series



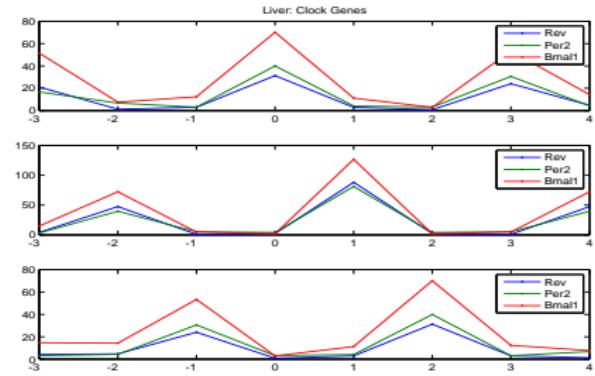
## Liver: Time Series



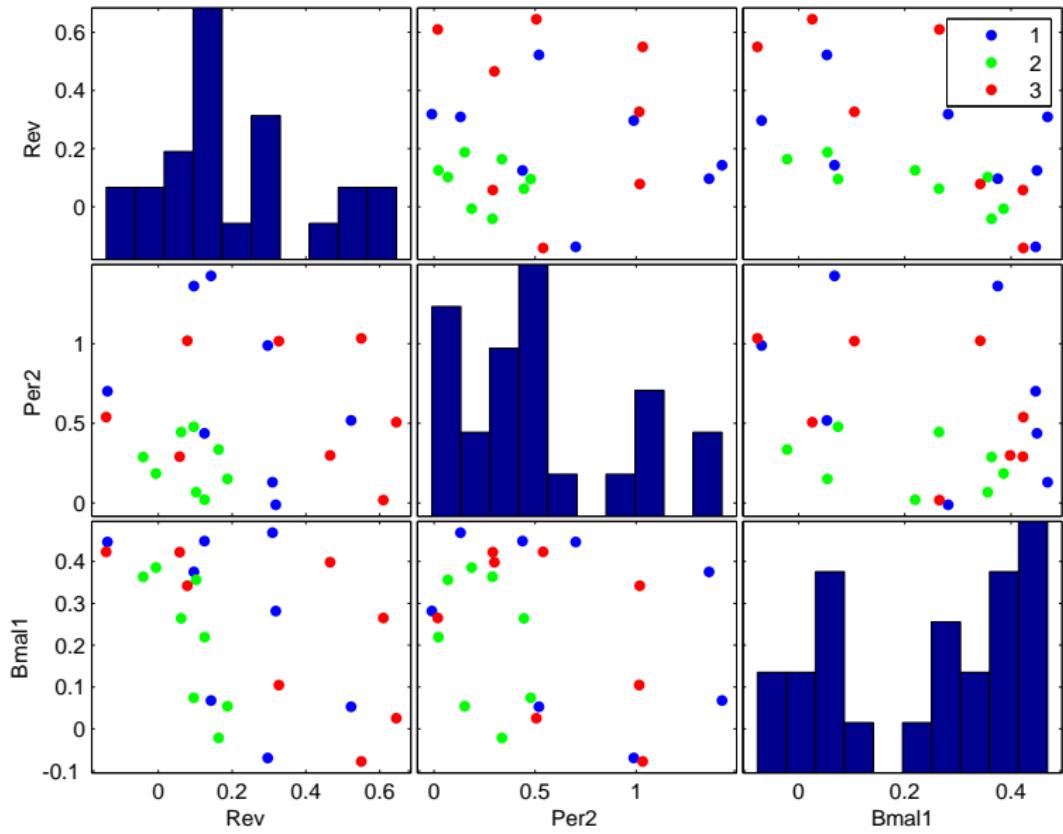
## Fourier Transform



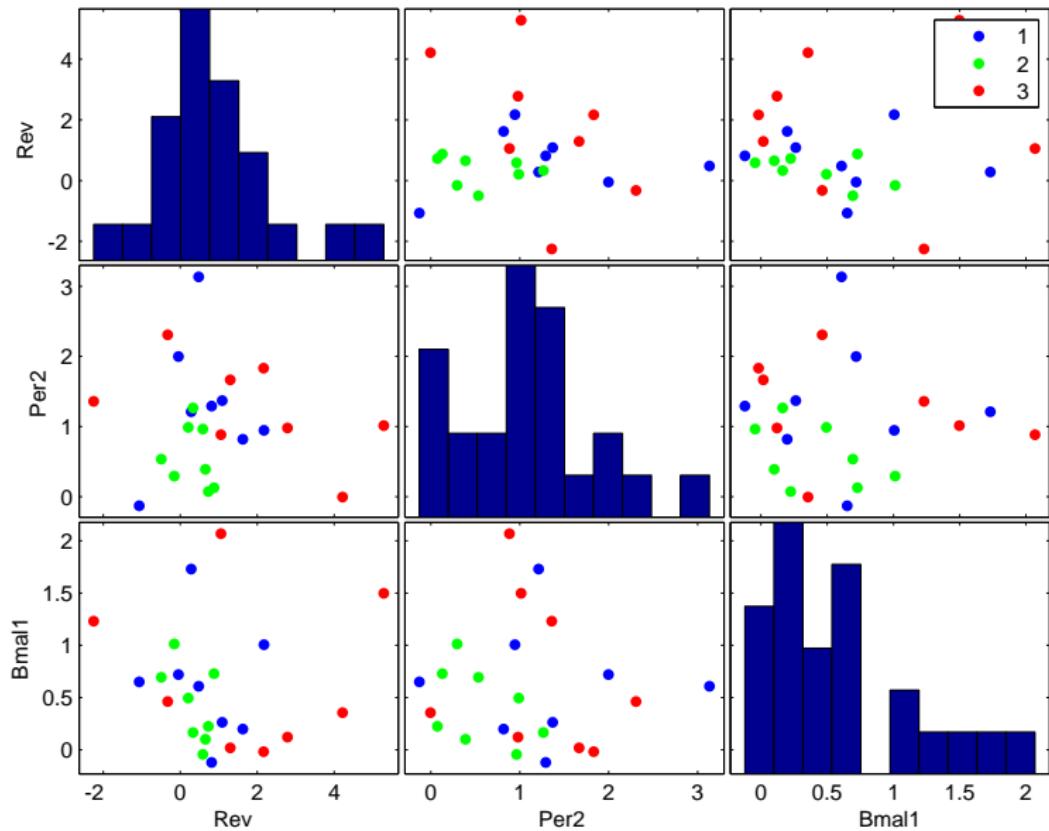
## Fourier Transform



# Genes Colon Time series: Clock

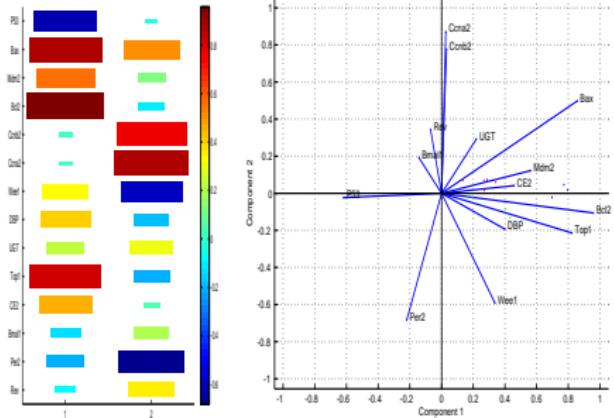


# Genes Liver Time series: Clock

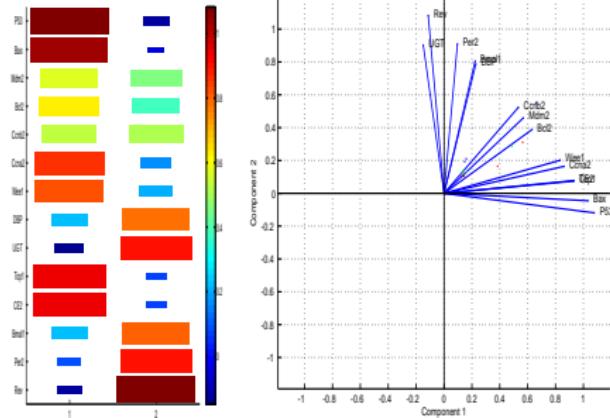


# Factor Analysis: 2 factors: Colon

Time series



FT Amplitudes



# How to determine the number of factors

- ▶ When  $N$  is given:

$$p(\mathbf{A}, \mathbf{f}|\mathbf{g}) \propto p(\mathbf{g}|\mathbf{A}, \mathbf{f}) p(\mathbf{A}) p(\mathbf{f})$$

Different choices for  $p(\mathbf{A})$  and  $p(\mathbf{f})$  and

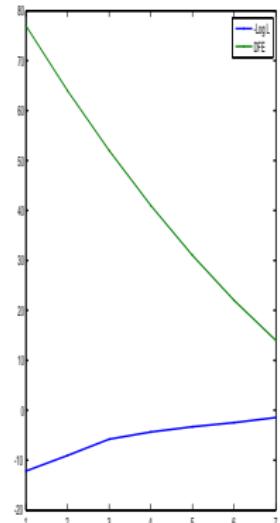
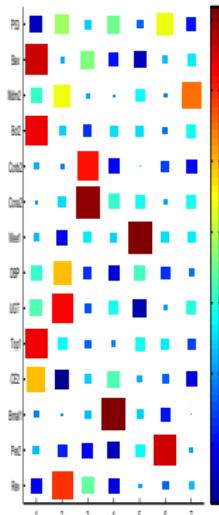
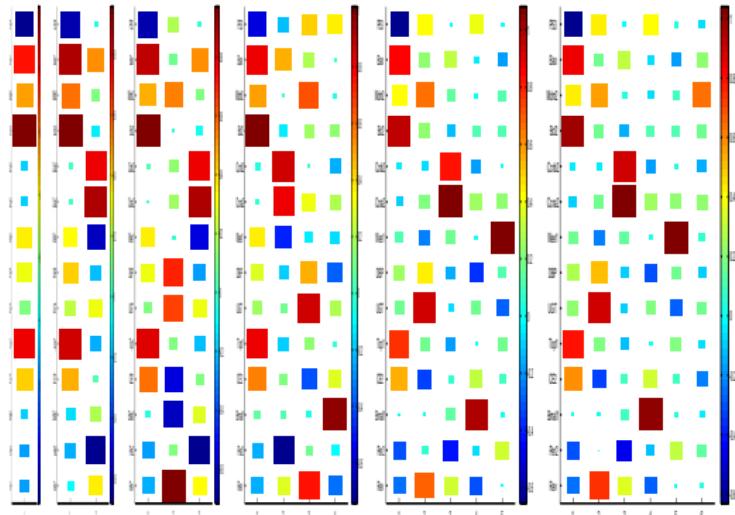
Different methods to estimate both  $\mathbf{A}$  and  $\mathbf{f}$ :

JMAP, EM, Variational Bayesian Approximation

- ▶ When  $N$  is not known:

- ▶ Model selection
- ▶ Bayesian or Maximum likelihood methods
- ▶ To determine the number of factors we do the analyze with different  $N$  factors and use two criteria:
  - ▶ -log likelihood –  $\ln p(\mathbf{g}|\mathbf{A}, N)$  of the observations and
  - ▶ DFE: Degrees of freedom error  $(N - M)^2 - (N + M))/2$  related to AIC or BIC model selection criteria.

# Factor Analysis: Time series, colon



# Dimension reduction: ML PPCA

- ▶ Forward model

$$\mathbf{g}(t) = \mathbf{A}\mathbf{f}(t) + \boldsymbol{\epsilon}(t) \text{ or } \mathbf{G} = \mathbf{A}\mathbf{F} + \mathbf{E}$$

- ▶ ML PPCA:

$$\boldsymbol{\Sigma}_g = \text{cov}[\mathbf{g}] = \mathbf{A}\text{cov}[\mathbf{f}]\mathbf{A}^t + \mathbf{v}_\epsilon \mathbf{I} = \mathbf{A}\mathbf{V}_f\mathbf{A}^t + \mathbf{v}_\epsilon \mathbf{I}$$

$$\boldsymbol{\Sigma}_g = \mathbf{A}\text{diag}[\mathbf{v}_f]\mathbf{A}^t + \mathbf{v}_\epsilon \mathbf{I}$$

- ▶ Likelihood

$$p(\mathbf{g}|\mathbf{A}, \mathbf{v}_f, \mathbf{v}_\epsilon) = \mathcal{N}(\mathbf{g}|0, \boldsymbol{\Sigma}_g)$$

$$p(\mathbf{g}|\mathbf{A}, \mathbf{v}_f, \mathbf{v}_\epsilon) \propto \det(\boldsymbol{\Sigma}_g)^{-1/2} \exp \left[ -\frac{1}{2\mathbf{v}_\epsilon} \sum_t \|\mathbf{g}(t) - \mathbf{A}\mathbf{f}(t)\|^2 \right]$$

- ▶ Alternate maximization of the -log likelihood

$$\mathcal{L}(\mathbf{A}, \mathbf{v}_f, \mathbf{v}_\epsilon) = \frac{1}{2} \ln \det(\mathbf{A}\text{diag}[\mathbf{v}_f]\mathbf{A}^t + \mathbf{v}_\epsilon \mathbf{I}) + \frac{1}{2\mathbf{v}_\epsilon} \sum_t \|\mathbf{g}(t) - \mathbf{A}\mathbf{f}(t)\|^2$$

with respect to its arguments can give the desired solution.

# Dimension reduction, PCA, Factor Analysis, ICA

- ▶  $M$  variables  $\mathbf{g}(t)$  are observed. They are redundant.  
Can we express them with  $N \leq M$  factors  $\mathbf{f}$  ?  
**How many factors (Principal Components, Independent Components) can describe the observed data?**

$$\begin{cases} \mathbf{g}_i(t) = \sum_{j=1}^N a_{ij} \mathbf{f}_j(t) + \epsilon_i(t) \\ \mathbf{g}(t) = \mathbf{A}\mathbf{f}(t) + \boldsymbol{\epsilon}(t) \end{cases} \quad \begin{cases} \mathbf{A} : (M \times N) \text{ Loading matrix , } N \leq M \\ \mathbf{f}(t) : \text{ factors, sources} \end{cases}$$

- ▶ How to find both  $\mathbf{A}$  and factors  $\mathbf{f}(t)$  ?

Three cases:

- ▶  $\mathbf{A}$  known, find  $\mathbf{f}$
- ▶  $\mathbf{f}$  known, find  $\mathbf{A}$
- ▶  $\mathbf{A}$  and  $\mathbf{f}$  both unknown

## Case 1: $A$ known, find $\hat{\mathbf{f}}$

$$\mathbf{g}(t) = \mathbf{A}\mathbf{f}(t) + \epsilon(t)$$

- ▶ First, assume the data iid. So  $\mathbf{g} = \mathbf{A}\mathbf{f} + \epsilon$ .
- ▶ Second, assume no noise. So  $\mathbf{g} = \mathbf{A}\mathbf{f}$
- ▶ Ideal case  $M = N$  and  $A$  invertible:  $\hat{\mathbf{f}} = \mathbf{A}^{-1}\mathbf{g}$
- ▶  $M < N$ :  $\mathbf{A}\mathbf{f} = \mathbf{g}$  has infinite number of solutions.  
Minimum Norme (MN) Solution:

$$\hat{\mathbf{f}} = \arg \min_{\{\mathbf{f}: \mathbf{A}\mathbf{f} = \mathbf{g}\}} \left\{ \|\mathbf{f}\|^2 \right\}$$

Lagrangian technique:  $\mathbf{A}\mathbf{A}^t\mathbf{f} = \mathbf{A}^t\mathbf{g}$

If  $A$  has full rank:  $\text{rank } (\mathbf{A}) = \min\{M, N\}$

$$\hat{\mathbf{f}} = \mathbf{A}^t[\mathbf{A}\mathbf{A}^t]^{-1}\mathbf{g}$$

## Case 1: $A$ known, find $\hat{f}$ : LS

- ▶  $M > N$ :  $A\hat{f} = \mathbf{g}$  may not have any solution.  
Least Square (LS) Solution:

$$\hat{\mathbf{f}} = \arg \min_{\mathbf{f}} \left\{ \|\mathbf{g} - A\mathbf{f}\|^2 \right\}$$

Gradient of  $J(\mathbf{f}) = \|\mathbf{g} - A\mathbf{f}\|^2$  is :

$$-2A^t(\mathbf{g} - A\mathbf{f}) = 0 \longrightarrow A^t A \mathbf{f} = A^t \mathbf{g}.$$

- ▶ If  $A$  has full rank:  $\text{rank}(A) = \min\{M, N\}$ :

$$\hat{\mathbf{f}} = [A^t A]^{-1} A^t \mathbf{g}$$

- ▶ General case ?
  - ▶ Truncated Singular Value Decomposition (TSVD)
  - ▶ Regularization
  - ▶ Bayesian

## Case 1: $A$ known, find $\hat{f}$ : TSVD

- ▶ SVD:  $A = U \Lambda V^t$ ,  $UU^t = I$ ,  $V^t V = I$ ,  $\Lambda$  diagonal.
- ▶ Right SVD:  $AA^t v_k = \lambda_k v_k$
- ▶ Left SVD:  $A^t A u_k = \lambda_k u_k$
- ▶ Full rank

$$\hat{f} = VS^{-1}U^t g = \sum_k v_k \frac{\langle g, u_k \rangle}{s_k}$$

- ▶ Rank deficient rank ( $A$ ) =  $K < \min\{M, N\}$

$$\hat{f} = VS^+U^t g = \sum_k^K v_k \frac{\langle g, u_k \rangle}{s_k}$$

- ▶ Main difficulty: How to decide which singular values are near to zero.

## Case 1: $\mathbf{A}$ known, find $\mathbf{f}$ : Regularization

- ▶ MN: Minimize  $\|\mathbf{f}\|^2$  subject to  $\mathbf{A}\mathbf{f} = \mathbf{g}$ ,
- ▶ LS: Minimize  $\|\mathbf{g} - \mathbf{A}\mathbf{f}\|^2$
- ▶ MN + LS:

$$\widehat{\mathbf{f}} = \arg \min_{\mathbf{f}} \{J(\mathbf{f})\}$$

$$\text{with } J(\mathbf{f}) = \|\mathbf{g} - \mathbf{A}\mathbf{f}\|^2 + \lambda \|\mathbf{f}\|^2$$

- ▶ Gradient of  $J(\mathbf{f})$  is :

$$-2\mathbf{A}^t(\mathbf{g} - \mathbf{A}\mathbf{f}) - 2\lambda\mathbf{I} = \mathbf{0} \longrightarrow (\mathbf{A}^t\mathbf{A} + \lambda\mathbf{I})\mathbf{f} = \mathbf{A}^t\mathbf{g}$$

$$\widehat{\mathbf{f}} = [\mathbf{A}^t\mathbf{A} + \lambda\mathbf{I}]^{-1}\mathbf{A}^t\mathbf{g}$$

- ▶  $\lambda \rightarrow 0$  LS.
- ▶ When  $M < N$ , we may also obtain

$$\widehat{\mathbf{f}} = \mathbf{A}^t[\mu\mathbf{A}\mathbf{A}^t + \mathbf{I}]^{-1}\mathbf{g}, \text{ with } \mu = \frac{1}{\lambda}.$$

## Case 1: $A$ known, find $f$ : Bayesian inference

$$\mathbf{g} = A\mathbf{f} + \boldsymbol{\epsilon}$$

- ▶ Prior knowledge on  $\boldsymbol{\epsilon}$ :

$$\boldsymbol{\epsilon} \sim \mathcal{N}(\boldsymbol{\epsilon}|0, v_\epsilon \mathbf{I}) \longrightarrow p(\mathbf{g}|\mathbf{f}, A) = \mathcal{N}(\mathbf{g}|A\mathbf{f}, v_\epsilon \mathbf{I}) \propto \exp\left[\frac{1}{2v_\epsilon} \|\mathbf{g} - A\mathbf{f}\|^2\right]$$

- ▶ Simple prior models for  $\mathbf{f}$ :  $p(\mathbf{f}|\alpha) \propto \exp\left[-\frac{1}{2v_f} \|\mathbf{f}\|^2\right]$
- ▶ Expression of the posterior law:

$$p(\mathbf{f}|\mathbf{g}, A) \propto p(\mathbf{g}|\mathbf{f}, A) p(\mathbf{f}) \propto \exp\left[-\frac{1}{2v_\epsilon} J(\mathbf{f})\right]$$

$$\text{with } J(\mathbf{f}) = \|\mathbf{g} - A\mathbf{f}\|^2 + \lambda \|\mathbf{f}\|^2, \quad \lambda = v_\epsilon/v_f$$

- ▶ Link between MAP estimation and regularization

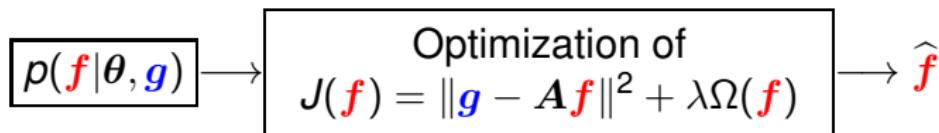
$$\hat{\mathbf{f}} = \arg \max_{\mathbf{f}} \{p(\mathbf{f}|\mathbf{g}, A)\} = \arg \min_{\mathbf{f}} \{J(\mathbf{f})\}$$

- ▶ Solution:  $\hat{\mathbf{f}} = (A'A + \lambda I)^{-1} A' \mathbf{g}$

## Bayesian inference for sources $\mathbf{f}$ when $\mathbf{A}$ is known

- More general prior model  $p(\mathbf{f}) \propto \exp[-\frac{\alpha}{2}\Omega(\mathbf{f})]$
- MAP:

$$J(\mathbf{f}) = \|\mathbf{g} - \mathbf{A}\mathbf{f}\|^2 + \lambda\Omega(\mathbf{f}), \quad \lambda = v_\epsilon\alpha$$



- Different priors=Different expressions for  $\Omega(\mathbf{f})$
- Solution can be obtained using appropriate optimization algorithm.

## MAP estimation with sparsity enforcing priors

- ▶ Gaussian:  $\Omega(\mathbf{f}) = \|\mathbf{f}\|^2 = \sum_j |f_j|^2$   
$$J(\mathbf{f}) = \|\mathbf{g} - \mathbf{A}\mathbf{f}\|^2 + \lambda \|\mathbf{f}\|^2 \longrightarrow \hat{\mathbf{f}} = [\mathbf{A}'\mathbf{A} + \lambda \mathbf{I}]^{-1} \mathbf{A}' \mathbf{g}$$

- ▶ Generalized Gaussian:

$$\Omega(\mathbf{f}) = \gamma \sum_j |\mathbf{f}_j|^\beta$$

- ▶ Student-t model:

$$\Omega(\mathbf{f}) = \frac{\nu + 1}{2} \sum_j \log \left( 1 + \mathbf{f}_j^2 / \nu \right)$$

- ▶ Elastic Net model:

$$\Omega(\mathbf{f}) = \sum_j \left[ \gamma_1 |\mathbf{f}_j| + \gamma_2 \mathbf{f}_j^2 \right]$$

For an extended list of such sparsity enforcing priors see:

A. Mohammad-Djafari, "Bayesian approach with prior models which enforce sparsity in signal and image processing," EURASIP Journal on Advances in Signal Processing, vol. Special issue on Sparse Signal Processing, 2012.

## Case 1: $\mathbf{A}$ known, find $\mathbf{f}$ : Regularization

- More general case:

$$\widehat{\mathbf{f}} = \arg \min_{\mathbf{f}} \{J(\mathbf{f})\}$$

with  $J(\mathbf{f}) = \|\mathbf{g} - \mathbf{A}\mathbf{f}\|^2 + \lambda \|\mathbf{D}\mathbf{f}\|^2$

- Gradient of  $J(\mathbf{f})$  is :

$$-2\mathbf{A}^t(\mathbf{g} - \mathbf{A}\mathbf{f}) - 2\mathbf{D}^t\mathbf{D} = 0 \longrightarrow (\mathbf{A}^t\mathbf{A} + \lambda\mathbf{D}^t\mathbf{D})\mathbf{f} = \mathbf{A}^t\mathbf{g}$$

$$\widehat{\mathbf{f}} = [\mathbf{A}^t\mathbf{A} + \lambda\mathbf{D}^t\mathbf{D}]^{-1}\mathbf{A}^t\mathbf{g}$$

- L1 Regularization:

$$J(\mathbf{f}) = \|\mathbf{g} - \mathbf{A}\mathbf{f}\|_2^2 + \lambda \|\mathbf{f}\|_1$$

with  $\|\mathbf{f}\|_1 = \sum_j |\mathbf{f}_j|$ .

- Still more general

$$J(\mathbf{f}) = \Delta_1(\mathbf{g}, \mathbf{A}\mathbf{f}) + \lambda \Delta_2(\mathbf{f}, \mathbf{f}_0))$$

## Case2: $f$ known, find $A$

- $\mathbf{g} = \mathbf{A}\mathbf{f} + \epsilon$  is linear in  $f$  and in  $A$ .

$$\mathbf{g} = \mathbf{A}\mathbf{f} \rightarrow g_i = \sum_j a_{ij} f_j \rightarrow \mathbf{g} = \sum_j a_{*j} f_j$$

$$\begin{bmatrix} g_1 \\ g_2 \end{bmatrix} = \begin{bmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{bmatrix} \begin{bmatrix} f_1 \\ f_2 \end{bmatrix} = \begin{bmatrix} f_1 & 0 & f_2 & 0 \\ 0 & f_1 & 0 & f_2 \end{bmatrix} \begin{bmatrix} a_{11} \\ a_{21} \\ a_{12} \\ a_{22} \end{bmatrix}$$
$$\mathbf{g} = \mathbf{A}\mathbf{f} = \mathbf{F}\mathbf{a} \quad \text{with} \quad \mathbf{F} = \mathbf{f} \odot \mathbf{I}, \quad \mathbf{a} = \text{vec}(\mathbf{A})$$

- $A$  known, estimation of  $f$ :  $\mathbf{g} = \mathbf{A}\mathbf{f} + \epsilon$
- $f$  known, estimation of  $A$ :  $\mathbf{g} = \mathbf{F}\mathbf{a} + \epsilon$
- Joint estimation of  $f$  and  $A$ :  $\mathbf{g} = \mathbf{A}\mathbf{f} + \epsilon = \mathbf{F}\mathbf{a} + \epsilon$

## Cases 1 and 2: Regularization

$$\mathbf{g} = \mathbf{A}\mathbf{f} \quad \text{or} \quad \mathbf{g} = \mathbf{A}\mathbf{f} = \mathbf{F}\mathbf{a}$$

- $\mathbf{A}$  known, find  $\mathbf{f}$

$$\widehat{\mathbf{f}} = \arg \min_{\mathbf{f}} \left\{ \|\mathbf{g} - \mathbf{A}\mathbf{f}\|^2 + \lambda \|\mathbf{f}\|^2 \right\} = (\mathbf{A}'\mathbf{A} + \lambda \mathbf{I})^{-1} \mathbf{A}'\mathbf{g}$$

- $\mathbf{f}$  known, find  $\mathbf{A}$

$$\widehat{\mathbf{A}} = \arg \min_{\mathbf{A}} \left\{ \|\mathbf{g} - \mathbf{A}\mathbf{f}\|^2 + \lambda \|\mathbf{A}\|^2 \right\} = \mathbf{g}\mathbf{f}'(\mathbf{f}\mathbf{f}' + \lambda \mathbf{I})^{-1}$$

- Both  $\mathbf{A}$  and  $\mathbf{f}$  are unknown

$$(\widehat{\mathbf{f}}, \widehat{\mathbf{A}}) = \arg \min_{(\mathbf{f}, \mathbf{A})} \left\{ \|\mathbf{g} - \mathbf{A}\mathbf{f}\|^2 + \lambda_1 \|\mathbf{f}\|^2 + \lambda_2 \|\mathbf{A}\|^2 \right\}$$

Alternate optimisation ?

# Estimation of $\mathbf{A}$ when the sources $\mathbf{f}$ are known

Source separation is a bilinear model:

$$\mathbf{g} = \mathbf{Af} = \mathbf{Fa}$$
$$\begin{bmatrix} g_1 \\ g_2 \end{bmatrix} = \begin{bmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{bmatrix} \begin{bmatrix} f_1 \\ f_2 \end{bmatrix} = \begin{bmatrix} f_1 & 0 & f_2 & 0 \\ 0 & f_1 & 0 & f_2 \end{bmatrix} \begin{bmatrix} a_{11} \\ a_{21} \\ a_{12} \\ a_{22} \end{bmatrix}$$
$$\mathbf{F} = \mathbf{f} \odot \mathbf{I}, \quad \mathbf{a} = \text{vec}(\mathbf{A})$$

- ▶ Problem is more ill-posed (underdetermined).
- ▶ We need absolutely to impose constraints on elements or the structure of  $\mathbf{A}$ , for example:
  - ▶ Positivity of the elements
  - ▶ Toeplitz or TBT structure
  - ▶ Symmetry
  - ▶ Sparsity
- ▶ The same Bayesian approach then can be applied.

# Estimation of $\mathbf{A}$ when the sources $f$ are known

$$\mathbf{g} = \mathbf{Af} + \epsilon = \mathbf{Fa} + \epsilon$$

- ▶ Prior on noise:

$$\begin{aligned} p(\mathbf{g}|f, \mathbf{A}) &= \mathcal{N}(\mathbf{g}|\mathbf{Af}, v_\epsilon \mathbf{I}) \propto \exp\left[-\frac{1}{2v_\epsilon} \|\mathbf{g} - \mathbf{Af}\|^2\right] \\ &\propto \exp\left[-\frac{1}{2v_\epsilon} \|\mathbf{g} - \mathbf{Fa}\|^2\right] \end{aligned}$$

- ▶ Simple prior models for  $a$ :

$$p(\mathbf{A}|\alpha) \propto \exp\left[-\alpha\|\mathbf{a}\|^2\right] \propto \exp\left[-\alpha\|\mathbf{A}\|^2\right]$$

- ▶ Expression of the posterior law:

$$p(\mathbf{A}|\mathbf{g}, f) \propto p(\mathbf{g}|f, \mathbf{A}) p(\mathbf{A}) \propto \exp[-J(\mathbf{A})]$$

$$\text{with } J(\mathbf{A}) = \frac{1}{2v_\epsilon} \|\mathbf{g} - \mathbf{Af}\|^2 + \alpha \|\mathbf{A}\|^2$$

- ▶ MAP estimation:

$$\hat{\mathbf{a}} = (\mathbf{F}'\mathbf{F} + \lambda \mathbf{I})^{-1} \mathbf{F}'\mathbf{g} \leftrightarrow \hat{\mathbf{A}} = \mathbf{g}\mathbf{f}'(\mathbf{f}\mathbf{f}' + \lambda \mathbf{I})^{-1}$$

$$\mathbf{g}(t) = \mathbf{Af}(t) + \epsilon(t) \longrightarrow \hat{\mathbf{A}} = \left( \sum_t \mathbf{g}(t)\mathbf{f}'(t) \right) \left( \sum_t \mathbf{f}(t)\mathbf{f}'(t) + \lambda \mathbf{I} \right)^{-1}$$

## Case 3: both $\mathbf{f}$ and $\mathbf{A}$ unknown

- Both  $\mathbf{A}$  and  $\mathbf{f}$  are unknown:

$$(\hat{\mathbf{f}}, \hat{\mathbf{A}}) = \arg \min_{(\mathbf{f}, \mathbf{A})} \left\{ \|\mathbf{g} - \mathbf{A}\mathbf{f}\|^2 + \lambda_1 \|\mathbf{f}\|^2 + \lambda_2 \|\mathbf{A}\|^2 \right\}$$

- Undeterminations:
  - Permutation:  $\mathbf{AP}, \mathbf{P}'\mathbf{f}$
  - Scale:  $k\mathbf{A}, \frac{1}{k}\mathbf{f}$
- Alternate optimisation

$$\begin{cases} \hat{\mathbf{f}} = \arg \min_{\mathbf{f}} \left\{ \|\mathbf{g} - \mathbf{A}\mathbf{f}\|^2 + \lambda_1 \|\mathbf{f}\|^2 \right\} = (\mathbf{A}'\mathbf{A} + \lambda_1 \mathbf{I})^{-1} \mathbf{A}'\mathbf{g} \\ \hat{\mathbf{A}} = \arg \min_{\mathbf{A}} \left\{ \|\mathbf{g} - \mathbf{A}\mathbf{f}\|^2 + \lambda_2 \|\mathbf{A}\|^2 \right\} = \mathbf{g}\mathbf{f}'(\mathbf{f}\mathbf{f}' + \lambda_2 \mathbf{I}) \end{cases}$$

- Importance of initialization and other constraints such as positivity
  - Non-negative Matrix decomposition

# Both $\mathbf{A}$ and $\mathbf{f}$ unknown: Bayesian approach

Three main steps:

1. Assigning priors:  $p(\mathbf{f})$  and  $p(\mathbf{A})$
2. Obtaining the expressions of the joint posterior:

$$p(\mathbf{f}, \mathbf{A} | \mathbf{g}) \propto p(\mathbf{g} | \mathbf{f}, \mathbf{A}) p(\mathbf{f}) p(\mathbf{A})$$

3. Doing the computations:

- Joint optimization of  $p(\mathbf{f}, \mathbf{A} | \mathbf{g})$ ;
- MCMC Gibbs sampling methods which need generation of samples from the conditionals  $p(\mathbf{f} | \mathbf{A}, \mathbf{g})$  and  $p(\mathbf{A} | \mathbf{f}, \mathbf{g})$ ;
- Marginalisation and EM algorithm:

$$p(\mathbf{A} | \mathbf{g}) = \int p(\mathbf{f}, \mathbf{A} | \mathbf{g}) d\mathbf{f} \longrightarrow \widehat{\mathbf{A}} \longrightarrow p(\mathbf{f} | \widehat{\mathbf{A}}, \mathbf{g}) \longrightarrow \widehat{\mathbf{f}}$$

- Bayesian Variational Approximation (BVA) methods which approximate  $p(\mathbf{f}, \mathbf{A} | \mathbf{g})$  by a separable one  $q(\mathbf{f}, \mathbf{A}) = q_1(\mathbf{f})q_2(\mathbf{A})$  and then using them for the estimation of  $\mathbf{f}$  and  $\mathbf{A}$ .

# Bayesian source separation: both $\mathbf{A}$ and $\mathbf{f}$ unknown

$$p(\mathbf{f}, \mathbf{A} | \mathbf{g}, \theta_1, \theta_2, \theta_3) = \frac{p(\mathbf{g} | \mathbf{f}, \mathbf{A}, \theta_1) p(\mathbf{f} | \theta_2) p(\mathbf{A} | \theta_3)}{p(\mathbf{g} | \theta_1, \theta_2, \theta_3)}$$

- ▶ Joint estimation (JMAP):

$$(\hat{\mathbf{f}}, \hat{\mathbf{A}}) = \arg \max_{(\mathbf{f}, \mathbf{A})} \{p(\mathbf{f}, \mathbf{A} | \mathbf{g}, \theta_1, \theta_2, \theta_3)\}$$

- ▶ JMAP with Gaussian priors:

$$(\hat{\mathbf{f}}, \hat{\mathbf{A}}) = \arg \min_{(\mathbf{f}, \mathbf{A})} \left\{ \|\mathbf{g} - \mathbf{Af}\|^2 + \lambda_1 \|\mathbf{f}\|^2 + \lambda_2 \|\mathbf{A}\|^2 \right\}$$

- ▶ Permutation and scale indeterminations:  
needs good choices for priors
- ▶ Alternate optimisation

$$\begin{cases} \hat{\mathbf{f}} = \arg \min_{\mathbf{f}} \{ \|\mathbf{g} - \mathbf{Af}\|^2 + \lambda_1 \|\mathbf{f}\|^2 \} = (\mathbf{A}'\mathbf{A} + \lambda_1 \mathbf{I})^{-1} \mathbf{A}'\mathbf{g} \\ \hat{\mathbf{A}} = \arg \min_{\mathbf{A}} \{ \|\mathbf{g} - \mathbf{Af}\|^2 + \lambda_2 \|\mathbf{A}\|^2 \} = \mathbf{g}\mathbf{f}'(\mathbf{f}\mathbf{f}' + \lambda_2 \mathbf{I}) \end{cases}$$

- ▶ Importance of initialization and constraints such as positivity

# Joint MAP Estimation of $\mathbf{A}$ and $\mathbf{f}$ with Gaussian priors

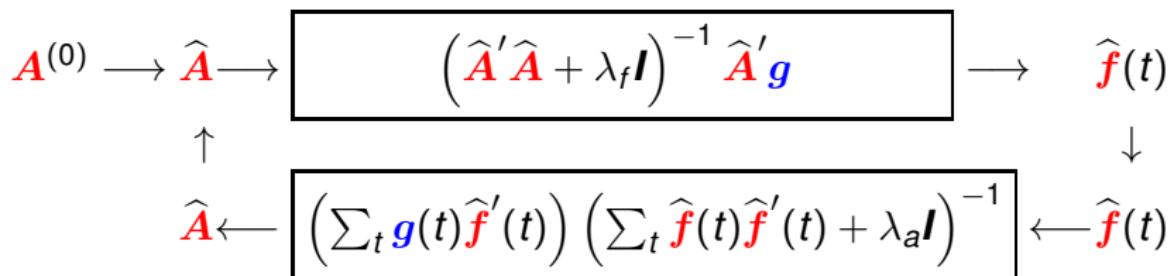
$$\mathbf{g}(t) = \mathbf{A} \mathbf{f}(t) + \epsilon(t), \quad t = 1, \dots, T, \quad iid$$

$$\begin{aligned} p(\mathbf{f}_{1..T}, \mathbf{A} | \mathbf{g}_{1..T}) &\propto p(\mathbf{g}_{1..T} | \mathbf{A}, \mathbf{f}_{1..T}, v_\epsilon) p(\mathbf{f}_{1..T}) p(\mathbf{A} | \mathbf{A}_0, \mathbf{V}_0) \\ &\propto \prod_t p(\mathbf{g}(t) | \mathbf{A}, \mathbf{f}(t), v_\epsilon) p(\mathbf{f}(t) | \mathbf{z}(t)) p(\mathbf{A} | \mathbf{A}_0, \mathbf{V}_0) \end{aligned}$$

Joint MAP: Alternate optimization

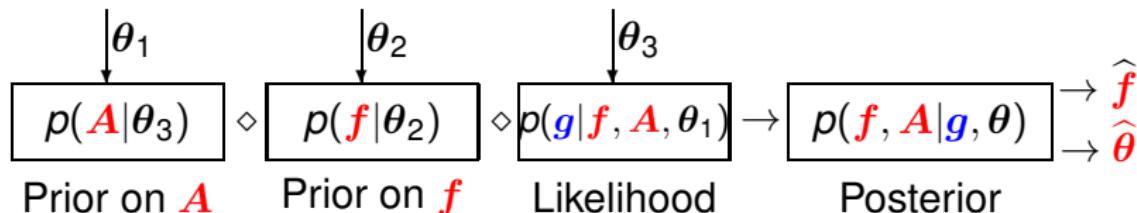
$$\left\{ \begin{array}{l} \widehat{\mathbf{f}}(t) = (\widehat{\mathbf{A}}' \widehat{\mathbf{A}} + \lambda_f \mathbf{I})^{-1} \widehat{\mathbf{A}}' \mathbf{g}(t), \quad \lambda_f = v_\epsilon / v_f \\ \widehat{\mathbf{A}} = \sum_t \mathbf{g}(t) \widehat{\mathbf{f}}'(t) \left( \sum_t \widehat{\mathbf{f}}(t) \widehat{\mathbf{f}}'(t) + \lambda_a \mathbf{I} \right)^{-1} \quad \lambda_a = v_\epsilon / v_a \end{array} \right.$$

Alternate optimization Algorithm:



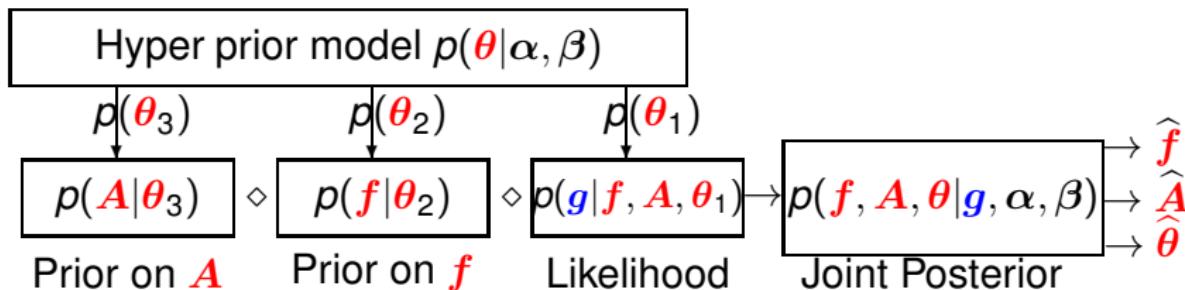
# Summary of Bayesian estimation with different levels

- ▶ Simple Bayesian Model and Estimation



- ▶ Full Bayesian Model and Hyperparameter Estimation scheme

$\downarrow \alpha, \beta$



# Summary of Bayesian estimation with different levels

- ▶ Marginalization over  $\mathbf{f}$



Joint Posterior Marginalize over  $\mathbf{f}$

- ▶ Marginalization over  $\mathbf{A}$



Joint Posterior Marginalize over  $\mathbf{A}$

- ▶ Joint MAP

$$p(\mathbf{f}, \mathbf{A} | \mathbf{g}) \rightarrow \begin{cases} \tilde{\mathbf{f}} = \arg \max_{\mathbf{f}} \{p(\mathbf{f}, \tilde{\mathbf{A}} | \mathbf{g})\} \\ \tilde{\mathbf{A}} = \arg \max_{\mathbf{A}} \{p(\tilde{\mathbf{f}}, \mathbf{A} | \mathbf{g})\} \end{cases} \rightarrow \hat{\mathbf{A}} \rightarrow \hat{\mathbf{f}}$$

Joint Posterior Alternate optimization

- ▶ Other solutions:

- ▶ MCMC for exploring the joint posterior law and computing mean values
- ▶ Variational Bayesian Approximation (VBA)

# Variational Bayesian Approximation

- ▶ Main idea: Approximate the joint pdf  $p(\mathbf{A}, \mathbf{f}|\mathbf{g})$  difficult to handle by a simpler one. For example a separable one  $q(\mathbf{A}, \mathbf{f}|\mathbf{g}) = q_1(\mathbf{A}) q_2(\mathbf{f})$  or even  $q_1(\mathbf{A}) \prod_j q_2(\mathbf{f}_j)$ .
- ▶ Criterion: minimize

$$\text{KL}(q|p) = \int q \ln \frac{q}{p} = -H(q_1) - H(q_2) - \langle \ln p(\mathbf{A}, \mathbf{f}|\mathbf{g}) \rangle_q$$

- ▶ Solution obtained by alternate optimization:

$$\begin{cases} q_1(\mathbf{f}) \propto \exp \left[ \langle \ln p(\mathbf{f}, \mathbf{A}|\mathbf{g}) \rangle_{q_2(\mathbf{A})} \right] \\ q_2(\mathbf{A}) \propto \exp \left[ \langle \ln p(\mathbf{f}, \mathbf{A}|\mathbf{g}) \rangle_{q_1(\mathbf{f})} \right] \end{cases}$$

- ▶ Use  $q_1(\mathbf{f})$  for inferring on  $\mathbf{f}$  and  $q_2(\mathbf{A})$  for inferring on  $\mathbf{A}$

# VBA and links with JMAP and EM

Three possibilities:

- $q_1(\mathbf{f}) = \delta(\mathbf{f} - \tilde{\mathbf{f}})$  and  $q_2(\mathbf{A}) = \delta(\mathbf{A} - \tilde{\mathbf{A}}) \rightarrow \text{JMAP}$

$$\begin{cases} q_1(\mathbf{f}) \propto p(\mathbf{f}, \tilde{\mathbf{A}} | \mathbf{g}) \\ q_2(\mathbf{A}) \propto p(\tilde{\mathbf{f}}, \mathbf{A} | \mathbf{g}) \end{cases} \rightarrow \begin{cases} \tilde{\mathbf{f}} = \arg \max_{\mathbf{f}} \{ p(\mathbf{f}, \tilde{\mathbf{A}} | \mathbf{g}) \} \\ \tilde{\mathbf{A}} = \arg \max_{\mathbf{A}} \{ p(\tilde{\mathbf{f}}, \mathbf{A} | \mathbf{g}) \} \end{cases}$$

- $q_1(\mathbf{f})$  free and  $q_2(\mathbf{A}) = \delta(\mathbf{A} - \tilde{\mathbf{A}}) \rightarrow \text{EM}$

$$\begin{cases} q_1(\mathbf{f}) \propto \exp \left[ \langle \ln p(\mathbf{f}, \mathbf{A} | \mathbf{g}) \rangle_{q_2(\mathbf{A})} \right] \propto p(\mathbf{f}, \tilde{\mathbf{A}} | \mathbf{g}) = p(\mathbf{f} | \tilde{\mathbf{A}}, \mathbf{g}) \\ q_2(\mathbf{A}) \propto \exp \left[ \langle \ln p(\mathbf{f}, \mathbf{A} | \mathbf{g}) \rangle_{q_1(\mathbf{f})} \right] \propto \exp \left[ Q(\mathbf{A}, \tilde{\mathbf{A}}) \right] \end{cases}$$

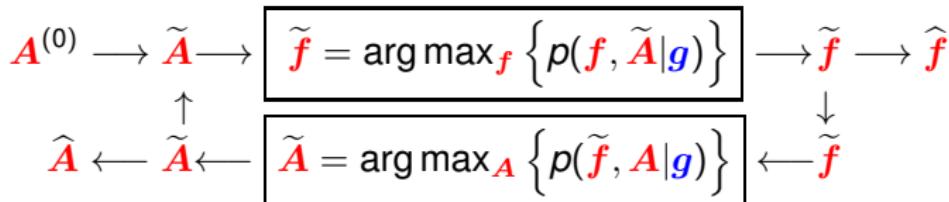
$$\rightarrow \text{EM} \begin{cases} Q(\mathbf{A}, \tilde{\mathbf{A}}) = \langle \ln p(\mathbf{f}, \mathbf{A} | \mathbf{g}) \rangle_{p(\mathbf{f} | \tilde{\mathbf{A}}, \mathbf{g})} \\ \tilde{\mathbf{A}} = \arg \max_{\mathbf{A}} \{ Q(\mathbf{A}, \tilde{\mathbf{A}}) \} \end{cases}$$

- $q_1(\mathbf{f})$  and  $q_2(\mathbf{A})$  free forms:

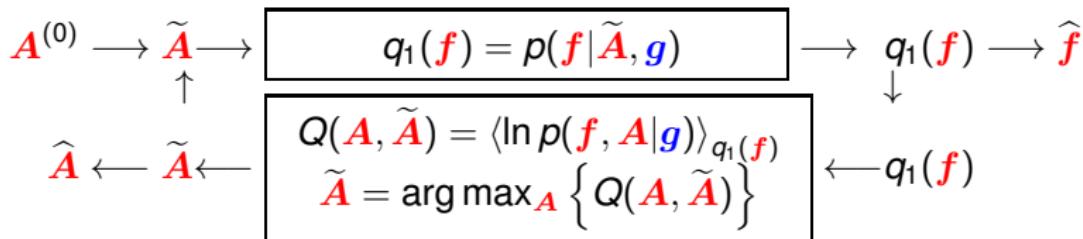
Affordable if exponential families and conjugate priors

# JMAP, EM and VBA

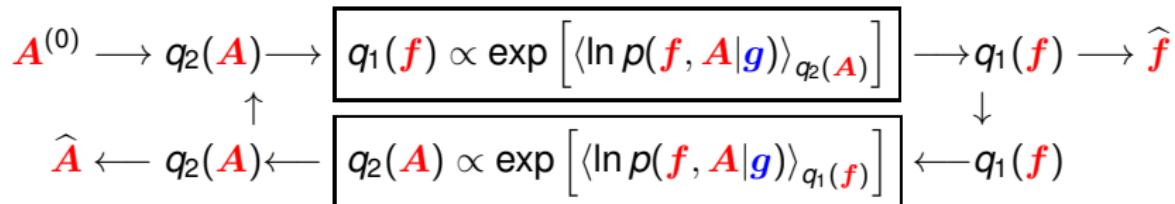
JMAP Alternate optimization Algorithm:



EM:



VBA:



# Summary of Bayesian estimation with different levels

- ▶ Marginalization over  $\mathbf{f}$

$$p(\mathbf{f}, \mathbf{A}|\mathbf{g}) \rightarrow p(\mathbf{A}|\mathbf{g}) \rightarrow \hat{\mathbf{A}} \rightarrow p(\mathbf{f}|\hat{\mathbf{A}}, \mathbf{g}) \rightarrow \hat{\mathbf{f}}$$

Joint Posterior Marginalize over  $\mathbf{f}$

- ▶ Marginalization over  $\mathbf{A}$

$$p(\mathbf{f}, \mathbf{A}|\mathbf{g}) \rightarrow p(\mathbf{f}|\mathbf{g}) \rightarrow \hat{\mathbf{f}} \rightarrow p(\mathbf{A}|\hat{\mathbf{f}}, \mathbf{g}) \rightarrow \hat{\mathbf{A}}$$

Joint Posterior Marginalize over  $\mathbf{A}$

- ▶ Joint MAP

$$p(\mathbf{f}, \mathbf{A}|\mathbf{g}) \rightarrow \begin{cases} \tilde{\mathbf{f}} = \arg \max_{\mathbf{f}} \{p(\mathbf{f}, \tilde{\mathbf{A}}|\mathbf{g})\} \\ \tilde{\mathbf{A}} = \arg \max_{\mathbf{A}} \{p(\tilde{\mathbf{f}}, \mathbf{A}|\mathbf{g})\} \end{cases} \rightarrow \hat{\mathbf{A}} \rightarrow \hat{\mathbf{f}}$$

Joint Posterior

Alternate optimization

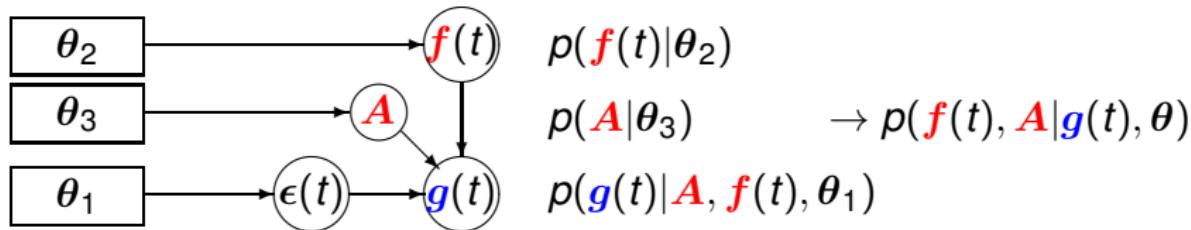
- ▶ VBA

$$p(\mathbf{f}, \mathbf{A}|\mathbf{g}) \rightarrow \begin{cases} q_1(\mathbf{f}) \propto \exp \left[ \langle \ln p(\mathbf{f}, \mathbf{A}|\mathbf{g}) \rangle_{q_2(\mathbf{A})} \right] \\ q_2(\mathbf{A}) \propto \exp \left[ \langle \ln p(\mathbf{f}, \mathbf{A}|\mathbf{g}) \rangle_{q_1(\mathbf{f})} \right] \end{cases} \rightarrow q_2(\mathbf{A}) \rightarrow \hat{\mathbf{A}} \\ \rightarrow q_1(\mathbf{f}) \rightarrow \hat{\mathbf{f}}$$

Joint Posterior

VB Approximation

## General case and Link with other methods



Different scenarios:

- ▶ IID (strict stationnarity)  $g(t) = A f(t) + \epsilon(t) \rightarrow g = A f + \epsilon$
- ▶ Second order stationnarity (Gaussian):  
All the probability laws are Gaussian
- ▶ Non Gaussian but IID: ICA
- ▶ Non-Gaussian, Time dependent, ...

# Gaussian white case: PCA, MNF, PMF and NMF

- White and Gaussian signals  $\mathbf{f}(t), \epsilon(t) \rightarrow \mathbf{g}(t)$ :

$$\mathbf{g}(t) = \mathbf{A} \mathbf{f}(t) + \epsilon(t) \rightarrow \mathbf{g} = \mathbf{A} \mathbf{f} + \epsilon$$

- Likelihood:

$$p(\mathbf{g} | \mathbf{A}, \mathbf{f}) = \mathcal{N}(\mathbf{g} | \mathbf{A}\mathbf{f}, \boldsymbol{\Sigma}_\epsilon), \quad p(\mathbf{f}) = \mathcal{N}(\mathbf{f} | 0, \boldsymbol{\Sigma}_f)$$

$$\rightarrow p(\mathbf{g} | \mathbf{A}) = \mathcal{N}(\mathbf{g} | 0, \mathbf{A}\boldsymbol{\Sigma}_f\mathbf{A}' + \boldsymbol{\Sigma}_\epsilon)$$

- PCA : Estimate  $\boldsymbol{\Sigma}_g$  by  $\frac{1}{T} \sum_t \mathbf{g}(t)\mathbf{g}'(t)$ , svd and keep all the non-zero svd:  $\boldsymbol{\Sigma}_g = \mathbf{A}\boldsymbol{\Sigma}_f\mathbf{A}'$
- Minimum Norm Factorization (MNF) :  
Estimate  $\boldsymbol{\Sigma}_g$ , svd and keep all svd  $\geq \sigma_\epsilon$ :  $\boldsymbol{\Sigma}_g = \mathbf{A}\boldsymbol{\Sigma}_f\mathbf{A}' + \boldsymbol{\Sigma}_\epsilon$
- Positive Matrix Factorization (MNF) :  
Decompose  $\boldsymbol{\Sigma}_g$  in positive definite matrices  
[Paatero & Tapper, 94]
- Non-negative Matrix Factorization (NMF) :  
Decompose  $\boldsymbol{\Sigma}_g$  in Non-negative definite matrices  
[Lee & Seung, 99]

## Non Gaussian white case: ICA

- ▶ White Non Gaussian signals and Exact model (no noise):

$$\mathbf{f}(t) \longrightarrow \mathbf{g}(t) \longrightarrow \mathbf{y}(t) = \mathbf{A}^{-1}\mathbf{g}(t) \longrightarrow \mathbf{y}(t) = \mathbf{B}\mathbf{g}(t)$$

- ▶ ICA: Find  $\mathbf{B}$  in such a way that the components of  $\mathbf{y}$  be the most independent
- ▶ Different measures of independencies:  
Entropy:

$$H(\mathbf{y}) = - \int p(y_i) \ln p(y_i) \, dy_i$$

Relative entropy or Kullback-Leibler divergence:

$$KL(p(\mathbf{y})) : \prod_i p(y_i)) = \int p(y_i) \ln \frac{\prod_i p(y_i)}{p(\mathbf{y})} \, dy_i$$

- ▶ Different choices and approximations for  $p(y_i) \rightarrow$  contrast functions, cumulants basis criteria

# Non Gaussian white case: Maximum Likelihood

- ▶ White Non Gaussian signals (Accounting for noise)

$$\mathbf{g}(t) = \mathbf{A} \mathbf{f}(t) + \boldsymbol{\epsilon}(t) \longrightarrow \mathbf{g} = \mathbf{A} \mathbf{f} + \boldsymbol{\epsilon}$$

- ▶ Likelihood:

$$p(\mathbf{g}|\mathbf{A}, \boldsymbol{\Sigma}_{\epsilon}) = \int p(\mathbf{g}|\mathbf{A}, \mathbf{f}, \boldsymbol{\Sigma}_{\epsilon}) p(\mathbf{f}) d\mathbf{f}$$

- ▶ ICA (Maximum Likelihood) :

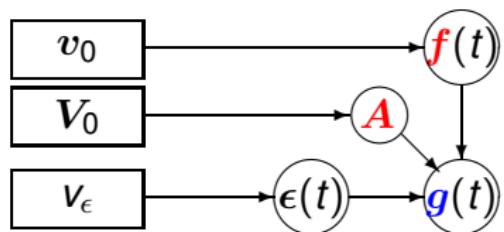
$$\hat{\boldsymbol{\theta}} = (\hat{\mathbf{A}}, \hat{\boldsymbol{\Sigma}_{\epsilon}}) = \arg \max_{\boldsymbol{\theta}} \{p(\mathbf{g}|\boldsymbol{\theta})\}$$

- ▶ EM iterative algorithm :

$$\begin{aligned} Q(\boldsymbol{\theta}, \boldsymbol{\theta}') &= E \left\{ \ln p(\mathbf{g}, \mathbf{f}, \boldsymbol{\theta} | \boldsymbol{\theta}') \right\} \\ \boldsymbol{\theta}' &= \arg \max_{\boldsymbol{\theta}} \{ Q(\boldsymbol{\theta}, \boldsymbol{\theta}') \} \end{aligned}$$

- ▶

# General Gaussian case: Joint Estimation of $\mathbf{A}$ and $\mathbf{f}$



$$p(\mathbf{f}_j(t)|v_{0j}) = \mathcal{N}(\mathbf{f}_j(t|0, v_{0j})$$

$$p(\mathbf{f}(t)|v_0) \propto \exp\left[-\frac{1}{2} \sum_j \mathbf{f}_j^2(t)/v_{0j}\right]$$

$$p(\mathbf{A}_{ij}|0, V_{0ij}) = \mathcal{N}(\mathbf{A}_{ij}|0, V_{0ij})$$

$$p(\mathbf{A}|0, V_0) = \mathcal{N}(\mathbf{A}|0, V_0)$$

$$p(\mathbf{g}(t)|\mathbf{A}, \mathbf{f}(t), v_\epsilon) = \mathcal{N}(\mathbf{Af}(t), v_\epsilon \mathbf{I})$$

$$p(\mathbf{f}_{1..T}, \mathbf{A}|\mathbf{g}_{1..T}) \propto p(\mathbf{g}_{1..T}|\mathbf{A}, \mathbf{f}_{1..T}, v_\epsilon) p(\mathbf{f}_{1..T}|v_0) p(\mathbf{A}|0, V_0)$$

$$\propto \prod_t p(\mathbf{g}(t)|\mathbf{A}, \mathbf{f}(t), v_\epsilon) p(\mathbf{f}(t)|v_0) p(\mathbf{A}|0, V_0)$$

$$p(\mathbf{f}(t)|\mathbf{g}_{1..T}, \mathbf{A}, v_\epsilon, v_0, V_0) = \mathcal{N}(\mathbf{f}(t)|\widehat{\mathbf{f}}(t), \widehat{\Sigma})$$

$$p(\mathbf{A}|\mathbf{g}_{1..T}, \mathbf{f}_{1..T}, v_\epsilon, v_0, V_0) = \mathcal{N}(\mathbf{A}|\widehat{\mathbf{A}}, \widehat{\mathbf{V}})$$

Two approaches:

- ▶ Alternate joint MAP (JMAP) estimation
- ▶ Bayesian Variational Approximation

# Joint Estimation of $\mathbf{A}$ and $\mathbf{f}$ : Alternate JMAP

Some simplification:

$$\mathbf{v}_0 = [v_f, \dots, v_f]', \quad \text{All sources a priori same variance } v_f$$

$$\mathbf{v}_\epsilon = [v_\epsilon, \dots, v_\epsilon]', \quad \text{All noise terms a priori same variance } v_\epsilon$$

$$\mathbf{A}_0 = 0, \quad \mathbf{V}_0 = v_a \mathbf{I}$$

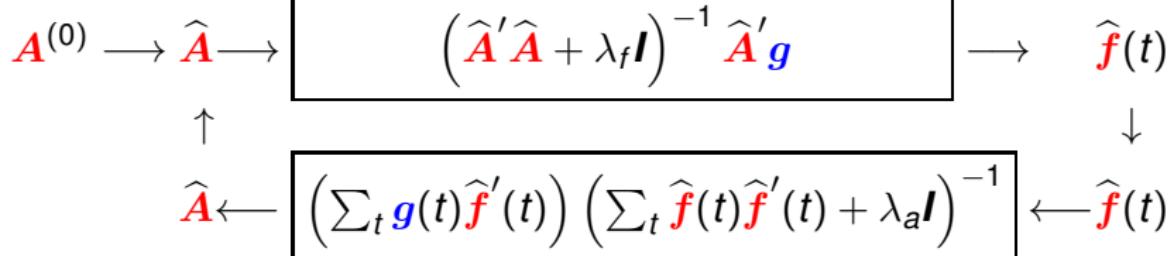
$$p(\mathbf{f}(t) | \mathbf{g}(t), \mathbf{A}, v_\epsilon, \mathbf{v}_0) = \mathcal{N}(\mathbf{f}(t) | \hat{\mathbf{f}}(t), \hat{\Sigma})$$

$$\begin{cases} \hat{\Sigma} = (\mathbf{A}' \mathbf{A} + \lambda_f \mathbf{I})^{-1} \\ \hat{\mathbf{f}}(t) = (\mathbf{A}' \mathbf{A} + \lambda_f \mathbf{I})^{-1} \mathbf{A}' \mathbf{g}(t), \quad \lambda_f = v_\epsilon / v_f \end{cases}$$

$$p(\mathbf{A} | \mathbf{g}(t), \mathbf{f}(t), v_\epsilon, \mathbf{A}_0, \mathbf{V}_0) = \mathcal{N}(\mathbf{A} | \hat{\mathbf{A}}, \hat{\mathbf{V}})$$

$$\begin{cases} \hat{\mathbf{V}} = (\mathbf{F}' \mathbf{F} + \lambda_f \mathbf{I})^{-1} \\ \hat{\mathbf{A}} = (\sum_t \mathbf{g}(t) \mathbf{f}'(t)) (\sum_t \mathbf{f}(t) \mathbf{f}'(t) + \lambda_a \mathbf{I})^{-1}, \quad \lambda_a = v_\epsilon / v_a \end{cases}$$

JMAP:



# Joint Estimation: Variational Bayesian Approximation

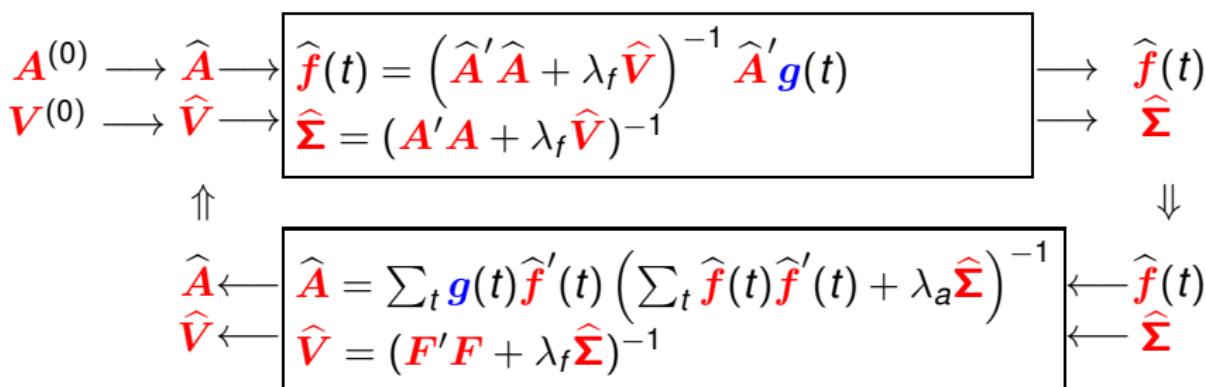
$$p(\mathbf{f}(t), \mathbf{A} | \mathbf{g}(t)) \longrightarrow q_1(\mathbf{f}(t)) q_2(\mathbf{A})$$

$$q_1(\mathbf{f}(t) | \mathbf{g}(t), \tilde{\mathbf{A}}, v_\epsilon, \mathbf{v}_0, \mathbf{V}_0) = \mathcal{N}(\mathbf{f}(t) | \tilde{\mathbf{f}}(t), \tilde{\Sigma})$$

$$\begin{cases} \tilde{\Sigma} = (\tilde{\mathbf{A}}' \tilde{\mathbf{A}} + \lambda_f \tilde{\mathbf{V}})^{-1} \\ \tilde{\mathbf{f}}(t) = (\tilde{\mathbf{A}}' \tilde{\mathbf{A}} + \lambda_f \tilde{\mathbf{V}})^{-1} \tilde{\mathbf{A}}' \mathbf{g}(t), \quad \lambda_f = v_\epsilon / v_f \end{cases}$$

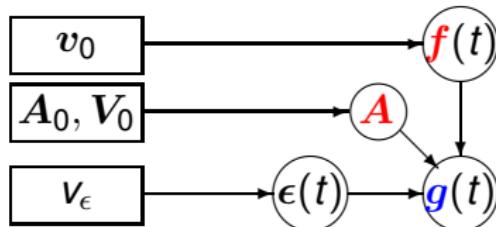
$$q_2(\mathbf{A} | \mathbf{g}(t), \tilde{\mathbf{f}}(t), v_\epsilon, \mathbf{A}_0, \mathbf{V}_0) = \mathcal{N}(\mathbf{A} | \tilde{\mathbf{A}}, \tilde{\mathbf{V}})$$

$$\begin{cases} \tilde{\mathbf{V}} = (\tilde{\mathbf{F}}' \tilde{\mathbf{F}} + \lambda_f \tilde{\Sigma})^{-1} \\ \tilde{\mathbf{A}} = \sum_t \mathbf{g}(t) \tilde{\mathbf{f}}'(t) \left( \sum_t \tilde{\mathbf{f}}(t) \tilde{\mathbf{f}}'(t) + \lambda_a \tilde{\Sigma} \right)^{-1}, \quad \lambda_a = v_\epsilon / v_a \end{cases}$$



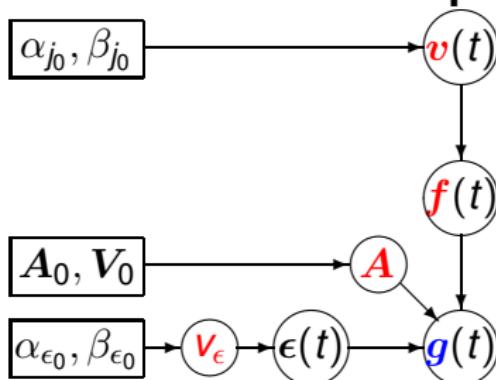
## Other more complexe models

Gaussian iid:



$$\begin{aligned} p(\mathbf{f}_j(t)|v_{0j}) &= \mathcal{N}(0, v_{0j}) \\ p(\mathbf{f}(t)|v_0) &\propto \exp\left[-\frac{1}{2}\sum_j \mathbf{f}_j^2(t)/v_{0j}\right] \\ p(\mathbf{A}_{jj}|A_{0jj}, V_{0jj}) &= \mathcal{N}(A_{0jj}, V_{0jj}) \\ p(\mathbf{A}|A_0, V_0) &= \mathcal{N}(A_0, V_0) \\ p(\mathbf{g}(t)|\mathbf{A}, \mathbf{f}(t), \nu_\epsilon) &= \mathcal{N}(\mathbf{A}\mathbf{f}(t), \nu_\epsilon \mathbf{I}) \end{aligned}$$

Variance modulated prior model inducing sparsity:

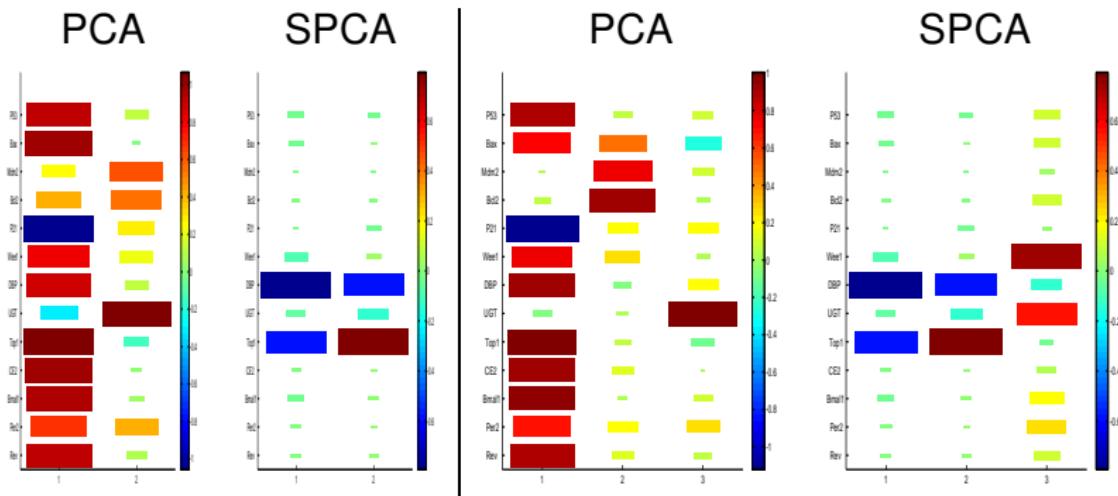


$$\begin{aligned} p(v_j(t)|\alpha_{j0}, \beta_{j0}) &= \text{IG}(\alpha_{j0}, \beta_{j0}) \\ p(\mathbf{f}_j(t)|v_j(t)) &= \mathcal{N}(0, v_j(t)) \\ p(\mathbf{f}(t)|v(t)) &\propto \exp\left[-\sum_j \mathbf{f}_j^2(t)/v_j(t)\right] \\ p(\mathbf{A}|A_0, V_0) &= \mathcal{N}(A_0, V_0) \\ p(\mathbf{g}(t)|\mathbf{A}, \mathbf{f}(t), \nu_\epsilon) &= \mathcal{N}(\mathbf{A}\mathbf{f}(t), \nu_\epsilon \mathbf{I}) \\ p(\nu_\epsilon|\alpha_{\epsilon_0}, \beta_{\epsilon_0}) &= \text{IG}(\alpha_{\epsilon_0}, \beta_{\epsilon_0}) \end{aligned}$$

$$p(\mathbf{f}_{1..T}, \mathbf{A}, \mathbf{v}_{1..T}, \nu_\epsilon | \mathbf{g}_{1..T}) \propto \prod_t p(\mathbf{g}(t)|\mathbf{A}, \mathbf{f}(t), \nu_\epsilon) p(\mathbf{f}(t)|v(t)) \prod_t \prod_j p(v_j(t)|\alpha_{j0}, \beta_{j0}) p(\mathbf{A}|A_0, V_0)$$

# Sparse PCA

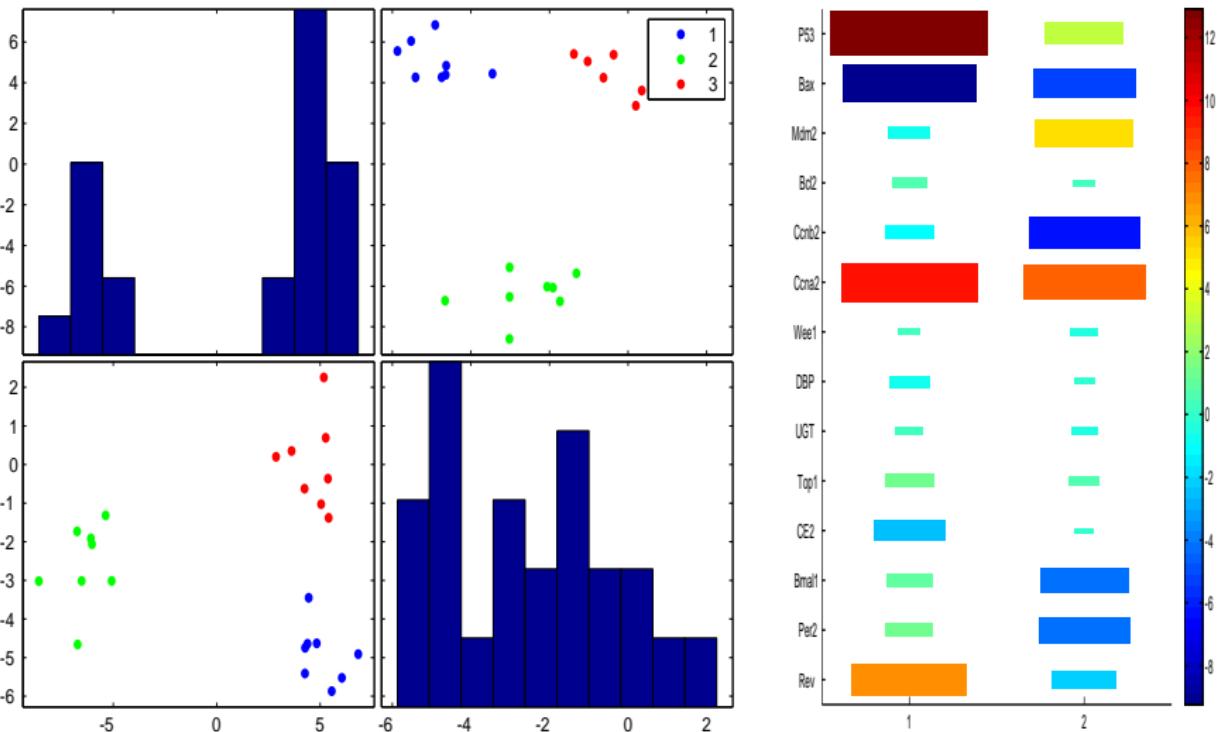
- In classical PCA, FA and ICA, one looks to obtain principal (uncorrelated or independent) components.
- In Sparse PCA or FA, one looks for the loading matrix  $A$  with sparsest components.
- This can be imposed via the prior  $p(A)$ . This leads to least variables selections.



# Discriminant Analysis

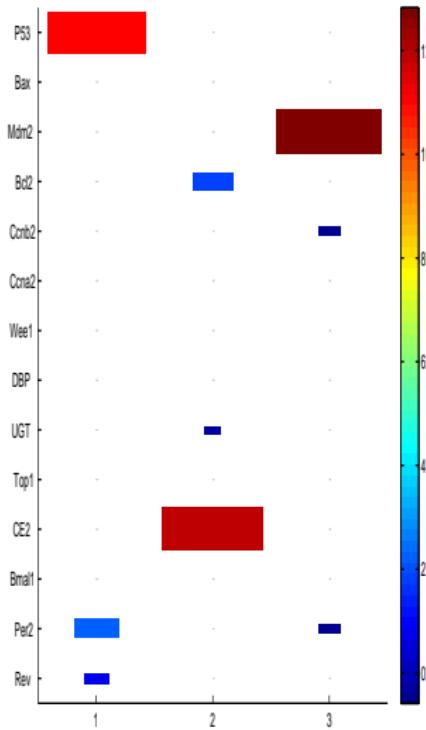
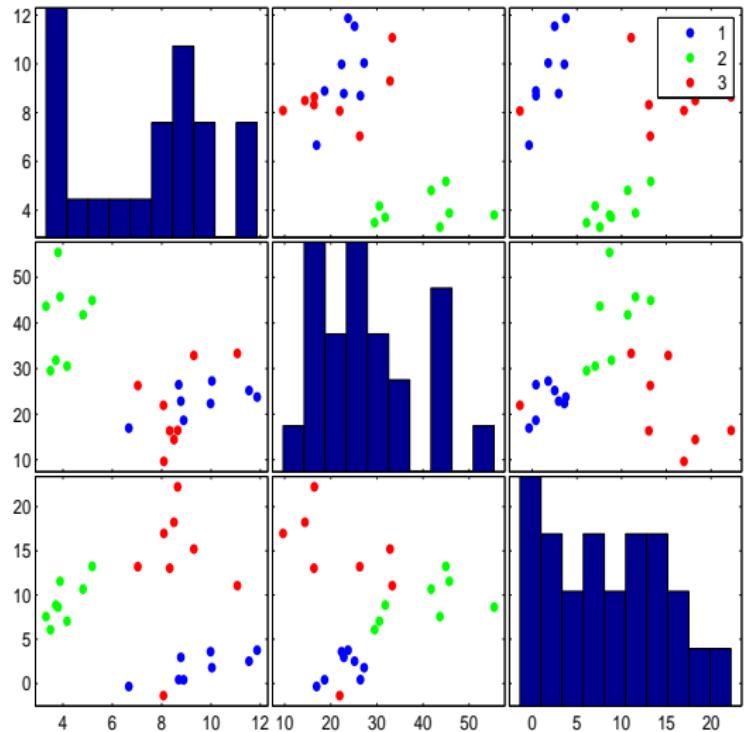
- ▶ When we have data and classes, the question to answer is:  
**What are the most discriminant factors?**
- ▶ There are many variants:
  - ▶ Linear Discriminant Analysis (LDA),
  - ▶ Quadratic Discriminant Analysis (QDA),
  - ▶ Exponential Discriminant Analysis (EDA),
  - ▶ Regularized LDA (RLDA), ...
- ▶ One can also ask for Sparsest Linear Discriminant factors (SLDA)
- ▶ Deterministic point of view (Geometrical distances)
- ▶ Probabilistic point of view (Mixture densities)
- ▶ Mixture of Gaussians models:  
Each classe is modelled by a Gaussian pdf

# Discriminant Analysis: Time series, Colon

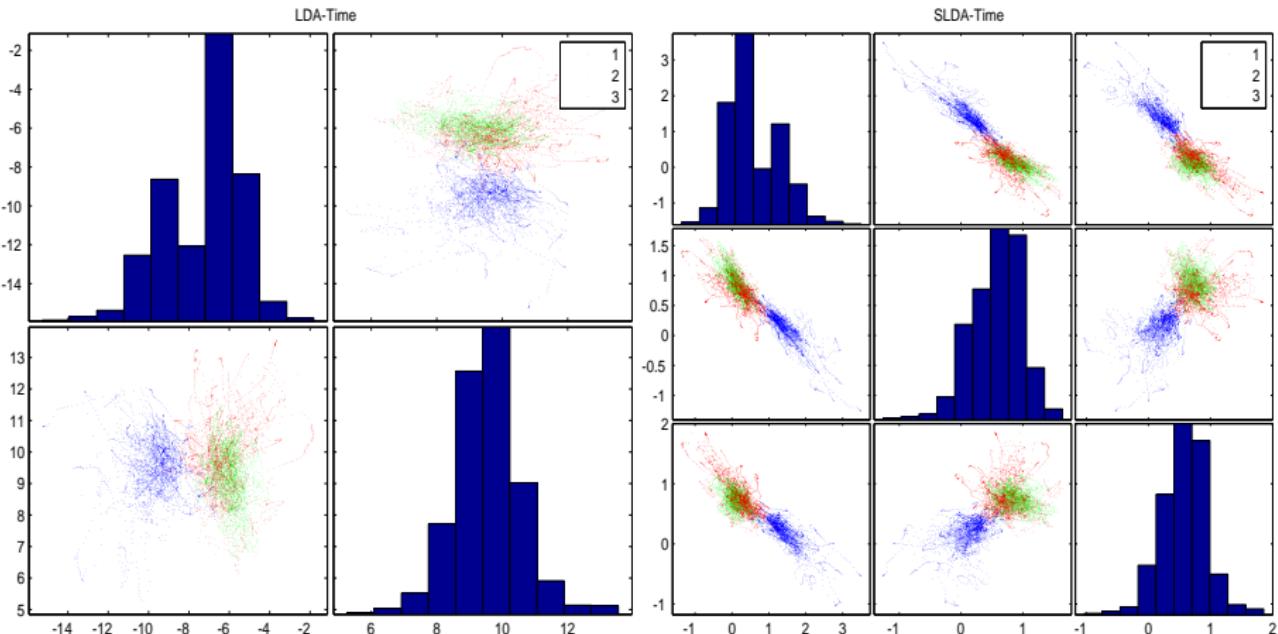


# Sparse Discriminant Analysis: Time series, colon

What are the sparsest discriminant factors?



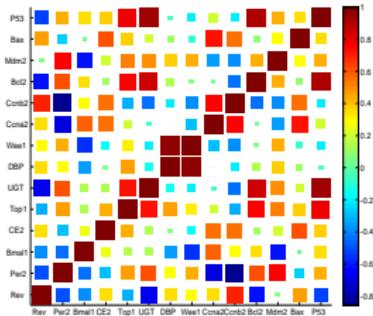
# LDA and SLDA study on time serie: 1:before, 2:during, 3:



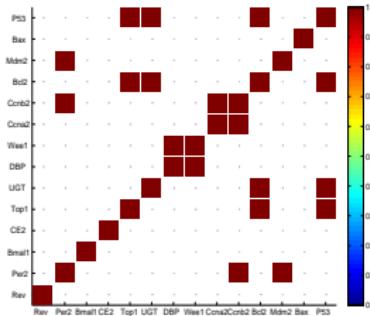
# Dependancy graphs

- The main objective here is to show the dependencies between variables
- Three different measures can be used: Pearson  $\rho$ , Spearman  $\rho_s$  and Kendall  $\tau$
- In this study we used  $\rho_s$
- A table of 2 by 2 mutual  $\rho_s$  are computed and used in different forms: Hinton, Adjacency table and Graphical network representation

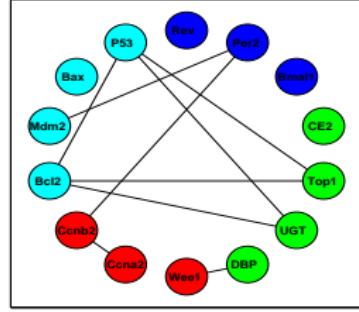
Hinton



Adjacency

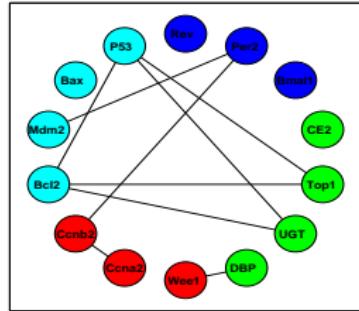
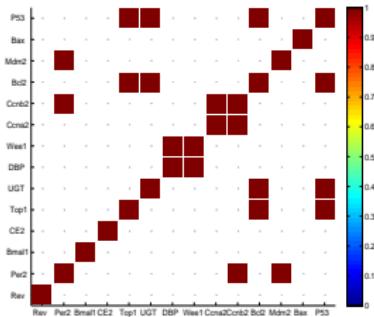
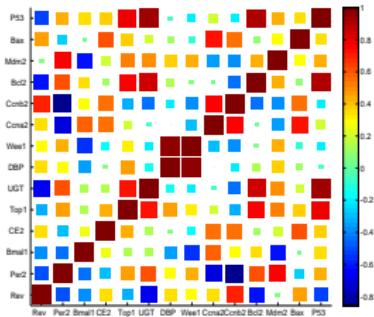


Network

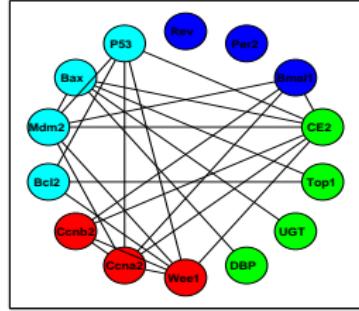
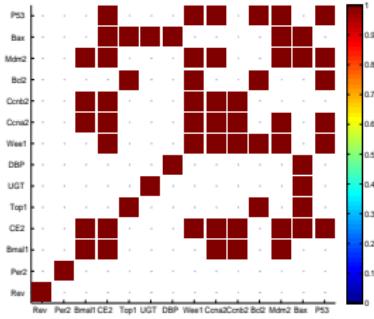
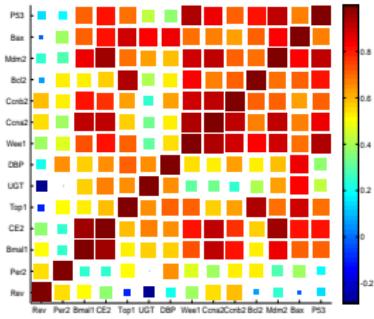


# Graph of Dependancies: Colon, Class 1

Time series



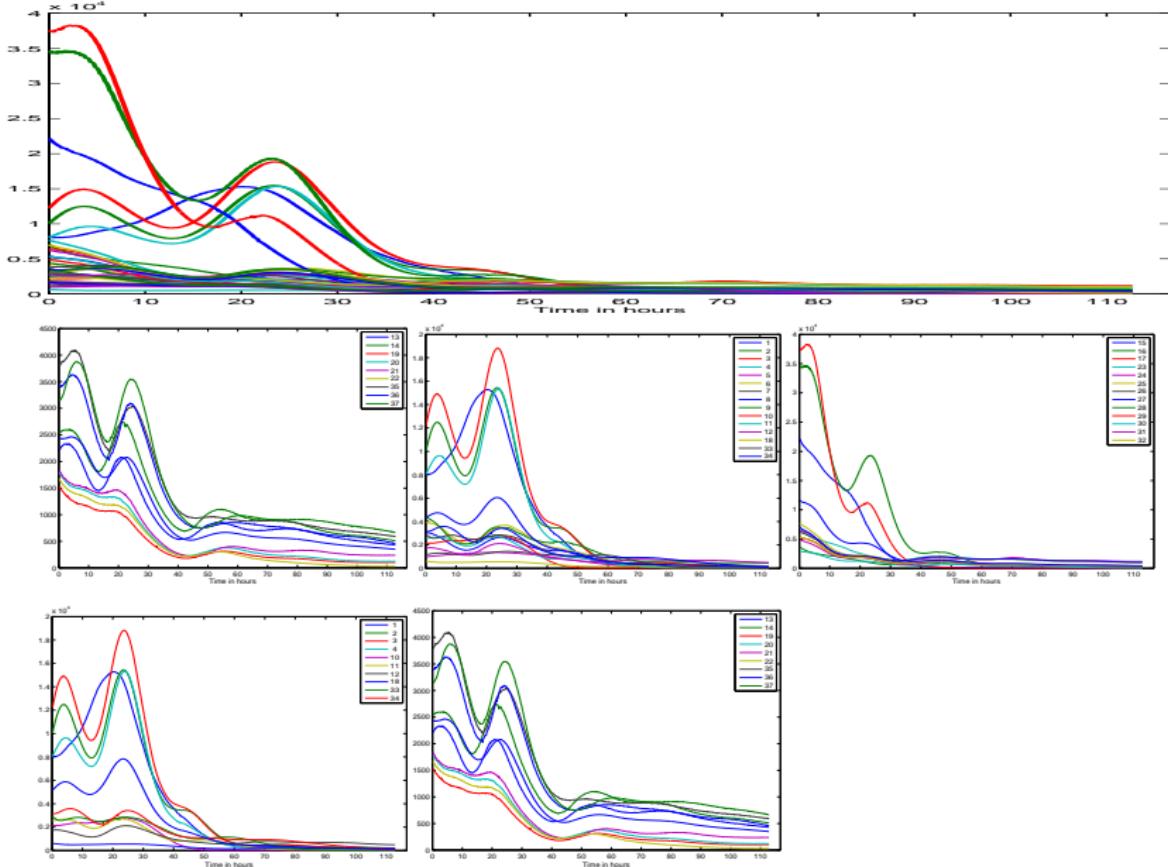
FT amplitudes



# Classification tools

- ▶ Supervised classification
  - ▶ K nearest neighbors methods
  - ▶ Needs Training sets data
  - ▶ Must be careful to measure the performances of the classification on a different set of data (Test set)
- ▶ Unsupervised classification
  - ▶ Mixture models
  - ▶ Expectation-Maximization methods
  - ▶ Bayesian versions of EM
  - ▶ Bayesian Variational Approximation (VBA)

# Classification tools



# Contents

1. Mixture models
2. Different problems related to classification and clustering
  - ▶ Training
  - ▶ Supervised classification
  - ▶ Semi-supervised classification
  - ▶ Clustering or unsupervised classification
3. Mixture of Student-t
4. Variational Bayesian Approximation
5. VBA for Mixture of Student-t
6. Conclusion

# Mixture models

- ▶ General mixture model

$$p(\mathbf{x}|\boldsymbol{a}, \Theta, K) = \sum_{k=1}^K a_k p_k(\mathbf{x}_k|\theta_k), \quad 0 < a_k < 1$$

- ▶ Same family  $p_k(\mathbf{x}_k|\theta_k) = p(\mathbf{x}_k|\theta_k)$ ,  $\forall k$
- ▶ Gaussian  $p(\mathbf{x}_k|\theta_k) = \mathcal{N}(\mathbf{x}_k|\mu_k, \Sigma_k)$  with  $\theta_k = (\mu_k, \Sigma_k)$
- ▶ Data  $\mathbf{X} = \{\mathbf{x}_n, n = 1, \dots, N\}$  where each element  $\mathbf{x}_n$  can be in one of these classes  $c_n$ .
- ▶  $a_k = p(c_n = k)$ ,  $\boldsymbol{a} = \{a_k, k = 1, \dots, K\}$ ,  
 $\Theta = \{\theta_k, k = 1, \dots, K\}$

$$p(\mathbf{X}_n, c_n = k | \boldsymbol{a}, \Theta) = \prod_{n=1}^N p(\mathbf{x}_n, c_n = k | \boldsymbol{a}, \Theta).$$

# Different problems

- ▶ Training:  
Given a set of (training) data  $\mathbf{X}$  and classes  $c$ , estimate the parameters  $a$  and  $\Theta$ .
- ▶ Supervised classification:  
Given a sample  $\mathbf{x}_m$  and the parameters  $K$ ,  $a$  and  $\Theta$  determine its class
$$k^* = \arg \max_k \{p(c_m = k | \mathbf{x}_m, a, \Theta, K)\}.$$
- ▶ Semi-supervised classification (Proportions are not known):  
Given sample  $\mathbf{x}_m$  and the parameters  $K$  and  $\Theta$ , determine its class
$$k^* = \arg \max_k \{p(c_m = k | \mathbf{x}_m, \Theta, K)\}.$$
- ▶ Clustering or unsupervised classification (Number of classes  $K$  is not known):  
Given a set of data  $\mathbf{X}$ , determine  $K$  and  $c$ .

# Training

- ▶ Given a set of (training) data  $X$  and classes  $c$ , estimate the parameters  $a$  and  $\Theta$ .
- ▶ Maximum Likelihood (ML):

$$(\hat{a}, \hat{\Theta}) = \arg \max_{(a, \Theta)} \{p(X, c|a, \Theta, K)\}.$$

- ▶ Bayesian: Assign priors  $p(a|K)$  and  $p(\Theta|K) = \prod_{k=1}^K p(\theta_k)$  and write the expression of the joint posterior laws:

$$p(a, \Theta|X, c, K) = \frac{p(X, c|a, \Theta, K) p(a|K) p(\Theta|K)}{p(X, c|K)}$$

where

$$p(X, c|K) = \iint p(X, c|a, \Theta|K) p(a|K) p(\Theta|K) da d\Theta$$

- ▶ Infer on  $a$  and  $\Theta$  either as the Maximum A Posteriori (MAP) or Posterior Mean (PM).

## Supervised classification

- Given a sample  $\mathbf{x}_m$  and the parameters  $K$ ,  $a$  and  $\Theta$  determine

$$p(c_m = k | \mathbf{x}_m, a, \Theta, K) = \frac{p(\mathbf{x}_m, c_m = k | a, \Theta, K)}{p(\mathbf{x}_m | a, \Theta, K)}$$

where  $p(\mathbf{x}_m, c_m = k | a, \Theta, K) = a_k p(\mathbf{x}_m | \theta_k)$  and

$$p(\mathbf{x}_m | a, \Theta, K) = \sum_{k=1}^K a_k p(\mathbf{x}_m | \theta_k)$$

- Best class  $k^*$ :

$$k^* = \arg \max_k \{p(c_m = k | \mathbf{x}_m, a, \Theta, K)\}$$

## Semi-supervised classification

- Given sample  $\mathbf{x}_m$  and the parameters  $K$  and  $\Theta$  (not the proportions  $a$ ), determine the probabilities

$$p(c_m = k | \mathbf{x}_m, \Theta, K) = \frac{p(\mathbf{x}_m, c_m = k | \Theta, K)}{p(\mathbf{x}_m | \Theta, K)}$$

where

$$p(\mathbf{x}_m, c_m = k | \Theta, K) = \int p(\mathbf{x}_m, c_m = k | a, \Theta, K) p(a | K) da$$

and

$$p(\mathbf{x}_m | \Theta, K) = \sum_{k=1}^K p(\mathbf{x}_m, c_m = k | \Theta, K)$$

- Best class  $k^*$ , for example the MAP solution:

$$k^* = \arg \max_k \{p(c_m = k | \mathbf{x}_m, \Theta, K)\}.$$

## Clustering or non-supervised classification

- ▶ Given a set of data  $\mathbf{X}$ , determine  $K$  and  $c$ .
- ▶ Determination of the number of classes:

$$p(K = L | \mathbf{X}) = \frac{p(\mathbf{X}, K = L)}{p(\mathbf{X})} = \frac{p(\mathbf{X}|K = L) p(K = L)}{p(\mathbf{X})}$$

and

$$p(\mathbf{X}) = \sum_{L=1}^{L_0} p(K = L) p(\mathbf{X}|K = L),$$

where  $L_0$  is the a priori maximum number of classes and

$$p(\mathbf{X}|K = L) = \int \int \prod_n \prod_{k=1}^L a_k p(\mathbf{x}_n, c_n = k | \boldsymbol{\theta}_k) p(a|K) p(\boldsymbol{\Theta}|K) da d\boldsymbol{\Theta}$$

- ▶ When  $K$  and  $c$  are determined, we can also determine the characteristics of those classes  $a$  and  $\boldsymbol{\Theta}$ .

# Mixture of Student-t model

- ▶ Student-t and its Infinite Gaussian Scaled Model (IGSM):

$$\mathcal{T}(\mathbf{x}|\nu, \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \int_0^\infty \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, z^{-1}\boldsymbol{\Sigma}) \mathcal{G}(z|\frac{\nu}{2}, \frac{\nu}{2}) dz$$

where

$$\begin{aligned}\mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) &= |2\pi\boldsymbol{\Sigma}|^{-\frac{1}{2}} \exp \left[ -\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})' \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}) \right] \\ &= |2\pi\boldsymbol{\Sigma}|^{-\frac{1}{2}} \exp \left[ -\frac{1}{2} \text{Tr} \left\{ (\mathbf{x} - \boldsymbol{\mu}) \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu})' \right\} \right]\end{aligned}$$

and

$$\mathcal{G}(z|\alpha, \beta) = \frac{\beta^\alpha}{\Gamma(\alpha)} z^{\alpha-1} \exp[-\beta z].$$

- ▶ Mixture of Student-t:

$$p(\mathbf{x}|\{\nu_k, \mathbf{a}_k, \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k, k = 1, \dots, K\}, K) = \sum_{k=1}^K a_k \mathcal{T}(\mathbf{x}_n|\nu_k, \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k).$$

## Mixture of Student-t model

- ▶ Introducing  $z_{nk}$ ,  $z_k = \{z_{nk}, n = 1, \dots, N\}$ ,  $Z = \{z_{nk}\}$ ,  
 $c = \{c_n, n = 1, \dots, N\}$ ,  
 $\theta_k = \{\nu_k, a_k, \mu_k, \Sigma_k\}$ ,  $\Theta = \{\theta_k, k = 1, \dots, K\}$
- ▶ Assigning the priors  
 $p(\Theta) = \prod_k p(\theta_k)$ , we can write:

$$p(X, c, Z, \Theta | K) = \prod_n \prod_k a_k \mathcal{N}(\mathbf{x}_n | \mu_k, z_{n,k}^{-1} \Sigma_k) \mathcal{G}(z_{nk} | \frac{\nu_k}{2}, \frac{\nu_k}{2}) p(\theta_k)$$

- ▶ Joint posterior law:

$$p(c, Z, \Theta | X, K) = \frac{p(X, c, Z, \Theta | K)}{p(X | K)}.$$

- ▶ The main task now is to propose some **approximations** to it in such a way that we can use it easily in all the above mentioned tasks of classification or clustering.

## Variational Bayesian Approximation (VBA)

- ▶ Main idea: to propose easy computational approximation  $q(c, Z, \Theta)$  for  $p(c, Z, \Theta | X, K)$ .
- ▶ Criterion:  $\text{KL}(q : p)$
- ▶ Interestingly, by noting that

$$p(c, Z, \Theta | X, K) = p(X, c, Z, \Theta | K) / p(X | K)$$

we have:

$$\text{KL}(q : p) = -\mathcal{F}(q) + \ln p(X | K)$$

where

$$\mathcal{F}(q) = \langle -\ln p(X, c, Z, \Theta | K) \rangle_q$$

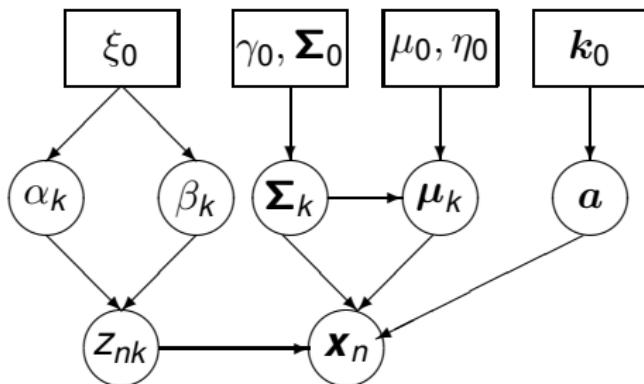
is called **free energy** of  $q$  and we have the following properties:

- Maximizing  $\mathcal{F}(q)$  or minimizing  $\text{KL}(q : p)$  are equivalent and both give an upper bound to the evidence of the model  $\ln p(X | K)$ .
- When the optimum  $q^*$  is obtained,  $\mathcal{F}(q^*)$  can be used as a criterion for model selection.

## VBA: choosing the good families

- ▶ Using  $\text{KL}(q : p)$  has the very interesting property that using  $q$  to compute the **means** we obtain the same values if we have used  $p$  (**Conservation of the means**).
- ▶ Unfortunately, this is not the case for variances or other moments.
- ▶ If  $p$  is in the exponential family, then choosing appropriate conjugate priors, the structure of  $q$  will be the same and we can obtain appropriate **fast optimization algorithms**.

# Hierarchical graphical model



Graphical representation of the model.

## VBA for mixture of Student-t

- ▶ In our case, noting that

$$p(\mathbf{X}, \mathbf{c}, \mathbf{Z}, \boldsymbol{\Theta} | K) = \prod_{n=1}^N \prod_{k=1}^K p(\mathbf{x}_n, c_n, z_{nk} | a_k, \mu_k, \boldsymbol{\Sigma}_k, \nu_k) \\ \prod_{k=1}^K [p(\alpha_k) p(\beta_k) p(\mu_k | \boldsymbol{\Sigma}_k) p(\boldsymbol{\Sigma}_k)]$$

with

$$p(\mathbf{x}_n, c_n, z_{nk} | a_k, \mu_k, \boldsymbol{\Sigma}_k, \nu_k) = \mathcal{N}(\mathbf{x}_n | \mu_k, z_{n,k}^{-1} \boldsymbol{\Sigma}_k) \mathcal{G}(z_{nk} | \alpha_k, \beta_k)$$

is separable, in one side for  $[c, Z]$  and in other size in components of  $\boldsymbol{\Theta}$ , we propose to use

$$q(c, Z, \boldsymbol{\Theta}) = q(c, Z) q(\boldsymbol{\Theta}).$$

## VBA for mixture of Student-t

- With this decomposition, the expression of the Kullback-Leibler divergence becomes:

$$\text{KL}(q_1(c, Z)q_2(\Theta) : p(c, Z, \Theta | X, K) = \\ \sum_c \int \int q_1(c, Z)q_2(\Theta) \ln \frac{q_1(c, Z)q_2(\Theta)}{p(c, Z, \Theta | X, K)} d\Theta dZ$$

- The expression of the Free energy becomes:

$$\mathcal{F}(q_1(c, Z)q_2(\Theta)) = \sum_c \int \int q_1(c, Z)q_2(\Theta) \ln \frac{p(X, c, Z | \Theta, K)p(\Theta | K)}{q_1(c, Z)q_2(\Theta)} d\Theta$$

## Proposed VBA for Mixture of Student-t priors model

- ▶ Using a generalized Student-t obtained by replacing  $\mathcal{G}(z_{n,k} | \frac{\nu_k}{2}, \frac{\nu_k}{2})$  by  $\mathcal{G}(z_{n,k} | \alpha_k, \beta_k)$  it will be easier to propose conjugate priors for  $\alpha_k, \beta_k$  than for  $\nu_k$ .

$$p(\mathbf{x}_n, c_n = k, z_{nk} | \boldsymbol{a}_k, \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k, \alpha_k, \beta_k, K) = a_k \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, z_{n,k}^{-1} \boldsymbol{\Sigma}_k) \mathcal{G}(z_{n,k} | \alpha_k, \beta_k)$$

- ▶ In the following, noting by  
 $\Theta = \{(\boldsymbol{a}_k, \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k, \alpha_k, \beta_k), k = 1, \dots, K\}$ ,  
we propose to use the factorized prior laws:

$$p(\Theta) = p(\boldsymbol{a}) \sum_k [p(\alpha_k) p(\beta_k) p(\boldsymbol{\mu}_k | \boldsymbol{\Sigma}_k) p(\boldsymbol{\Sigma}_k)]$$

with the following components:

$$\begin{cases} p(\boldsymbol{a}) = \mathcal{D}(\boldsymbol{a} | \boldsymbol{k}_0), & \boldsymbol{k}_0 = [k_0, \dots, k_0] = k_0 \mathbf{1} \\ p(\alpha_k) = \mathcal{E}(\alpha_k | \zeta_0) = \mathcal{G}(\alpha_k | 1, \zeta_0) \\ p(\beta_k) = \mathcal{E}(\beta_k | \zeta_0) = \mathcal{G}(\beta_k | 1, \zeta_0) \\ p(\boldsymbol{\mu}_k | \boldsymbol{\Sigma}_k) = \mathcal{N}(\boldsymbol{\mu}_k | \boldsymbol{\mu}_0 \mathbf{1}, \eta_0^{-1} \boldsymbol{\Sigma}_k) \\ p(\boldsymbol{\Sigma}_k) = \mathcal{IW}(\boldsymbol{\Sigma}_k | \gamma_0, \gamma_0 \boldsymbol{\Sigma}_0) \end{cases}$$

# Proposed VBA for Mixture of Student-t priors model

where

$$\mathcal{D}(\boldsymbol{a}|\boldsymbol{k}) = \frac{\Gamma(\sum_I k_I)}{\prod_I \Gamma(k_I)} \prod_I a_I^{k_I - 1}$$

is the Dirichlet pdf,

$$\mathcal{E}(t|\zeta_0) = \zeta_0 \exp[-\zeta_0 t]$$

is the Exponential pdf,

$$\mathcal{G}(t|a, b) = \frac{b^a}{\Gamma(a)} t^{a-1} \exp[-bt]$$

is the Gamma pdf and

$$\mathcal{IW}(\boldsymbol{\Sigma}|\gamma, \gamma\boldsymbol{\Delta}) = \frac{|\frac{1}{2}\boldsymbol{\Delta}|^{\gamma/2} \exp\left[-\frac{1}{2}\text{Tr}\left\{\boldsymbol{\Delta}\boldsymbol{\Sigma}^{-1}\right\}\right]}{\Gamma_D(\gamma/2)|\boldsymbol{\Sigma}|^{\frac{\gamma+D+1}{2}}}.$$

is the inverse Wishart pdf.

With these prior laws and the likelihood: joint posterior law:

$$p_k(\boldsymbol{c}, \boldsymbol{Z}, \boldsymbol{\Theta}|\boldsymbol{X}) = \frac{p(\boldsymbol{X}, \boldsymbol{c}, \boldsymbol{Z}, \boldsymbol{\Theta})}{p(\boldsymbol{X})}.$$

## Expressions of $q$

$$q(c, Z, \Theta) = q(c, Z) q(\Theta) = \prod_n \prod_k [q(c_n = k | z_{nk}) q(z_{nk})] \\ \prod_k [q(\alpha_k) q(\beta_k) q(\mu_k | \Sigma_k) q(\Sigma_k)] q(a).$$

with:

$$\left\{ \begin{array}{l} q(a) = \mathcal{D}(a | \tilde{k}), \quad \tilde{k} = [\tilde{k}_1, \dots, \tilde{k}_K] \\ q(\alpha_k) = \mathcal{G}(\alpha_k | \tilde{\zeta}_k, \tilde{\eta}_k) \\ q(\beta_k) = \mathcal{G}(\beta_k | \tilde{\zeta}_k, \tilde{\eta}_k) \\ q(\mu_k | \Sigma_k) = \mathcal{N}(\mu_k | \tilde{\mu}, \tilde{\eta}^{-1} \Sigma_k) \\ q(\Sigma_k) = \mathcal{IW}(\Sigma_k | \tilde{\gamma}, \tilde{\gamma} \tilde{\Sigma}) \end{array} \right.$$

With these choices, we have

$$\mathcal{F}(q(c, Z, \Theta)) = \langle \ln p(X, c, Z, \Theta | K) \rangle_{q(c, Z, \Theta)} = \prod_k \prod_n \mathcal{F}_{1_{kn}} + \prod_k \mathcal{F}_{2_k}$$

$$\mathcal{F}_{1_{kn}} = \langle \ln p(x_n, c_n, z_{nk}, \theta_k) \rangle_{q(c_n=k | z_{nk}) q(z_{nk})}$$

$$\mathcal{F}_{2_k} = \langle \ln p(x_n, c_n, z_{nk}, \theta_k) \rangle_{q(\theta_k)}$$

## VBA Algorithm step

Expressions of the updating expressions of the tilded parameters are obtained by following three steps:

- ▶ **E step:** Optimizing  $\mathcal{F}$  with respect to  $q(c, Z)$  when keeping  $q(\Theta)$  fixed, we obtain the expression of  $q(c_n = k | z_{nk}) = \tilde{a}_k$ ,  $q(z_{nk}) = \mathcal{G}(z_{nk} | \tilde{\alpha}_k, \tilde{\beta}_k)$ .
- ▶ **M step:** Optimizing  $\mathcal{F}$  with respect to  $q(\Theta)$  when keeping  $q(c, Z)$  fixed, we obtain the expression of  $q(a) = \mathcal{D}(a | \tilde{k})$ ,  $\tilde{k} = [\tilde{k}_1, \dots, \tilde{k}_K]$ ,  $q(\alpha_k) = \mathcal{G}(\alpha_k | \tilde{\zeta}_k, \tilde{\eta}_k)$ ,  $q(\beta_k) = \mathcal{G}(\beta_k | \tilde{\zeta}_k, \tilde{\eta}_k)$ ,  $q(\mu_k | \Sigma_k) = \mathcal{N}(\mu_k | \tilde{\mu}, \tilde{\eta}^{-1} \Sigma_k)$ , and  $q(\Sigma_k) = \mathcal{IW}(\Sigma_k | \tilde{\gamma}, \tilde{\gamma} \tilde{\Sigma})$ , which gives the updating algorithm for the corresponding tilded parameters.
- ▶  **$\mathcal{F}$  evaluation:** After each E step and M step, we can also evaluate the expression of  $\mathcal{F}(q)$  which can be used for **stopping rule** of the iterative algorithm.
- ▶ Final value of  $\mathcal{F}(q)$  for each value of  $K$ , noted  $\mathcal{F}_k$ , can be used as a criterion for **model selection**, i.e.; the **determination of the number of clusters**.

## Conclusions

- ▶ Clustering and classification of a set of data are between the most important tasks in statistical researches for many applications such as data mining in biology.
- ▶ Mixture models and in particular Mixture of Gaussians are classical models for these tasks.
- ▶ We proposed to use a **mixture of generalised Student-t distribution** model for the data via a hierarchical graphical model.
- ▶ To obtain **fast algorithms** and be able to handle large data sets, we used conjugate priors everywhere it was possible.
- ▶ The proposed algorithm has been used for clustering, classification and discriminant analysis of some biological data (**Cancer research related**), but in this paper, we only presented the main algorithm.