Multicomponent Multivariate data and signal analysis and processing in Cancer Biology

Ali Mohammad-Djafari

Groupe Problèmes Inverses Laboratoire des Signaux et Systèmes UMR 8506 CNRS - SUPELEC - Univ Paris Sud 11 Supélec, Plateau de Moulon, 91192 Gif-sur-Yvette, FRANCE.

djafari@lss.supelec.fr
http://djafari.free.fr
http://www.lss.supelec.fr

ISSSMA, June 3-4, 2013, IHP, Paris, France

June 5, 2013

A. Mohammad-Djafari, ISSSMA

ISSSMA, June 3-4, 2013, IHP, Paris, France,

1/28

イロト イヨト イヨト

Summary

- Different experiences
 - Individual cells, Population of cells, Small animals, Human
 - In vitro and In Vivo
- A great number of data, variables, time series, signals, images, $\dots \longrightarrow$ multivariate and multicomponent
- Need for Visualization tools
 - Time domain
 - Transformed domain: Fourier, Wavelets, Time-Frequency...
 - Scatter plots, histograms, statistics, ...
- Modeling time series
 - Parametric:

Superposition of sinusoids, Gaussians shapes, ...

Non Parametric:

Fourier, Wavelets, Time-frequency, time-scale,...

Probabilistic:

Moving Average (MA), Autoregressive (AR), ARMA, Markovian models, ...

(日)

2/28

- Modeling the relation between data/signals
 - Linear / Non linear
 - Training and test data ad-Djafari, ISSSMA, June 3-4, 2013, IHP, Paris, France,

A. Mohammad-Djafari,

Summary

- Simple Analysis: Estimating periods, Computing harmonic components, spectra, , ...
- Multicomponent/Multivariate data analysis: Dimensional Reduction
 PCA, FA, ICA, Sparse PCA for dimensional reduction and main factors extraction
- Multicomponent/Multivariate Discriminant Analysis with classification:

LDA, EDA, RDA, Sparse LDA for finding the most discriminant factors

- Blind sources separation
- Correlation (Pearson or Spearman) computation and dependency graph visualization

3/28

Modelling input-output relations

Temperature and activity Time series before, during and after some treatment



Temperatures, before, during and after changes



5/28

Simple Analysis tools: Period estimation. What do we mean by period ? Case of 3 sinusoids









Case of few sinusoids+noise



A. Mohammad-Djafari,

ISSSMA, June 3-4, 2013, IHP, Paris, France,



6/28

How to define a period from Spectra $S(\omega)$?

• Consider $S(\omega)$ as a distribution



How to estimate Spectra $S(\omega)$?

- Fast Fourier Transform (FFT):
 - $g(t) \longrightarrow FFT \longrightarrow f(\omega) \longrightarrow S(\omega) = |f(\omega)|^2$
 - Advantages: Well-known and understood, fast
 - Drawbacks: linear in frequencies v. but not equidistance in periods

$$\nu = [0, \cdots, N-1] \longrightarrow p = [\infty, 1, \cdots, 1/(N-1)]$$

• Autocorrelation function: $\gamma(\tau)$

• If g(t) is periodic, then $\gamma(\tau)$ is also periodic, but much smoother

•
$$\gamma(\mathbf{0}) = \mathbf{1} \ \gamma(\tau) \leq \gamma(\mathbf{0}), \forall \tau$$

- Autocorrelation function and FT: $\gamma(\tau) \longrightarrow$ FFT S(ω)
- Inverse problem approach: Compute spectra for a given interval of possible periods)
 - $f(p) \longrightarrow g(t)$ is a linear forward operation

$$g(t) = \sum_n f(2\pi/p_n) \exp\left[j2\pi nt/p_n\right]$$

8/28

 $\boldsymbol{q} = \boldsymbol{H}\boldsymbol{f} + \boldsymbol{\epsilon} \longrightarrow \widehat{\boldsymbol{f}}$ ISSSMA, June 3-4, 2013, IHP, Paris, France, A. Mohammad-Djafari,

How to estimate Spectra? Inverse Problem Approach

$$oldsymbol{g} = oldsymbol{H}oldsymbol{f} + oldsymbol{\epsilon}$$

Regularization:

$$\widehat{f} = rg\min_{f} \left\{ \|m{g} - m{H}m{f}\|^2 + \lambda \|m{f}\|^2
ight\}
ightarrow \widehat{f} = (m{H}'m{H} + \lambda m{I})^{-1}m{H}'m{g}$$

- Bayesian approach:
 - Assign the Likelihood : p(g|f)
 - Assign the prior law: p(f)
 - Use the Bayes rule : $p(f|g) \propto p(g|f) p(f)$
 - Use this posterior law to infer on *f*.
 - For example MAP:

$$\widehat{m{f}} = rg\max_{m{f}} \left\{ p(m{f}|m{g})
ight\} = rg\min_{m{f}} \left\{ J(m{f})
ight\}$$

but there are other possibilities: Posterior mean, median, ...

Assuming Gaussian noise and Gaussian prior

$$J(\boldsymbol{f}) = \|\boldsymbol{g} - \boldsymbol{H}\boldsymbol{f}\|^2 + \lambda \|\boldsymbol{f}\|^2$$

► Different priors (Gaussian, Generalized Gaussian, Cauchy,...) $J(f) = ||g - Hf||^2 + \lambda \Omega(f)$

Bayesian estimation with priors enforcing sparsity

- Sparsity: For any periodic signal, the spectrum is a set of Diracs
- Biological signals related to clocks: a few independent oscillators
- Spectrum has a few non zero elements in any given interval
- To translate this information use a heavy tailed prior law like Double exponential or Cauchy with its hierarchical structure and hidden variables

$$(\mathbf{f}|\nu) = \prod_{j} \mathcal{S}t(\mathbf{f}_{j}|\nu) \text{ with } \mathcal{S}t(\mathbf{f}_{j}|\nu) = \int_{0}^{\infty} \mathcal{N}(\mathbf{f}_{j}|0, 1/\tau_{j}) \mathcal{G}(\tau_{j}|\nu/2, \nu/2) \,\mathrm{d}\tau_{j}$$

Hierarchical structure with hidden variables τ_j:

$$\begin{cases} p(f_j|\tau_j) = \mathcal{N}(f_j|0, 1/\tau_j) \propto \exp\left[-\frac{1}{2}\tau_j f_j^2\right] \\ p(\tau_j|\alpha, \beta) = \mathcal{G}(\tau_j|\alpha, \beta) \propto \tau_j^{(\alpha-1)} \exp\left[-\beta\tau_j\right] \text{ with } \alpha = \beta = \nu/2 \end{cases}$$

► Posterior $p(f, \tau | g) \propto p(g | f) p(f | \tau) p(\tau)$

Multicomponent period estimation



 $\boldsymbol{g}_k = \boldsymbol{H}\boldsymbol{f}_k + \boldsymbol{\epsilon}_k$

f_k have some common spectra.

A. Mohammad-Djafari,

ISSSMA, June 3-4, 2013, IHP, Paris, France,

Dimension reduction, PCA, Factor Analysis, ICA

M variables g(t) are observed. They are redundant.
 Can we express them with N ≤ M factors f ?
 How many factors (Principal Components, Independent Components) can describe the observed data?

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

- How to find both A and factors f(t)?
- Deterministic methods:

$$(\widehat{m{A}},\widehat{m{f}}) = rg\min_{(m{A},m{f})} \left\{ \|m{g} - m{A}m{f}\|^2
ight\}$$
 s.t. constraints on $m{A}$ and $m{f}$

Bayesian methods:

$$(\widehat{A}, \widehat{f}) = \arg \max_{(A,f)} \{p(A, f|g)\} = \arg \min_{(A,f)} \{\ln p(g|A, f) - \ln p(A) - \ln p(f)\}$$

12/28

How to determine the number of factors

When N is given:

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

Different choices for p(A) and p(f) and Different methods to estimate both A and 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 $-\ln p(g|A, N)$ 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: 2 factors: Colon



A. Mohammad-Djafari,

ISSSMA, June 3-4, 2013, IHP, Paris, France,

14/28

Image: A matching of the second se

Factor Analysis: Time series, colon



A. Mohammad-Djafari,

ISSSMA, June 3-4, 2013, IHP, Paris, France,

・ ロ ト ・ (一 ト ・ 三 ト ・ 三 ・ つ へ) 15/28

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.



16/28

A. Mohammad-Djafari,

ISSSMA, June 3-4, 2013, IHP, Paris, France,

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)

(日) (個) (E) (E) (E)

17/28

- 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



A. Mohammad-Djafari,

ISSSMA, June 3-4, 2013, IHP, Paris, France,

18/28

A B A B A
 B A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A

Sparse Discriminant Analysis: Time series, colon

What are the sparsest discriminant factors?



A. Mohammad-Djafari,

ISSSMA, June 3-4, 2013, IHP, Paris, France,

19/28

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



A. Mohammad-Djafari, ISSSMA, June 3-4, 2013, IHP, Paris, France,

20/28

Dependancy graphs

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



21/28

A. Mohammad-Djafari,

ISSSMA, June 3-4, 2013, IHP, Paris, France,

Graph of Dependancies: Colon, Class 1

Time series



FT amplitudes











A. Mohammad-Djafari,

ISSSMA, June 3-4, 2013, IHP, Paris, France,

22/28

(日)

Graph of Dependancies: Colon, Class 3

Time series



FT amplitudes











・ロト ・ 母 ト ・ ヨ ト ・ ヨ

A. Mohammad-Djafari,

ISSSMA, June 3-4, 2013, IHP, Paris, France,

23/28

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)

24/28

- Unsupervised classification
 - Mixture models
 - Expectation-Maximization methods
 - Bayesian versions of EM
 - Bayesian Variational Approximation (VBA)

Classification tools



A. Mohammad-Djafari,

ISSSMA, June 3-4, 2013, IHP, Paris, France,

25/28

Input-Output modeling using training data and test data

- Linear models: $g_k = Af_k + \epsilon_k$, $k = 1, \cdot, K$
- Bayesian framework, MAP estimation with hyperparameter estimation

$$p(\boldsymbol{A}|\{\boldsymbol{g}_k,\boldsymbol{f}_k\}) \propto \prod_k p(\boldsymbol{g}_k|\boldsymbol{A},\boldsymbol{f}_k) p(\boldsymbol{A})$$

• Gaussian priors for ϵ_k and for A and MAP solution: $\hat{A} = \arg \max_A \{ p(A | \{g_k, f_k\}) \}$

$$\widehat{oldsymbol{A}} = \left(\sum_k oldsymbol{g}_k oldsymbol{f}_k'
ight) \left(\sum_k oldsymbol{f}_k oldsymbol{f}_k' + \lambda oldsymbol{I}
ight)^{-1}$$

- Other priors to enforce sparsity or bloc-sparsity of the prediction matrix A
- See the poster of Mircea Dumitru et al for weight loss prediction from the two genes expressions of Bmal1 and Rev-erb-alpha

26/28

Conclusions

- A lot to do to answer the questions of biologists
- Forward modeling and Bayesian inference are natural tools to answer these questions
- Very often the questions are ill-posed inverse problems which need prior knowledge
- Appropriate translation of prior knowledge to prior laws is very important
- Carefull computational algorithms have to be developped
- Carfeful presentation and interpretation of the inference results are very important
- Constant dialogue between "biologists" and "Data and Signal processors" is of great importance.

白人名德人名英人名英人

27/28

Publications

 J. Lapuyade-Lahorgue and A. Mohammad-Djafari, Nearest neighbors and correlation dimension for dimensionality estimation. Application to factor analysis of real biological time series data, in ESANN 2011 Proceedings. Gent. Belgium

in ESANN 2011 Proceedings, Gent, Belgium. ISBN 978-2-87419-044-5

 A. Mohammad-Djafari, G. Khodabandelou and J. Lapuyade-Lahorgue,
 A Matteh tealbary for data reduction visualization of

A Matlab toolbox for data reduction, visualization, classification and knowledge extraction of complex biological data, BIOCOMP2011, Las Vegas, USA

A. Mohammad-Djafari,

Bayesian approach with prior models which enforce sparsity in signal and image processing,

Review paper accepted for publication in European Association for Signal, Speech, and Image Processing (EURASIP) special issue on Sparsity in signal and image processing.