



**HAL**  
open science

## A method for statistical learning in large databases of heterogeneous imaging, cognitive and behavioral data.

Luigi Antelmi, Marco Lorenzi, Valeria Manera, Philippe Robert, Nicholas Ayache

### ► To cite this version:

Luigi Antelmi, Marco Lorenzi, Valeria Manera, Philippe Robert, Nicholas Ayache. A method for statistical learning in large databases of heterogeneous imaging, cognitive and behavioral data.. EPICLIN 2018 - 12ème Conférence Francophone d'Epidémiologie Clinique / CLCC 2018 - 25èmes Journées des statisticiens des Centre de Lutte Contre le Cancer, May 2018, Nice, France. Elsevier, Revue d'Épidémiologie et de Santé Publique, 66 (3), pp.S180, 2018, 12e Conférence francophone d'Épidémiologie clinique 25e Journée des statisticiens des Centres de lutte contre le cancer. 10.1016/j.respe.2018.03.306 . hal-01827389

**HAL Id: hal-01827389**

**<https://inria.hal.science/hal-01827389v1>**

Submitted on 2 Jul 2018

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

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

## EPICLIN 2018

### **A method for statistical learning in large databases of heterogeneous imaging, cognitive and behavioural data.**

**Authors:** Luigi Antelmi<sup>1,\*</sup>, Marco Lorenzi<sup>1</sup>, Valeria Manera<sup>2,3</sup>, Philippe Robert<sup>3,4</sup>, Nicholas Ayache<sup>1</sup>

\* corresponding author: luigi.antelmi@inria.fr

#### **Affiliations:**

1 UCA, Inria Sophia Antipolis, Epione Research Project.

2 UCA, Inria Sophia Antipolis, Stars Research Project.

3 CoBTeK, University of Nice Sophia Antipolis.

4 Centre Memoire, CHU de Nice.

#### **The authors declare no conflict of interest**

#### **Abstract**

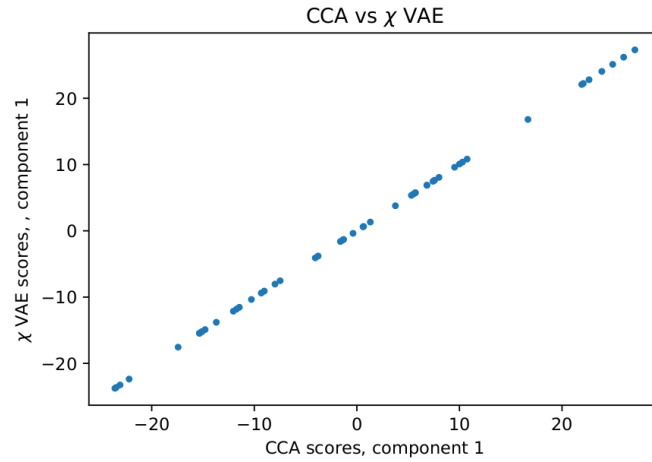
*Introduction* The aim of this study is to develop a generative and probabilistic statistical learning model for the joint analysis of heterogeneous biomedical data. The model will be applied to the investigation of neurological disorders from collections of brain imaging, body sensors, biological and clinical data available in current large-scale health databases. The resulting methodological framework will be tested on the UK Biobank, as well as on pathology-specific clinical data, as provided by the ADNI, or INSIGHT initiatives.

*Methods* We propose a variational approximation of Bayesian Canonical Correlation Analysis (CCA). The proposed formulation is inspired by current advanced in variational learning, and offers the potential to scale to high-dimensional observations, such as medical images and arrays of biological data.. We proved that the variational lower bound can be optimized through modern learning libraries such as Torch and TensorFlow.

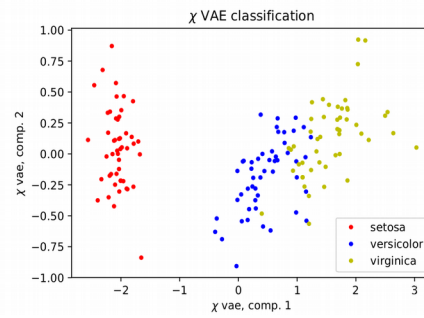
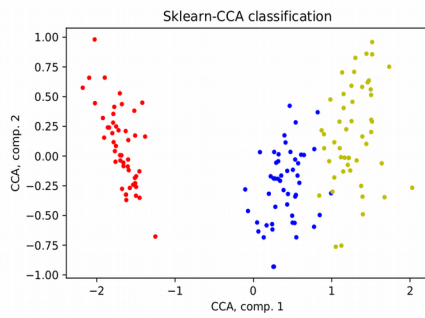
*Results* We currently benchmarked the method with respect to classical CCA on both synthetic data and on the classical benchmarking datasets in machine learning (IRIS dataset). With respect to the synthetic dataset (Fig. 1A), we observed a strong agreement between the score components computed with classical CCA and our method. Moreover, the classification results on IRIS showed that the two methods essentially provide the same latent representation (Fig. 1B).

*Conclusion* Our method shows promising results for the future application to medical data. The method is computationally efficient and scalable, hence able to process complex multivariate multidimensional datasets. We expect to highlight meaningful relationship among biomarkers that could be used to develop optimal strategies for disease classification, quantification, and prediction. In the future the proposed approach will be tested in several experimental settings:

- 1) Classification/stratification;
- 2) Prediction and imputation from a set of observed data (e.g. predict biological and clinical output from medical imaging information).



(A)



(B)

**Fig 1. (A) Comparison of score components computed with classical CCA and our method from the synthetic dataset. The perfect linear relationship between two methods applied to this datasets means that they both capture the same latent representation. (B) Classification of the IRIS dataset. Classical CCA on the left; our method on the right.**