



**HAL**  
open science

## Plsda versus pca on barycenters applied to metabolomics in a context of discrimination

Marion Brandolini-Bunlon, Benoît Jaillais, Mohamed Hanafi

### ► To cite this version:

Marion Brandolini-Bunlon, Benoît Jaillais, Mohamed Hanafi. Plsda versus pca on barycenters applied to metabolomics in a context of discrimination. 11th Colloquium Chemiometricum Mediterraneum (CCM XI 2023), Jun 2023, Padoue / Padova, Italy. hal-04151913

**HAL Id: hal-04151913**

**<https://hal.science/hal-04151913>**

Submitted on 5 Jul 2023

**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.

# PLSDA VERSUS PCA ON BARYCENTERS APPLIED TO METABOLOMICS IN A CONTEXT OF DISCRIMINATION

**Marion Brandolini-Bunlon<sup>1</sup>, Benoît Jaillais<sup>2</sup>, Mohamed Hanafi<sup>2</sup>**

*Addresses: <sup>1</sup>Université Clermont Auvergne, INRAE, UNH, Plateforme d'Exploration du Métabolisme, MetaboHUB Clermont, 63000 Clermont-Ferrand, France*

*<sup>2</sup>Oniris, INRAE, StatSC, 44300 Nantes, France*

*E-mail: marion.brandolini-bunlon@inrae.fr*

Abstract for **oral communication**

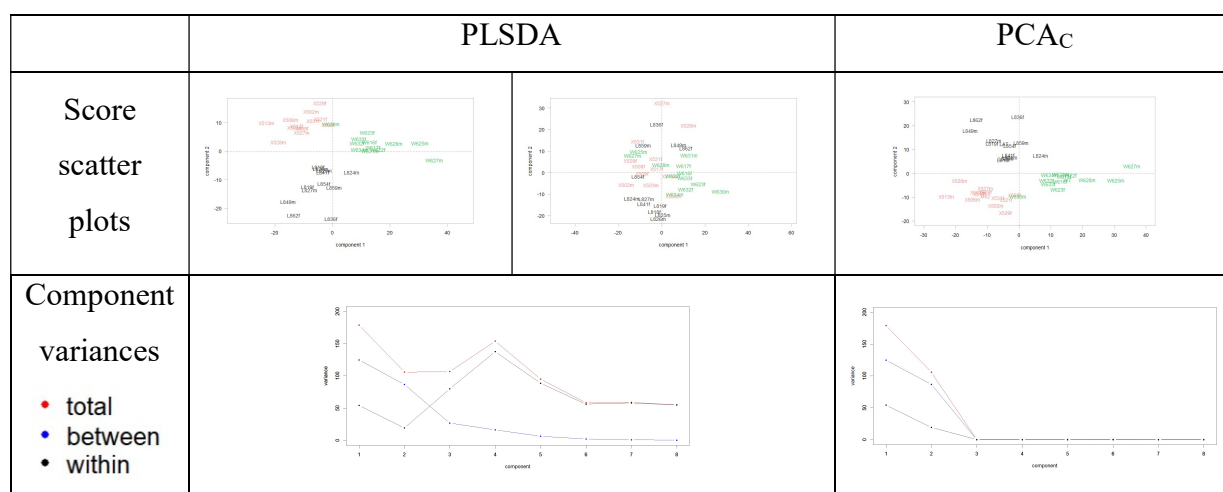
---

**Introduction.** Untargeted metabolomics is a powerful phenotyping tool to better understand the biological mechanisms involved in the physiopathological processes, and identify biomarkers of metabolic status. The complex data need dedicated preparation and treatments to extract meaningful information. The major specificity of metabolomics data is the large number of variables compared to the number of samples, as well as their high degree of correlation. The common analysis strategy consists in performing univariate and multivariate statistics to highlight variables of interest. In a discriminant context, partial least squares-discriminant analysis (PLSDA) is one of the most effective multivariate tools currently used, because of its ability to analyze collinear and noisy data. Another multivariate method that could be used is the Principal Component Analysis (PCA) of the matrix of barycenters of the observation groups (called here “PCA<sub>C</sub>”). The objective of our study is to compare these approaches in terms of explained variances and important variables.

**Material and methods.** Published data from a project on the impact of genetic mutations in mice (ProMetIS) were used as a case study (Imbert, 2021). Mice (n=42), males and females, belonged to one of the three genotype groups (wild type, lacking the linker for activation of T cells gene, or lacking the MX dynamin-like GTPase 2 gene). The metabolomics dataset we used, was obtained from the analysis of plasma samples using a mass spectrometry-based untargeted approach (LC-MS), and contained 6104 variables after preparation. In the present work, data analysis was performed with the R-package “rchemo”. Six atypical mice were removed to have a balanced experiment design, before Pareto scaling. Due to the sex effect,

without interaction with the genotype, the data were centered by sex before applying PLSDA, and PCA<sub>C</sub>, to discriminate genotype groups. On one hand, the optimal number of PLS components was determined according to the error in repeated cross-validation (30 repetitions of 10-fold cross-validation) and application of the one-standard-error-rule. On the other hand, using PCA<sub>C</sub> model, subjects were projected onto the components. For each PLSDA or PCA<sub>C</sub> component, the total, inter- and intra-group variances were then calculated, and the group effect was assessed by ANOVA. ANOVA were also performed for metabolomics variables becoming important in the discrimination in the PLSDA model with the optimal number of components (called here “PLSDA<sub>opt</sub>”), compared to the one with 2 components.

**Results.** The optimal number of PLS components was 3, and the number of PCA<sub>C</sub> components was 2. Only these components had a significant p-value in the ANOVA. As expected, the 1<sup>st</sup> component of both methods were the same, and, as shown below, in our study, the 2<sup>nd</sup> components were closely similar. The 3<sup>rd</sup> component of the PLSDA model was also of interest because it still significantly explained intergroup variability and highlighted other important discriminant variables.



**Discussion and conclusion.** The PLSDA and the PCA<sub>C</sub> components maximize the intergroup variability. When 3 groups are to be discriminated, the PCA<sub>C</sub> finds 2 components while the PLSDA can find more. Presumably, each component of the PLSDA<sub>opt</sub> can discriminate one group from one or both others, and would thus allow a better discrimination.

## References

Imbert A, *et al.* (2021) ProMetIS, deep phenotyping of mouse models by combined proteomics and metabolomics analysis. *Sci Data*, 8(1), 311.