Combining ASCA and mixed models to analyse high dimensional designed data

M. Martin¹ P. De Tullio² B. Govaerts¹

¹ Institute of Statistics, Biostatistics and Actuarial Sciences (ISBA), Université catholique de Louvain (UCL), Belgium

²Laboratoire de Chimie Pharmaceutique, Université de Liège (ULg), Liège 1, Belgium

January 30, 2018







Challenging data analysis

Application to life science data:

- From -omics sciences: genomics, transcriptomics, metabolomics
- High Dimensional: multivariate; often *m* variables > *n* samples
- Biological variability, instrumental noise/artifacts
- Multicollinearity between variables

Nontrivial Design Of Experiments (DOE):

- ex.: longitudinal, multi-centre, cross-over studies, etc.
- Presence of random factors (day, lab variation, ...)

Introduction	ASCA-related methods	мімонд³ 000000000000000	In summary	Acknowledgments	References
Table of	content				

- High Dimensional designed data analysis
- ASCA methodology
- Extension to mixed models: description & application



Methods needed to analyse HD designed data

Multivariate projection methods

- PCA, ICA, (O)PLS, ...
- Disadvantages:
 - Simple exp. designs (e.g. 2 groups comparison)
 - Few statistical modelling and tests

Statistical regression methods

Depends on the response dimension:

- $y_{(1 \times m)}$: linear or logistic regression, ANOVA, mixed models
- $Y_{(n \times m)}$ with m < n: MANOVA, multivariate-GLM

But often $m \gg n \Rightarrow$ need to combine dimension reduction and statistical modelling.







4/24





STEP 4: Global measure of effect significance based on permutation tests



1. Balanced ANOVA designs with fixed factors

 \Rightarrow Generalisation to unbalanced data with fixed and random factors

2. // ANOVA does not take into account the correlation between the responses

⇒ Prior PCA reduction & back-transformations

3. Permutation tests implementation is challenging for advanced DOE [Anderson and Braak, 2003]

 \Rightarrow Use of alternative test strategies: Likelihood ratio tests & bootstrap



APCA+ for unbalanced designs (Thiel et al. [2017], Guisset et al. [Submitted])





The Metabiose repeatability datasets

Context:

- Urine and Serum ¹H-NMR spectral data from control and endometriosis patients
- Statistical perspective of spectral reproducibility/repeatability and quality control: not yet well studied in metabolomics
- Unbalanced design for the Serum dataset

3 factors:

- 1 fixed: groups (G; 2)
- 2 random: patients (P; 7/group) & repetitions over weeks (W; 3)

Main research questions:

- Compare the groups
- Quantify the variability of the repetitions and the patients
- Test the significance of these random/fixed effects

		W1	W2	W3
	P 1	marken	mult	munta
Endometriosis	P2	hr	mult	munta
	P 8	munhm	muln	induction
Control	P9	manlo	mult	mulm

Experimental design





Step 0: Prior PCA dimension reduction (2)

Transform the highly correlated response matrix into a reduced number of orthogonal components without information loss

PCA dimension reduction on the response matrix $Y = T_C P'_C + F$

Keep C = 11 first Principal Components (PC) with $\sum_{C} var(PC) \ge 99\%$





General framework for mixed models

Fixed + random factors vs linear regression (only 1 source of random variation) The mixed model for one response t_c can be written as:

- Drop the hypothesis of independence between the samples
 - \Rightarrow model advanced designs
- Coding: Sum coding for fixed and dummy coding for random effects
- Typical applications: Multicentre study, Multilevel data, Repeated data, etc.
- Parameters are estimated with the Restricted Maximum Likelihood (REML) method

PCA on fixed/random pure effect matrices

PCA on fixed/random Residual-Augmented effect matrices

In APCA: E added to M_f

But which variance components should be added for mixed models?

Solution: RA effect matrices based on the ANOVA F-tests (Expected Mean Squares ratio)

 $egin{array}{l} ilde{M}_G = M_G + M_P + E \ ilde{M}_P = M_P + E \ ilde{M}_R = M_R + E \end{array}$

Step 3: Quantification of effect importance

For each response t_c , c = 1, .., C:

Random effects: Variance components

$$\hat{\sigma}_{P,c}^2; \hat{\sigma}_{R,c}^2; \hat{\sigma}_c^2$$

Fixed effects [Nakagawa and Schielzeth, 2013]:

$$\hat{\sigma}^2_{G,c} = ext{var}(\hat{eta}_{G,c} x_G)$$

For all t_c responses:

Total variance:

$$\hat{\sigma}_{tot}^2 = \sum_{c=1}^{C} (\hat{\sigma}_{G,c}^2 + \hat{\sigma}_{P,c}^2 + \hat{\sigma}_{R,c}^2 + \hat{\sigma}_c^2)$$

Medium

	Serum	Urine
Group	14.32	2.4
Patient	64.03	91.48
Repetition	0.6	0
Residuals	21.05	6.12

Effect importance percentage

Introduction	ASCA-related methods	MiMoHD ³ ○○○○○○○○○○○●○○	In summary	Acknowledgments	References
Step 4:	Global measur				

• Likelihood Ratio Test (LRT)

Test the significance of a fixed/random effect in a (mixed) linear model Compare the likelihoods L between nested models

LRT statistic: $2[\log(L_{full}) - \log(L_{null})] \sim_{H_0} \chi^2_{df}$

• a Global LRT (GLRT)

The test statistic for a fixed/random effect matrix

GLRT statistic: $2\left[\sum_{c=1}^{C} (\log(L_{full,c}) - \log(L_{null,c}))\right] \sim_{H_0} \chi^2_{C \times df}$

- Assess the significance of an effect based on:
 - The χ^2 distribution with known *df* (fixed effects)
 - bootstrap simulations (fixed/random effects)

Introduction	ASCA-related methods	МіМоНD³ 000000000000000	In summary	Acknowledgments	References
In summary					

MiMoHD³, a combination of mixed models & multivariate projection methods

- Innovative extension and generalisation in the ASCA(+) framework
- Enable to model unbalanced designs with random factors
- Take into account the correlation between the response variables
- Global test of effect significance
- Quantification and comparison of the mixed variability sources
- Targeted applications: repeatability/reproductibility study & longitudinal data

MiMoHD³

In summary

Acknowledgments

Introduction	ASCA-related methods	мімонд³ 000000000000000	In summary	Acknowledgments	References
Referen	ces				

Marti Anderson and Cajo Ter Braak. Permutation tests for multi-factorial analysis of variance. *Journal of Statistical Computation and Simulation*, 73(2):85–113, 2003. doi: 10.1080/00949650215733. URL http://dx.doi.org/10.1080/00949650215733.

Séverine Guisset, Bernadette Govaerts, and Manon Martin. Comparison of parafasca, acomdim, and amopls approaches in the multivariate glm modelling of multi-factorial designs. *Chemometrics and Intelligent Laboratory Systems*, Submitted.

Jeroen J. Jansen, Huub C. J. Hoefsloot, Jan van der Greef, Marieke E. Timmerman, Johan A. Westerhuis, and Age K. Smilde. Asca: analysis of multivariate data obtained from an experimental design. *Journal of Chemometrics*, 19(9):469–481, 2005. ISSN 1099-128X. doi: 10.1002/cem.952. URL http://dx.doi.org/10.1002/cem.952.

Shinichi Nakagawa and Holger Schielzeth. A general and simple method for obtaining r2 from generalized linear mixed-effects models. *Methods in Ecology and Evolution*, 4(2):133–142, 2013. ISSN 2041-210X. doi: 10.1111/j.2041-210x.2012.00261.x. URL http://dx.doi.org/10.1111/j.2041-210x.2012.00261.x.

Michel Thiel, Baptiste Féraud, and Bernadette Govaerts. Asca+ and apca+: Extensions of asca and apca in the analysis of unbalanced multifactorial designs. *Journal of Chemometrics*, 31(6): e2895–n/a, 2017. ISSN 1099-128X. doi: 10.1002/cem.2895. URL http://dx.doi.org/10.1002/cem.2895. e2895 cem.2895.