



HAL
open science

MBA-GUI: A chemometric graphical user interface for multi-block data visualisation, regression, classification, variable selection and automated pre-processing

Puneet Mishra, Jean Michel Roger, Douglas Rutledge, Alessandra Biancolillo, Federico Marini, Alison Nordon, Delphine Jouan-Rimbaud Bouveresse

► To cite this version:

Puneet Mishra, Jean Michel Roger, Douglas Rutledge, Alessandra Biancolillo, Federico Marini, et al.. MBA-GUI: A chemometric graphical user interface for multi-block data visualisation, regression, classification, variable selection and automated pre-processing. *Chemometrics and Intelligent Laboratory Systems*, 2020, 205, 10.1016/j.chemolab.2020.104139 . hal-02959982

HAL Id: hal-02959982

<https://hal.inrae.fr/hal-02959982>

Submitted on 8 Apr 2021

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.

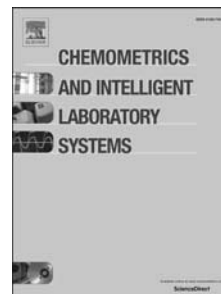
L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

Copyright

Journal Pre-proof

MBA-GUI: A chemometric graphical user interface for multi-block data visualisation, regression, classification, variable selection and automated pre-processing

Puneet Mishra, Jean Michel Roger, Douglas N. Rutledge, Alessandra Biancolillo, Federico Marini, Alison Nordon, Delphine Jouan-Rimbaud-Bouveresse



PII: S0169-7439(20)30395-6

DOI: <https://doi.org/10.1016/j.chemolab.2020.104139>

Reference: CHEMOM 104139

To appear in: *Chemometrics and Intelligent Laboratory Systems*

Received Date: 2 July 2020

Revised Date: 12 August 2020

Accepted Date: 16 August 2020

Please cite this article as: P. Mishra, J.M. Roger, D.N. Rutledge, A. Biancolillo, F. Marini, A. Nordon, D. Jouan-Rimbaud-Bouveresse, MBA-GUI: A chemometric graphical user interface for multi-block data visualisation, regression, classification, variable selection and automated pre-processing, *Chemometrics and Intelligent Laboratory Systems* (2020), doi: <https://doi.org/10.1016/j.chemolab.2020.104139>.

This is a PDF file of an article that has undergone enhancements after acceptance, such as the addition of a cover page and metadata, and formatting for readability, but it is not yet the definitive version of record. This version will undergo additional copyediting, typesetting and review before it is published in its final form, but we are providing this version to give early visibility of the article. Please note that, during the production process, errors may be discovered which could affect the content, and all legal disclaimers that apply to the journal pertain.

© 2020 Published by Elsevier B.V.

| | | |
|----|--|-----------|
| 1 | Table of Contents | |
| 2 | Abstract | 4 |
| 3 | Introduction | 4 |
| 4 | Similar works | 8 |
| 5 | Software description | 8 |
| 6 | Software architecture and brief mathematical background of techniques available | 10 |
| 7 | <i>Pre-processing</i> | 10 |
| 8 | Smoothing operations | 10 |
| 9 | Scatter correction, Baseline correction and Normalisation | 11 |
| 10 | Derivatives | 11 |
| 11 | <i>Pattern recognition and data visualisation tool</i> | 12 |
| 12 | ComDim- based methods | 12 |
| 13 | <i>Regression</i> | 14 |
| 14 | SO-PLS Regression..... | 14 |
| 15 | ComDim regression..... | 16 |
| 16 | <i>Classification</i> | 16 |
| 17 | SO-PLS-LDA | 16 |
| 18 | ComDim - based linear discriminant analysis | 17 |
| 19 | <i>Variable selection</i> | 17 |
| 20 | SO-CovSel | 17 |
| 21 | <i>SPORT</i> | 18 |
| 22 | <i>Standard one block chemometric analysis</i> | 18 |
| 23 | Datasets for MBA-GUI demonstration | 18 |
| 24 | Operating procedure and demo analysis | 21 |
| 25 | <i>Data loading and pre-processing</i> | 21 |
| 26 | <i>Data visualisation</i> | 24 |
| 27 | <i>Regression</i> | 26 |
| 28 | <i>Classification</i> | 29 |
| 29 | <i>Variable selection</i> | 31 |
| 30 | <i>SPORT</i> | 33 |
| 31 | Conclusion | 35 |
| 32 | Validation | 35 |
| 33 | Disclaimer | 37 |
| 34 | Acknowledgments | 37 |

35 **References**.....**38**

36 **Supplementary** Error! Bookmark not defined.

37

38

39

40

Journal Pre-proof

41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65
66
67

MBA-GUI: A chemometric graphical user interface for multi-block data visualisation, regression, classification, variable selection and automated pre-processing

Puneet Mishra^{1,2}, Jean Michel Roger^{3,4}, Douglas N. Rutledge^{5,6}, Alessandra Biancolillo⁷,
Federico Marini⁸, Alison Nordon², Delphine Jouan-Rimbaud-Bouveresse⁹

¹*Food and Biobased Research, Wageningen University and Research, Bornse Weilanden 9, 6708 WG, Wageningen, The Netherlands.*

²*WestCHEM, Department of Pure and Applied Chemistry and Centre for Process Analytics and Control Technology, University of Strathclyde, Glasgow, G1 1XL, United Kingdom*

³*ITAP, IRSTEA, Montpellier SupAgro, University Montpellier, Montpellier, France*

⁴*ChemHouse Research Group, Montpellier, France*

⁵*Université Paris-Saclay, INRAE, AgroParisTech, UMR SayFood, 75005, Paris, France*

⁶*National Wine and Grape Industry Centre, Charles Stuart University, Wagga Wagga, Australia*

⁷*Department of Physical and Chemical Sciences, University of L'Aquila, Via Vetoio 67100, Coppito, L'Aquila, Italy*

⁸*Department of Chemistry, University of Rome "La Sapienza", P.le Aldo Moro 5, 00185, Rome, Italy*

⁹*UMR PNCA. AgroParisTech, INRA. Université Paris-Saclay, 75005 Paris, France*

Corresponding author: puneet.mishra@wur.nl; jean-michel.roger@inrae.fr

68 Abstract

69 In recent years, due to advances in sensor technology, multi-modal measurement of process
70 and products properties has become easier. However, multi-modal measurements are only of
71 use if the data from adding new sensors is worthwhile, especially in the case of industrial
72 applications where financial justification is needed for new sensor purchase and integration,
73 and if the multi-modal data generated can be properly utilised. Several multi-block methods
74 have been developed to do this; however, their use is largely limited to chemometricians, and
75 non-experts have little experience with such methods. To deal with this, we present the first
76 version of a MATLAB-based graphical user interface (GUI) for multi-block data analysis
77 (MBA), capable of performing data visualisation, regression, classification and variable
78 selection for up to 4 different sensors. The MBA-GUI can also be used to implement a recent
79 technique called sequential pre-processing through orthogonalization (SPORT). Data sets are
80 supplied to demonstrate how to use the MBA-GUI. In summary, the developed GUI makes
81 the implementation of multi-block data analysis easier, so that it could be used also by
82 practitioners with no programming skills or unfamiliar with the MATLAB environment. The
83 fully functional GUI can be downloaded from ([https://github.com/puneetmishra2/Multi-](https://github.com/puneetmishra2/Multi-block.git)
84 [block.git](https://github.com/puneetmishra2/Multi-block.git)) and can be either installed to run in the MATLAB environment or as a standalone
85 executable program. The GUI can also be used for analysis of a single block of data (standard
86 chemometrics).

87 *Keywords: data fusion; multi-sensor; chemometrics; graphical user interface*

88 Introduction

89 Sensing technologies play a major role in chemical industries where they are implemented to
90 monitor and optimise process and product properties [1]. Sensing technologies do this by
91 rapid estimation of key critical quality attributes of the process and products. However,

92 sometimes the products or processes are so complex that a single technique fails to obtain
93 sufficient information about the samples. One such case, in the framework of process
94 monitoring applications, is the use of Raman and mid-infrared (MIR) spectroscopy. Both
95 techniques are complementary to one another as they are sensitive to different vibration
96 modes corresponding to molecular groups. Furthermore, they complement one another's
97 drawbacks as Raman signal may get influenced by fluorescence, but it will work with high
98 moisture samples while MIR spectroscopy will be affected by the presence of moisture, but it
99 will work well with fluorescent samples. In such a case, data from both techniques can be
100 utilised in a complementary way: a combination of MIR and Raman spectroscopies could
101 yield better results as each one will compensate for the drawbacks of the other.

102 Innovations in measurement technologies such as combining multi-spectral techniques for
103 non-destructive estimation of process and product properties is now a major research domain
104 requiring multi-block data analysis techniques. A possible application could be the
105 integration of several different PATs such as MIRS, Raman, near-infrared spectroscopy
106 (NIRS) and fluorescence spectroscopy (FS) into a single measurement probe (Fig. 1A). In the
107 case of a process monitoring application, such a combined probe could be inserted into the
108 process vessel with signals being obtained from multiple sensors, thus enabling the recording
109 of continuous data from multiple techniques (Fig. 1B). However, this will require not only
110 advances in hardware but also in chemometric procedures, such as variable selection
111 methodologies for identification of key spectral regions, and data fusion algorithms for the
112 combination of data coming from multiple techniques.

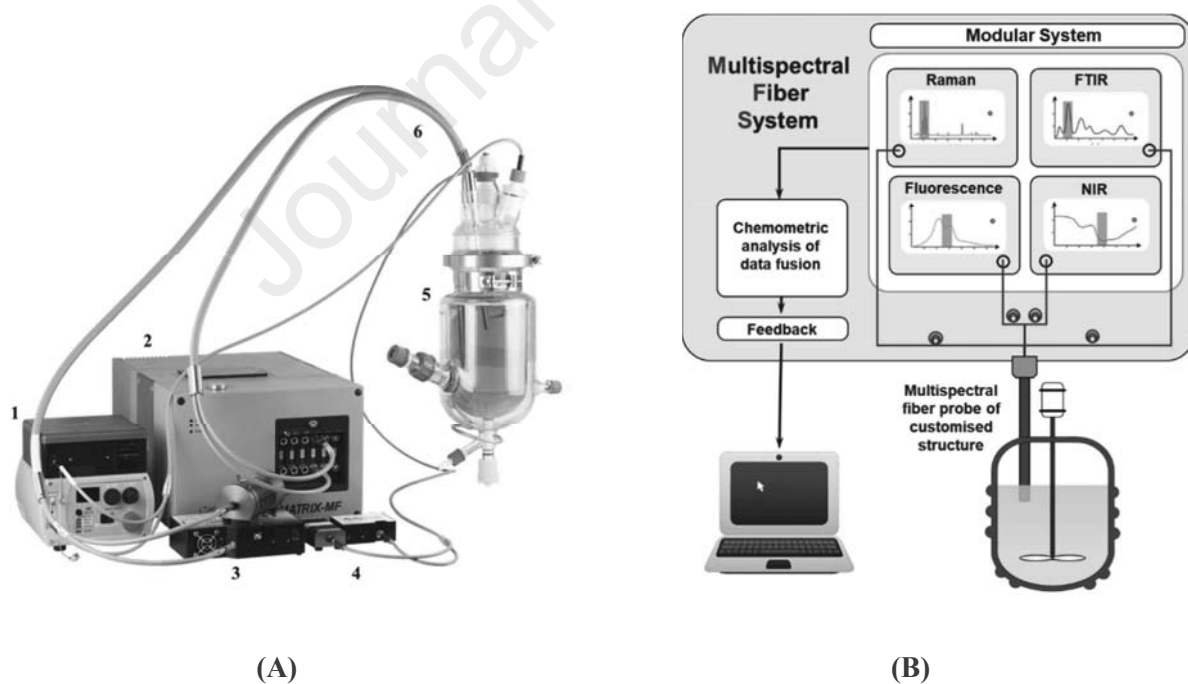
113 In some cases, it could be preferable for the data analysis to be performed using a sequential
114 approach so as to highlight the added value coming from including new measurement modes,
115 especially for industrial applications where financial justification is needed for new sensor
116 purchase and integration. Data fusion is a feasible approach to combine all the information

117 from multiple sensors [2-7]. Data fusion can be dealt with in different ways, depending on the
118 scientific domain. In the machine learning domain, data fusion can be performed at three
119 levels, i.e. low, mid and high-level [8]. Low-level data fusion implies taking all the data
120 together and performing an analysis in a similar way as for data from a single sensor. Mid-
121 level fusion involves some preliminary data refining step where interesting features are
122 extracted from each dataset and subsequently fused together. High-level data fusion is the
123 fusion of the conclusions drawn from the output of the models created on each dataset, using
124 decision rules such as majority voting or averaging. The main aim of data fusion in the
125 machine learning domain is to improve the model accuracy with less attention being paid to
126 the background process. However, in the chemometrics domain, data fusion, or multi-block
127 data analysis, aims not only to improve the model accuracy but also to have a better
128 understanding of the underlying characteristics involved. The aim is to identify the common
129 and the specific hidden factors between and within multiple blocks of measurements, and to
130 subsequently use them to build explainable and explanatory models [9-14].

131 There are two main tasks that need to be accomplished in multi-block data analysis, i.e., to
132 enhance the data visualisation and improve the predictive performance of models. For data
133 visualisation, a summary and comparison of methods can be found in [9, 11, 12]. Recently, a
134 new data visualisation approach was presented for exploration of designed experiments in a
135 multi-block scenario [10]. For multi-block predictive modelling, partial least squares
136 regression-based methods are summarised in [15]. Multi-block methods have also been
137 extended to incorporate variable selection [16, 17], which is important, for example, when
138 different spectral sensors are used to study the same set of objects. In such a case, variable
139 selection looks for subsets of variables that are important in each of the spectral techniques,
140 which can then be useful to have a better understanding of the process, and to orient the
141 development of cheap multi-spectral sensors. Multi-block methods are also emerging to

142 perform fusion of data of very different types, such as when fusing a 3 way data array (3D
 143 tensor) with a 2-way matrix (2D tensor) [18, 19].

144 In the present work, the first version of a MATLAB-based graphical user interface for multi-
 145 block data analysis (MBA-GUI) is presented. The MBA-GUI can perform data visualisation,
 146 sequential regression and variable selection on up to 4 different data sources. However, the
 147 algorithms implemented in the GUI are not limited to 4 data sources and can be used for any
 148 number of data sources. The MBA-GUI can also be used to implement a technique called
 149 sequential pre-processing through orthogonalization (SPORT). The MBA-GUI can also be
 150 used for a single block scenario. In addition, cases are presented showing how to use the
 151 MBA-GUI for data visualisation, regression, classification, variable selection and SPORT in
 152 the multi-block scenario. This is the first version of the MBA-GUI and, as multi-block
 153 analysis methods progress, the toolbox will be updated to incorporate new algorithms.



156 *Fig. 1: Scheme of multispectral fibre system (figure courtesy of Art Photonics GmbH,*
 157 *Germany). (A) Raman system (1); FTIR absorption System (2); NIR reflection System (3);*

158 *fluorescent System (4); chemical Reactor (5); fibre optic probes (6), and (B) A schematic of*
159 *the multi-block data generated in a four blocks scenario.*

160 Similar works

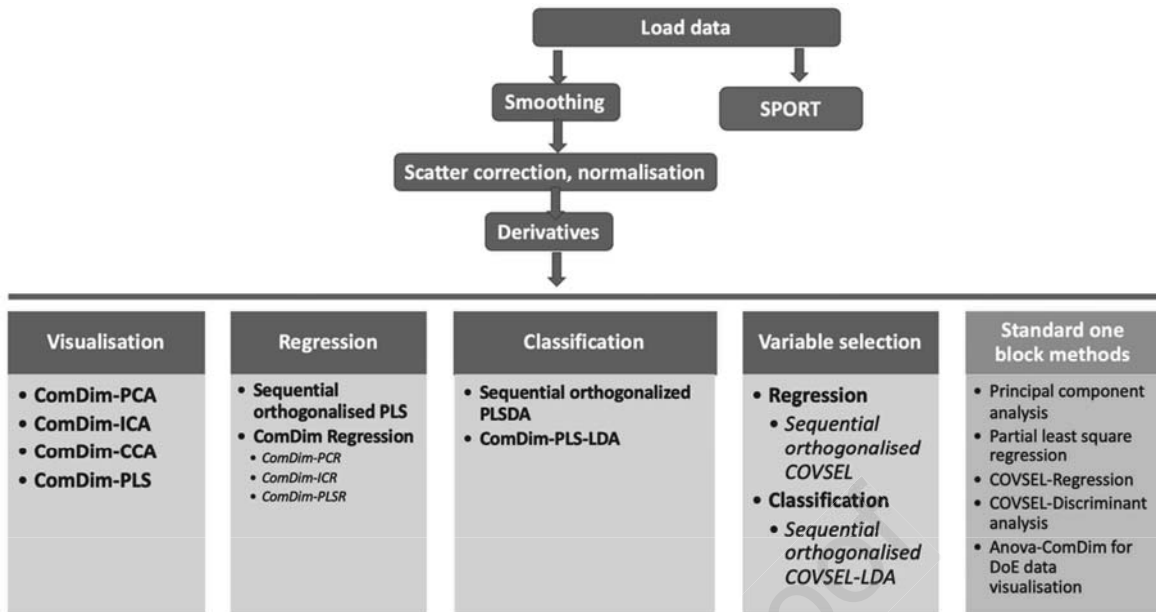
161

162 When looking at the available resources to perform multi-block data analysis, three main
163 toolboxes can actually be found. multi-blockThe first one is the multi-block toolbox by the
164 University of Copenhagen, Denmark
165 (<http://www.models.life.ku.dk/~courses/MBtoolbox/mbtmain.htm>), which, having been last
166 revised in 2001, focuses only on the two data fusion approaches which were most popular in
167 that period, i.e, multi-block principal component analysis and multi-block partial least
168 squares regression. The second one is the ‘multi-block regression by parallel and sequential
169 partial least-squares regression’ toolbox by NOFIMA [20]. Both these toolboxes provide
170 command line functionalities (which may be difficult to implement for people unfamiliar
171 with the Matlab environment) and, anyway, consist of a limited number of tools. There is also
172 a basic GUI available for performing multi-block component analysis in the domain of
173 Behaviour research [21]. However, the GUI can only perform multi-block component
174 analysis for data visualisation. Therefore, there exists a need for a GUI which has updated
175 tools and a complete set of functionalities to perform multi-block data analysis.

176 Software description

177 The MBA-GUI was built utilising the application builder in MATLAB version 2018b
178 (Natick, MA, USA). The application can be downloaded and installed in MATLAB
179 (preferred version 2018b or higher) or can be used as a stand-alone executable or can be run
180 through the ‘.mlapp’ files in MATLAB command line. If user does not have MATLAB
181 version 2018b or greater then it is recommended to install free MATLAB runtime tool and
182 run the app as standalone. All the executable and MATLAB function can be downloaded

183 from (<https://github.com/puneetmishra2/Multi-block.git>). In the GitHub repository, the
184 standalone toolbox executable files can be downloaded as '*Multi-block_toolbox.zip*' and the
185 function for running the tools in command line as '*Toolbox.zip*'. The dataset demonstrated in
186 this article can be downloaded as '*Dataset.zip*' from the same GitHub repository. All three
187 files are available in the link (10/June/2020). To run the toolbox from command line user
188 should use the toolbox folder as the current folder and type T1 on the command line which
189 will start the main GUI interface. The user should put the password: 'welovedata' without
190 comma and click run. Then user can load data and run the analysis. See also supplementary
191 file to have a visual understating on how to download and setup the GUI. The GUI supports
192 data format of .csv, .xlsx and .mat. A summary of the functionalities is presented in Fig. 2. In
193 summary, the toolbox has options for loading data, three levels of pre-processing, i.e.,
194 smoothing, scatter correction and normalisation, and derivative estimation, as well as multi-
195 block data visualisation, regression, classification, variable selection and SPORT. Multi-
196 block variable selection methods are available for both regression and classification cases.
197 Two main types of regression and classification are available, i.e., sequential
198 orthogonalization [15], and common components and specific weights analysis (CCSWA),
199 aka common dimensions (ComDim) [22].



200

201 *Fig. 2: Schematic of the tasks that can be performed with the MBA-GUI.*

202

203 Software architecture and brief mathematical background of
 204 techniques available

205

206 Pre-processing

207 Data pre-processing is an important step to clean and homogenise the data prior to analysis.

208 Proper data pre-processing can improve data modelling dramatically. There can be multiple

209 steps in data pre-processing such as smoothing, scatter correction, normalisation and many

210 others [23-25]. In the toolbox, we have provided a collection of common pre-processing

211 methods.

212 Smoothing operations

213 Data smoothing reduces high-frequency noise from datasets prior to modelling. In the

214 toolbox, several techniques are provided for performing smoothing in the variable domain.

215 Three window-based smoothing techniques, i.e., Savitzky-Golay (SAVGOL), moving

216 average and moving polynomial are provided. Further, two data decomposition and

217 reconstruction techniques, i.e., principal components reconstruction and independent

218 components reconstruction are provided. In this toolbox, it is up to the user to choose a
219 technique and decide, based on the model performance, which is best for their data. Pre-
220 processing can also be explored automatically with the SPORT approach. All the spectral
221 smoothing techniques are implemented using the codes explained in [23].

222 Scatter correction, Baseline correction and Normalisation

223 Multivariate data, and especially spectral data, suffer from a range of physical and chemical
224 factors leading to baseline, additive and multiplicative effects. Prior to data modelling, it is
225 always recommended to perform correction and normalisation of the data. However,
226 depending on the data, the correction or normalisation method may be different. In the
227 toolbox, the user may select several scatter and spectral normalisation techniques, including
228 detrending [26], offset correction, multiplicative scatter correction [27], spline correction
229 (where spline fitting is used to approximate the baseline which is then subtracted from the
230 signal), asymmetric least-squares (ALS) correction [28], standard normal variate (SNV) [26],
231 variable sorting for normalisation (VSN) [29], probabilistic quotient normalisation (PQN)
232 [30], robust normal variate (RNV) [31], logarithm transform, autoscaling, 1st derivative
233 (Savitzky-Golay), 2nd derivative (Savitzky-Golay), min-max, norm, range and max
234 correction. All the correction and normalisation methods were implemented using the codes
235 presented in [23]. It is recommended to perform the smoothing step, if required, before
236 baseline correction and normalisation as these techniques can be affected by high-frequency
237 noise. Detailed understanding of pre-processing methods in chemometrics can be found in
238 [23-25].

239 Derivatives

240 Many normalisation and baseline correction techniques can remove effects like baseline shift,
241 additive and multiplicative scatter. Derivatives, however, are also able to reveal underlying
242 peaks. Therefore, the user can choose to perform 1st and 2nd derivatives or define the order of

243 the derivative as a 3rd pre-processing step. The algorithm for this operation is based on the
244 codes provided in [23].

245 Pattern recognition and data visualisation tool

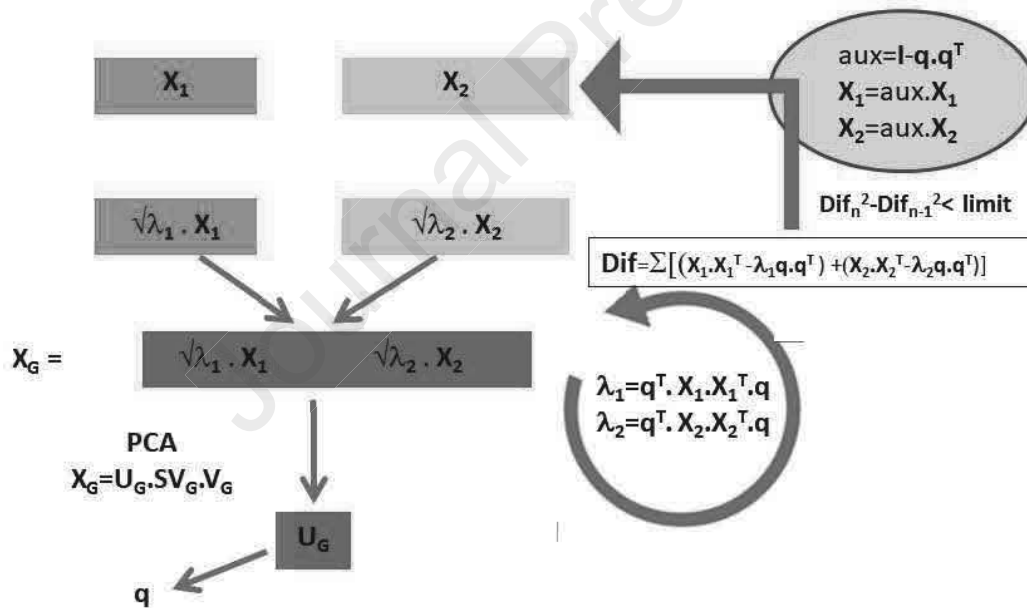
246 Data visualisation is the primary step prior to data analysis to gain an insight into the data
247 [32]. Many factors, such as the presence of outliers, any groupings of samples or any strange
248 patterns, can be detected by visualising the data. Based on the type of data, its visualisation
249 can be very simple or very complex. If the data comprises only a few variables then major
250 patterns can be observed via 1D, 2D and 3D plots. However, if the data consists of many
251 variables, then dimension reduction techniques are used to transform the high-dimensional
252 data into a space defined by interesting properties such as the variance. Once projected into a
253 subspace (dimensionality reduction), the samples can be visualised via 1D, 2D and 3D plots
254 together with plots of the transformed variables. Since spectral data usually contain many
255 variables, ranging from hundreds to thousands depending on the spectral resolution, data
256 transformation is almost always required [33]. In the toolbox, we implement a range of data
257 transformation tools to enhance data visualisation task in such cases. The techniques are
258 mainly based on capturing the major sources of variance in the sample domain. The
259 multivariate data transformation techniques are specific to the multi-block scenario as they
260 can highlight the contributions from the different blocks [9].

261 ComDim- based methods

262

263 ComDim belongs to the family of multi-block methods, which aim to extract the global
264 components that highlight the important dimensions as well as the local (block specific)
265 components [22]. Originally, the ComDim method was developed with the name common
266 components and specific weights analysis (CCSWA) [34]. ComDim extracts the global and
267 local components from multiple blocks of data in a sequential way. Each data block has a
268 specific contribution to each common component which is called its '*Salience*'. ComDim

269 starts by normalising each block, X_i , by its Frobenius' norm so that they all have the same
 270 total variance. Details regarding the algorithm can be found in [22, 35]. Application of the
 271 original ComDim (which, for reasons which will become apparent later, we will call
 272 ComDim-PCA here) in the case of 2-blocks is presented in Fig. 3. The common components
 273 (CCs) are extracted sequentially in an iterative fashion, by extracting the eigenvector
 274 associated with the largest eigenvalue from a matrix W , which is the created by
 275 concatenating the weighted individual matrices, X_i . The weightings (the *saliences*) are
 276 initially all set to 1, but they are recalculated during the iterations to reflect the contribution
 277 of each block to the dispersion of the individuals along that CC. After one CC is extracted,
 278 the X_i matrices are deflated, and the procedure is repeated to obtain the next CC.



279

280 *Fig. 3: The ComDim algorithm applied in the 2-block case ([36]): In the first step of the*
 281 *algorithm, X_1 and X_2 have been normalised.*

282

283 In the MBA-GUI, ComDim is provided in several variants, i.e., common dimensions-
 284 principal components analysis (ComDim-PCA), common dimensions-independent

285 components analysis (ComDim-ICA), common dimensions-common components analysis
286 (ComDim-CCA) (unsupervised decomposition methods), common dimensions-partial least
287 squares-independent components analysis (ComDim-PLS-ICA) (a semi-oriented
288 decomposition) and common dimensions-partial least squares (ComDim-PLS) (an oriented
289 decomposition). The difference between the first three ComDim variants is the way the
290 concatenated matrix of blocks is decomposed, i.e., using singular value decomposition
291 (SVD), independent components analysis (ICA) or common components analysis (CCA).
292 ComDim-PLS-ICA is a slightly more complex version of ComDim-ICA as it replaces the
293 PCA decomposition of \mathbf{W} by an ICA algorithm where the initial estimates of the independent
294 components have been determined using a partial least squares (PLS) regression. In
295 supervised ComDim-PLS, the PCA decomposition within ComDim is simply replaced by a
296 PLS regression. The iterative procedure is identical in all cases, but the resulting CCs differ
297 somewhat as a function of the criteria that determine the different decompositions of \mathbf{W} .

298 Regression

299

300 SO-PLS Regression

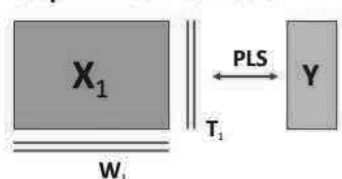
301

302 Sequential and orthogonalized PLS (SO-PLS) regression belongs to the family of multi-block
303 PLS methods; the centrepiece of the method is the orthogonalization step, which ensures the
304 removal of redundancies among modelled data blocks [15]. In SO-PLS regression, the
305 extraction of information is sequential, meaning that the aim is to incorporate blocks of data
306 one at a time and to assess the incremental contribution. A PLS regression is calculated
307 between the first block \mathbf{X}_1 and \mathbf{Y} , yielding scores \mathbf{T}_1 . Then, all the remaining blocks $\mathbf{X}_2, \dots, \mathbf{X}_k$
308 and the \mathbf{Y} block are orthogonalized with regards to \mathbf{T}_1 . Then, the process is repeated on the
309 second block, and so on for all the blocks. A scheme showing sequential extraction of
310 information by SO-PLS regression is presented in Fig. 4. The major advantages of SO-PLS

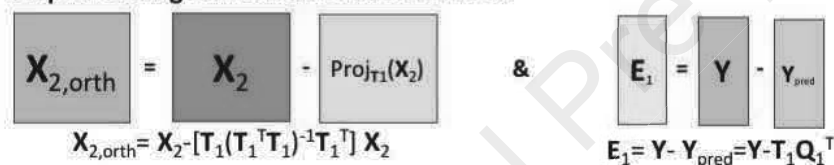
311 are linked to orthogonalization, which removes redundant information, and to its sequential
 312 nature, which allows the interpretation of the incremental contributions provided by each data
 313 block. For more details, the reader is directed to [15, 37, 38]. The SO-PLS regression
 314 function integrated into the toolbox is the freely available one at
 315 <https://www.chem.uniroma1.it/romechemometrics/research/algorithms/>

316

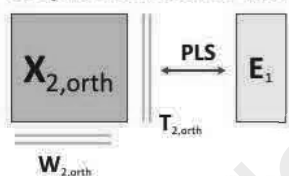
Step 1: First PLS model



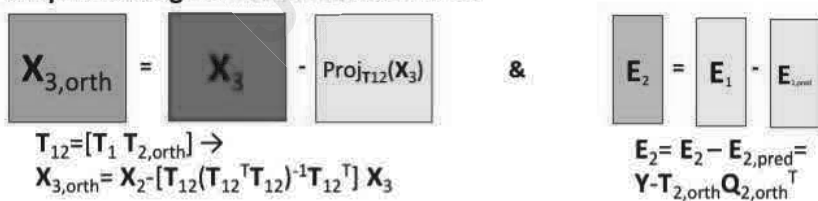
Step 2: Orthogonalization of second block



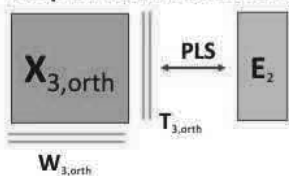
Step 3: Second PLS model



Step 4: Orthogonalization of third block



Step 5: Third PLS model



$$\text{Global model: } Y_{\text{pred}} = T_1 Q_1^T + T_{2,\text{orth}} Q_{2,\text{orth}}^T + T_{3,\text{orth}} Q_{3,\text{orth}}^T$$

317

318 *Fig. 4: A scheme presenting the sequential orthogonalized partial least squares (SO-PLS)*

319 *regression method.*

320 ComDim regression

321 The aim of ComDim regression methods is to sequentially extract the global component from
322 multi-block data and to later use the scores to perform the regression on a y vector. To
323 perform multi-block regression, four ComDim variants are integrated, i.e., ComDim-
324 Principal component regression (ComDim-PCR), ComDim-Independent component
325 regression (ComDim-ICR), ComDim-Partial least squares regression (ComDim-PLSR) and
326 ComDim-Partial least squares – Independent component regression (ComDim-PLS-ICR).

327 The difference between these regression methods is simply that the global components that
328 are used were obtained using the different multi-block ComDim methods presented above.

329 ComDim-PCR and ComDim-ICR use global scores that were extracted in an unsupervised
330 way to perform multi-linear regression (MLR). ComDim-PLSR and ComDim-PLS-ICR use
331 scores that were extracted in a supervised way by the PLS step within ComDim or within the
332 ICA which is nested inside ComDim. To use new data, the ComDim models developed on
333 the calibration dataset are used to calculate the scores for the new dataset and these are then
334 introduced into the trained MLR model for the prediction. MLR regression on the ComDim
335 scores is performed using in-house codes freely downloadable at
336 <https://www.chem.uniroma1.it/romechemometrics/research/algorithms/>.

337

338 Classification

339 SO-PLS-LDA

340 Sequential orthogonalized partial least squares linear discriminant analysis (SO-PLS-LDA) is
341 the natural extension of SO-PLS to the classification field [39]. To create a SO-PLS-LDA
342 model, a SO-PLS model is first created using a dummy class matrix as Y , and then LDA can
343 be applied either to the predicted Y or to the concatenated scores. Due to it being closely
344 related to the SO-PLS algorithm, the steps of the SO-PLS-LDA approach are the same as
345 those already sketched in Fig. 4; the main differences being that the response matrix Y should

346 be coded to account for class membership and that, as a further step, linear discriminant
347 analysis is applied to either the predicted response or the concatenated scores block. The
348 function integrated into the GUI is the one freely downloadable at
349 <https://www.chem.uniroma1.it/romechemometrics/research/algorithms/>.

350 ComDim - based linear discriminant analysis

351

352 The aim of ComDim based classification methods is to sequentially extract global
353 components from multi-block data and then to use the scores to perform a linear discriminant
354 analysis (LDA). In the GUI, two variants of ComDim-LDA are included. ComDim-PLS2-
355 LDA simply uses a PLS2 inside ComDim to orient the decomposition of the \mathbf{W} matrix. On
356 the other hand, ComDim-PLS2-ICA-LDA uses PLS2 to orient the extraction of the
357 independent components from the \mathbf{W} matrix by ICA, which has replaced PCA nested inside
358 ComDim. In both cases, the scores can be directly used for LDA. Test set prediction is
359 performed by first transforming the new data to the same space using the ComDim
360 calibration model and then inputting into the trained LDA model.

361 Variable selection

362 SO-CovSel

363

364 SO-CovSel is a multi-block variable selection technique recently developed by [16]. The SO-
365 CovSel is an extension of the CovSel technique to the multi-block scenario [40]. CovSel
366 extracts the ' k ' variables from a matrix \mathbf{X} that are most correlated to \mathbf{Y} and independent to
367 each other. SO-CovSel performs CovSel in a sequential orthogonalized way. It works as SO-
368 PLS, replacing the PLS-scores by the COVSEL selected variables. CovSel in such an
369 approach performs variable selection so that extraction of the variables from the consecutive
370 block improves the model. SO-CovSel is designed for both regression as well as
371 classification cases. SO-CovSel for regression and SO-CovSel-LDA are integrated into the

372 MBA-GUI utilising in-house codes freely downloadable at
373 <https://www.chem.uniroma1.it/romechemometrics/research/algorithms/>.

374 SPORT

375 Novel use of multi-block methods can also be understood as boosting of different pre-
376 processing techniques. Such a methodology called SPORT was recently developed by [41],
377 where the sequential orthogonalization approach was used to fuse the information from
378 different pre-processing techniques. Thus, instead of choosing between pre-treatments,
379 SPORT allows us to make optimal use of the advantages of all pre-treatments. In this GUI, up
380 to four different pre-processing techniques can be used for boosting. Boosting with sequential
381 orthogonalization is a recent approach and has proven to be of high value to improve the
382 prediction accuracies of the models. To perform SPORT, the same data can be loaded in four
383 blocks and then different pre-processing can be applied to them. Further, based on the case of
384 regression/classification, SO-PLS/SO-PLS-DA is used. More details on SPORT can be found
385 in [41].

386

387 Standard one block chemometric analysis

388 Apart from multi-block analysis, the MBA-GUI provides an option to perform standard
389 chemometric analyses. For a single block, the MBA-GUI has options to do PCA, PLS as well
390 as CovSel variable selection for regression and discriminant analysis. The one-block analysis
391 will automatically start when the MBA-GUI detects that only one block of data is loaded.

392 Datasets for MBA-GUI demonstration

393 Use of the MBA-GUI is demonstrated with three datasets. All the three datasets can be
394 accessed in the same GitHub repository. The first dataset for data visualisation consists of the
395 FTIR spectra of olive oils of 4 different origins (4 classes) measured in transmission mode
396 [42, 43]. The second dataset to demonstrate the regression and variable selection task relates

397 to dry matter prediction in olive fruits measured with a portable spectrometer in diffuse
398 reflectance. The reference dry matter was measured by the hot air oven method by noting the
399 fresh and the dry weight of the samples. More details on the olive fruits data can be found in
400 [44]. The third dataset to demonstrate the classification task consists of the NIR spectra of
401 mayonnaise samples made from oils of 6 different origins. The classification task can be
402 understood as a 6-class problem. The mayonnaise dataset was obtained from the official
403 website of ChemHouse (www.chemproject.org). To make the data fit for multi-block
404 analysis, each dataset was split into two blocks in the spectral domain. A further description
405 of the datasets is provided in Table 1.

406

407

Table 1: Description of the datasets used in the demonstration.

| Samples | Task | Calibration | Calibration | Test | Test | Y calibration | Y test |
|--|---|--------------------|--------------------|-----------------------|-----------------------|------------------|--------|
| | | Block 1 | Block 2 | Block 1 | Block 2 | | |
| Olive oils (Discrete response) | Visualisation | 83×300 | 83×270 | | | 83×1 | |
| | | (798-1375 nm) | (1377-1896 nm) | | | | |
| Olive fruits (Continuous response) | Regression and variable selection | 350×75 | 350×60 | 145×75 | 145×60 | 350×1 | 145×1 |
| | | (708-930 nm) | (933-1113 nm) | (708- 930 nm) | (933- 1113 nm) | | |
| Mayonnaise (Discrete response) | Classification | 72×150 | 72×201 | 72×150 | 72×201 | 72×1 | 72×1 |
| | | (1100- 1696 nm) | (1700- 2500 nm) | (1100- 1696 nm) | (1700- 2500 nm) | | |

408

409

410

411

412 Operating procedure and demo analysis

413 The operating procedures are presented in separate sections to demonstrate the use of the
414 MBA-GUI for the analysis of different datasets so that a person with minimal experience can
415 repeat the analysis and use the GUI in their day-to-day multi-block data analysis tasks. The
416 sections are data loading and pre-processing, data visualisation, regression, classification,
417 variable selection and SPORT. All the figures presented in the analysis come directly from
418 the MBA-GUI.

419 Data loading and pre-processing

420

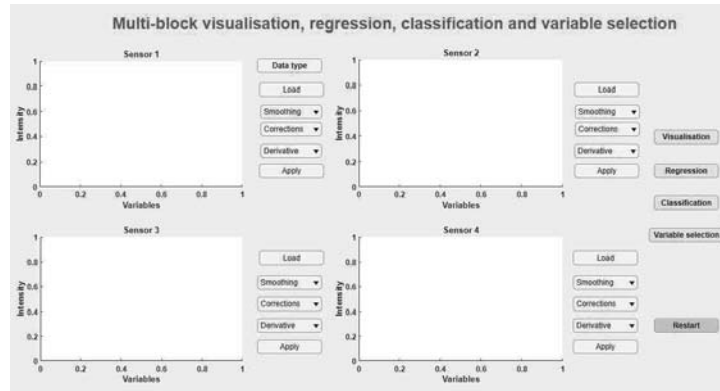
421 Fig. 5A shows the MBA-GUI interface for loading datasets. Currently, two data formats can
422 be loaded, i.e., .xls and .csv. It is currently possible to load up to 4 blocks of data. The
423 datasets can either come from several sensors or can be the same dataset imported multiple
424 times to perform SPORT fusion. The datasets should be loaded in a logical order since
425 sequential loading of the dataset may improve the performance of sequential methods where
426 the natural order of blocks can provide meaningful insights into the datasets. Once the data
427 are loaded, the figure in the MBA-GUI will be updated to show the data. Fig. 5B shows an
428 example where two data blocks are loaded. Each data block can be separately pre-processed.
429 The pre-processing can be performed in three steps - smoothing, normalisation or scatter
430 correction, and derivatives. Derivatives are of particular use for NIR data where they can
431 reveal underlying peaks. An example of pre-processing is shown in Fig. 5B. The pre-
432 processing choice in Fig. 5C involved SAVGOL smoothing, SNV normalisation and no-
433 derivative. Fig. 5C shows that once the pre-processing is selected the MBA-GUI will show
434 the new pre-processed spectra. The user is free to explore multiple pre-processing methods by
435 visualising how pre-processing affects the spectra. After pre-processing, the data are ready
436 for multi-block analysis. There are five push-button options provided - visualisation,

437 regression, classification, variable selection and SPORT. There is also an option to restart the
438 complete analysis by clearing all the previous data and operation logs. A point to note is that
439 if there is only one data block, the MBA-GUI options will take the user to standard
440 chemometric analysis where analyses such as PCA, PLS and variable selection can be
441 performed. Due to limited space, the one block chemometric analysis is not presented in this
442 article.

443

444

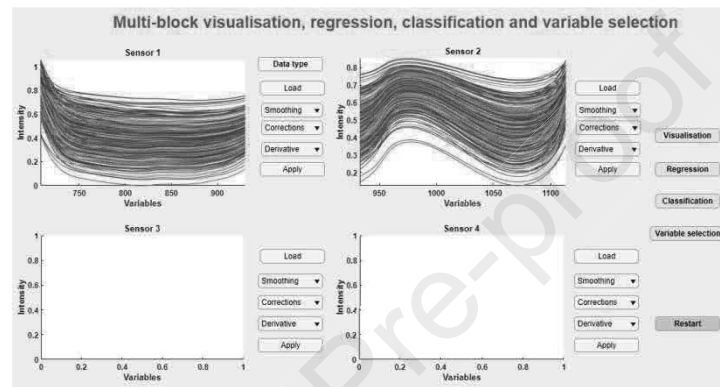
Journal Pre-proof



445

446

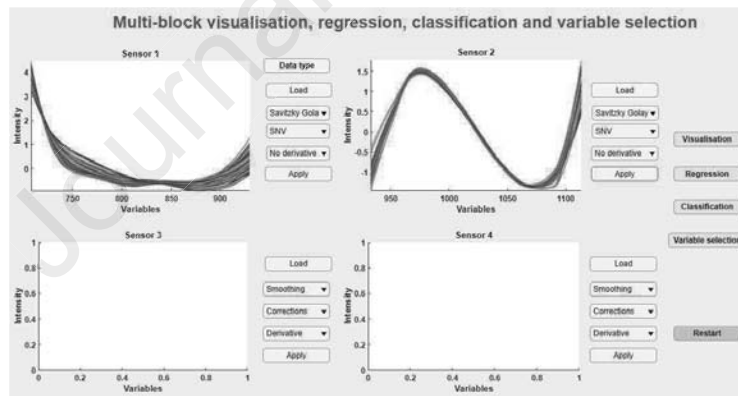
(A)



447

448

(B)



449

450

(C)

451 *Fig. 5: MBA-GUI interface for data loading. (A) Up to four different data blocks can be*
 452 *loaded and analysed. (B) GUI interface once the data are loaded using the Load button. The*
 453 *spectra can be loaded in a sequential order i.e. 1, 2, 3 and 4. Once the data are loaded, all of*
 454 *the pre-processing methods can be applied. (C) MBA-GUI interface after pre-processing.*
 455 *Once the pre-processing option is selected and the Apply button is pressed, the figures will*
 456 *contain the pre-processed spectra.*

457 Data visualisation

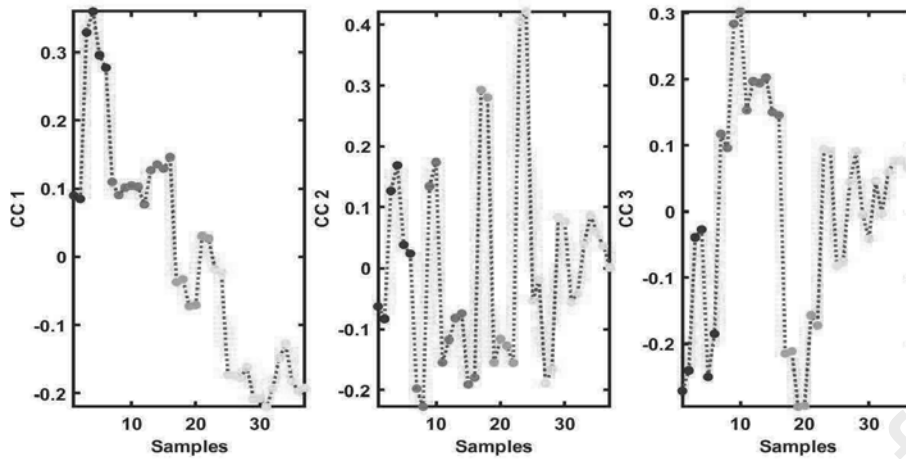
458

459 Multi-block data visualisation is a challenging task where the interest is not only in the
460 extraction of latent variables but also in how different blocks are linked with each other or
461 contribute to the extraction of the global LVs. In the present case, the objective is also to limit
462 the extraction of redundant information from the different blocks. In the MBA-GUI, multiple
463 multi-block data visualisation methods are implemented. The presentation here will be
464 limited to the use of ComDim-PLS, which is a supervised common dimension extraction
465 method requiring a response variable to orient the decomposition of the blocks. It will be
466 applied to the olive oil dataset.

467 Fig. 6A, 6B and 6C show the scores when a model with three common components (CCs) is
468 selected. Different coloured points in the figure indicate samples belonging to different
469 classes as defined by the Y vector. It can be seen that with 1st and the 3rd CCs, a clear
470 distinction of different classes is possible. Fig. 6D, 6E and 6F show the saliences for each
471 block that contributed to the CCs. In this case, the higher saliences for block 2 show that CC1
472 was dominated by the information from that block, whereas CC2 and CC3 were dominated
473 by the information from block 1. Since they were normalised, each block has a total salience
474 of 1 which means, in the present case, a total salience of two. In the figure, we can see the
475 amount of salience extracted by each block for each CC as well as the total salience extracted
476 by each CC. It is also possible to see that the total amount extracted from the 2 blocks by 3
477 CCs is about 1.75. Fig. 6G, 6H, 6I show the loadings for each CC presented in two figures
478 each, corresponding to the two-blocks. Such loading plots can be used to understand which
479 variables are of interest.

480

Global scores on common components with samples



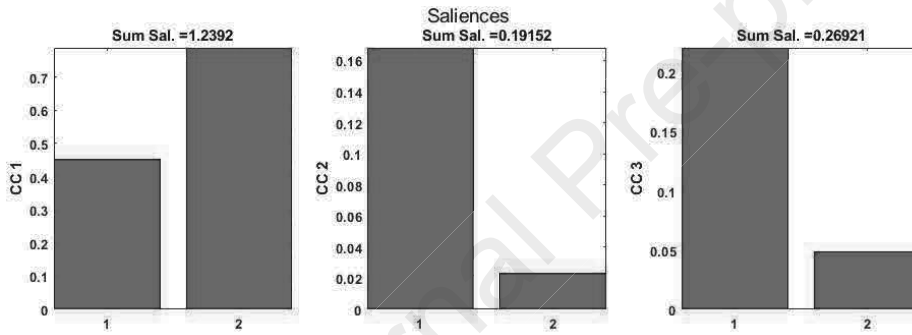
481

(A).

(B).

(C).

482



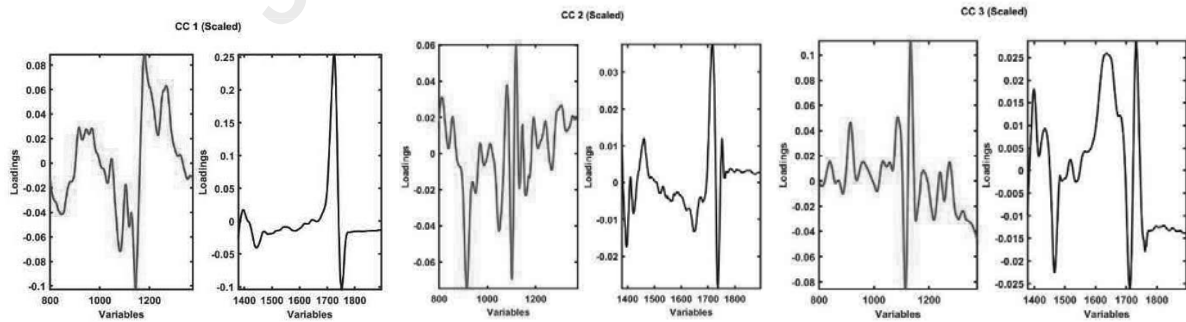
483

(D).

(E).

(F).

484



485

(G).

(H).

(I).

486

487 Fig. 6: Output from ComDim-PLS performed on the olive oil dataset. (A) Scores from CC1,

488 (B) scores from CC2, (C) scores from CC3, (D) saliences for CC1, (E) saliences for CC2, (F)

489 saliences for CC3, (G) loading for CC1, (H) loading for CC2, and (I) loading for CC3.

490 Regression

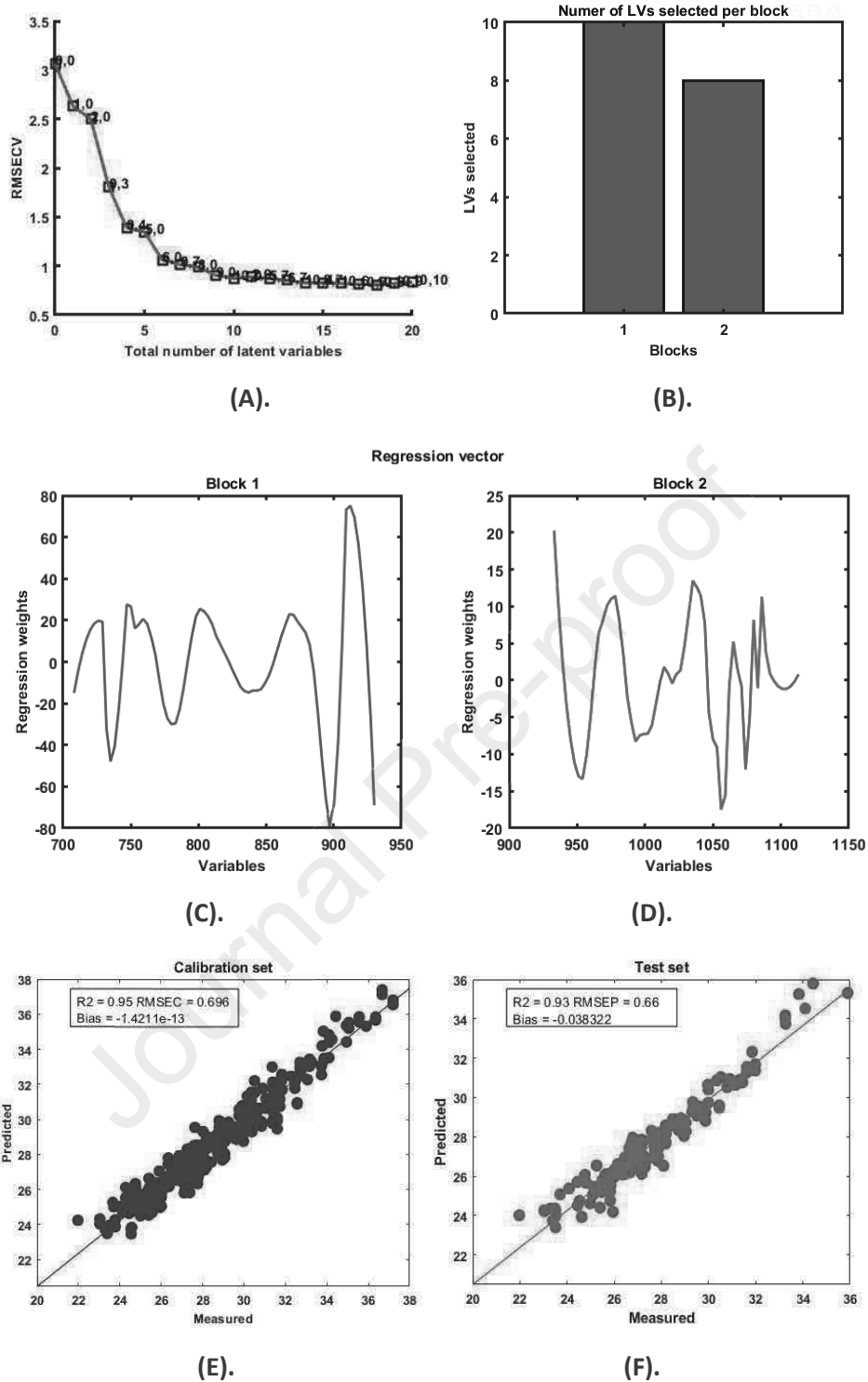
491

492 Multi-block regression is useful when multiple sensors are integrated to improve the
493 predictive performance of the model. In the GUI, two types of multi-block regression are
494 implemented, i.e., SO-PLS and MLR, both based on the ComDim-PLS scores. An example
495 using SO-PLS is presented here. Fig. 7 presents the output of the SO-PLS analysis, where
496 Fig. 8A shows the cross-validation error and Fig. 7B shows the number of LVs selected for
497 each data block, i.e., 10 LVs for block 1 and 8 LVs for block 2. Fig. 7C and 7D show the
498 final regression vector based on the LVs extracted from blocks 1 (Fig. 7C) and 2 (Fig. 7D).
499 The results from the SO-PLS modelling are presented as the calibration plots in Fig. 7E and
500 7F. The R^2 for calibration and prediction were 0.95 and 0.93, respectively. The calibration and
501 prediction errors were 0.69 and 0.66 % dried matter, respectively. It should be noted that
502 when only one block of data is used, the R^2 is lower and the error is higher (as shown in
503 single block analysis in Fig. 10), showing the benefit of multi-block regression.

504

505

506



507

508

509

510

511

512

513 *Fig. 7: The output from SO-PLS regression performed on the olive fruit dry matter dataset.*514 *(A) Cross-validation error, (B) number of LVs selected from each block, (C) regression*515 *vector from block 1, (D) regression vector from block 2, (E) calibration set results, and (F)*516 *test set results.*

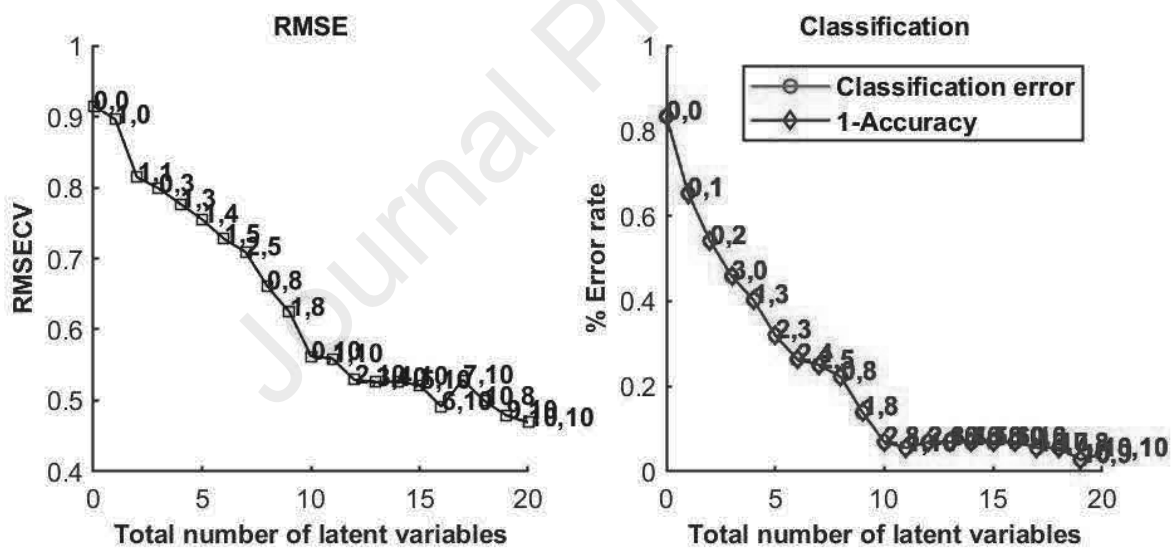
Journal Pre-proof

518 Classification

519

520 Multi-block classification involves the use of multiple sensor data to improve the
 521 classification accuracies. In the MBA-GUI toolbox, several multi-block classification
 522 techniques are implemented. Here an example of SO-PLS-LDA on the mayonnaise dataset is
 523 presented.

524 Fig. 8A and 8B shows the RMSE and the error evolution, respectively, as a function of the
 525 number of LVs. The RMSE and error plot were used for automatic selection of the number of
 526 latent variables for the two data blocks. The classification results from the calibration and test
 527 sets are presented in Fig. 8C and 8D respectively. An overall prediction accuracy of 93% was
 528 obtained.

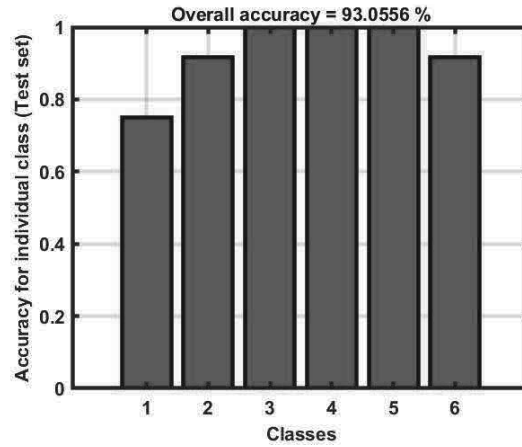
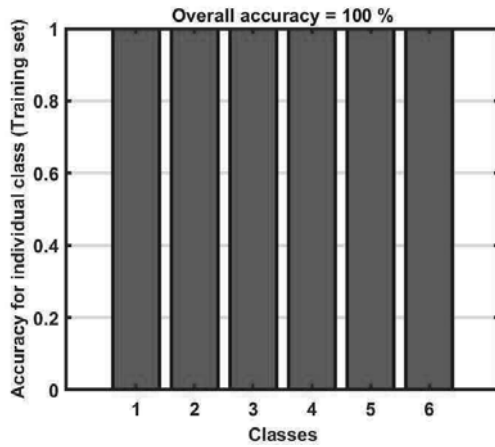


529

530

531

532



533

534

(C).

(D).

535

536 *Fig. 8: Results from SO-PLS-LDA of the mayonnaise dataset. (A) RMSE as a function of the*537 *number of variables selected, (B) error rate as function of number of variables selected, (C)*538 *accuracy on calibration data, and (D) accuracy on test data.*

539

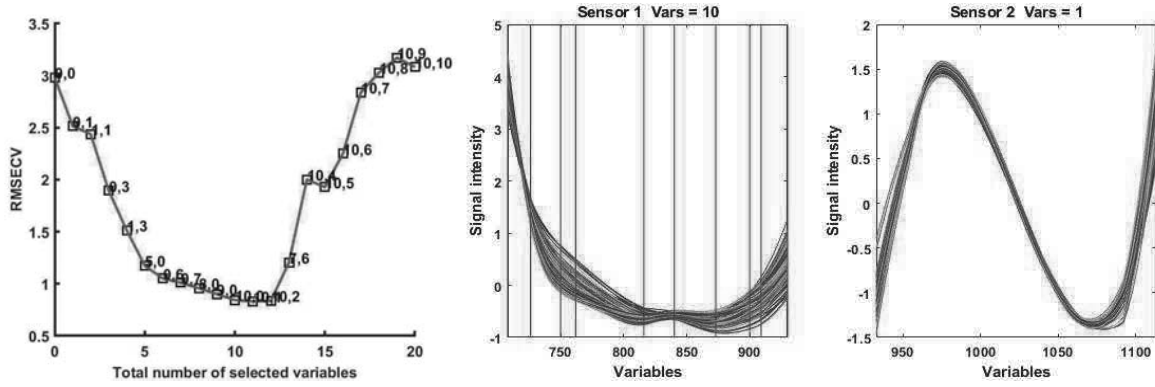
540 Variable selection

541

542 Variable selection is a key step to identify the predictive variables that are most responsible
543 for explaining the response variables. Variable selection can give a better understanding of
544 the important parameters and, in many cases, help in the development of cheap multi-spectral
545 sensor systems. In this work, a demonstration of multi-block variable selection using SO-
546 CovSel is given. The analysis was carried out on the olive fruits dry matter content dataset,
547 where two data blocks in the NIR range are used to predict the dry matter in olives.
548 Currently, only CovSel variable selection is integrated into the MBA-GUI, however, it can be
549 used for both continuous (regression) and discrete (classification) response variables. A
550 cross-validation option is also provided which supports the selection of key variables. Once
551 the response variables are loaded, the calibrate button can be used and the results will appear
552 in new figures. Fig. 9 shows the outcome of the SO-CovSel analysis where Fig. 9A shows the
553 cross-validation error, Fig. 9B shows the variables selected from block 1 and Fig. 9C shows
554 the variables selected from block 2. A total of 11 variables were selected, 10 from block 1
555 and only 1 from block 2. The number of selected variables is almost 1/10 of the initial 146 in
556 the two blocks. The calibration and prediction R^2 were 0.93 and 0.91, respectively (Fig. 9D
557 and 9E), and the RMSEC and RMSEP values were 0.78 and 0.77 % dried matter. Although
558 there was a slight decrease in R^2 and a slight increase in RMSEP with the models based on
559 selected variables compared to SO-PLS regression, it should be noted that the model is now
560 much simpler as it includes only 11 variables.

561

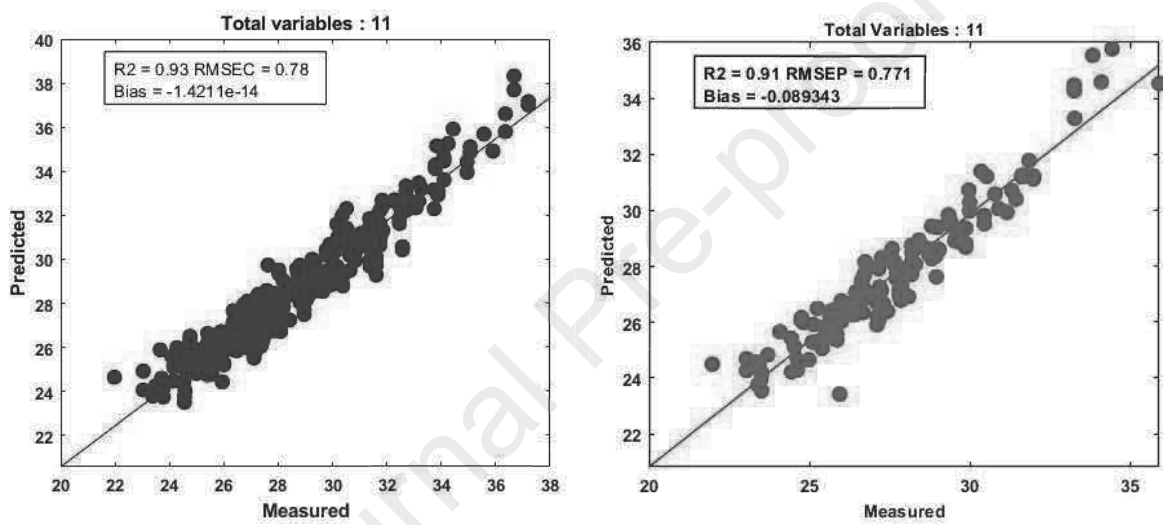
562



(A).

(B).

(C).



(D).

(E).

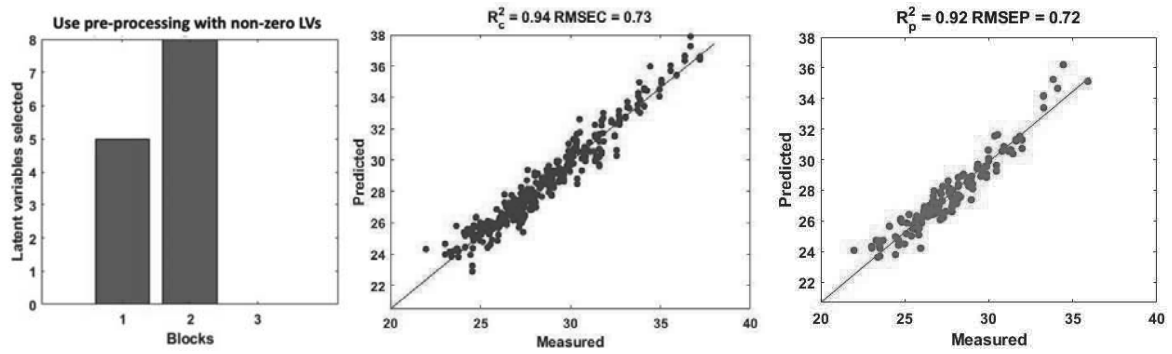
Fig. 9: Output results for SO-CovSel analysis performed on the olive fruits dry matter content dataset. (A) Error plot for variable extraction, (B) variables selected in block 1, (C) variables selected in block 2, (D) calibration set, and (E) test set.

573 SPORT

574

575 Choosing the best pre-processing technique can sometimes be a challenging task such as in
576 the case of spectral data. However, multiple pre-processing techniques can provide
577 complementary information; for example, pre-processing with a derivative can help in
578 revealing the underlying peaks and techniques such as SNV can help to reduce multiplicative
579 effects. In the present MBA-GUI, a newly developed approach called SPORT is also
580 integrated. Using SPORT is completely automated and just requires the user to load the same
581 data in multiple blocks. The user can load the same data up to four times, and therefore, can
582 apply a different combination of pre-processing which the user thinks are the best candidates
583 for the type of data.

584 In the present case, three-blocks were assigned with three different types of pre-processing.
585 The first block was pre-processed using SNV, the second block has a combination of
586 SAVGOL smoothing and VSN, and the third block had a combination of SAVGOL and
587 MSC. After performing SPORT, the best pre-processing options were selected highlighting
588 the LVs extracted from each block (Fig. 10A). The block having zero LVs is not useful, and
589 therefore, the pre-processing associated with this block provides no improvement over the
590 pre-processing previously used. In the present case, SNV and VSN pre-processing both had
591 non-zeros LVs, in contrast to MSC pre-processing which had zero LVs, meaning that MSC is
592 not useful in this case. Further, both SNV (5 LVs) and VSN (8 LVs) pre-processing has a
593 significant number of LVs in the final calibration. The model obtained had a R^2 of calibration
594 and prediction of 0.94 and 0.92, respectively. Further, the RMSEC and RMSEP were 0.73 %
595 and 0.72 %, respectively. In summary, SPORT identified the best pre-processing method.



596

597

(A)

(B).

(C).

598

Fig.10: Selection of the best pre-processing options and their fusion with the SPORT

599

methodology. (A) Latent variables selected from each block, (B) calibration set modelled with

600

selected pre-processing fusion, and (C) test set.

601

602

603

604 Conclusion

605 A MATLAB based GUI for multi-block data analysis (MBA-GUI) is presented. The toolbox
606 can perform a range of common pre-processing methods on blocks of multivariate data.
607 Multi-block data analysis for regression, classification, visualisation, variable selection and
608 SPORT are proposed. The performance of the MBA-GUI for each of the data analysis tasks
609 was demonstrated with several data sets. The results showed that the MBA-GUI performed
610 well, and all the options are fully functional. The main advantage of the toolbox is that it can
611 be easily understood and used by non-experts. The first version of the GUI can be
612 downloaded at (<https://github.com/puneetmishra2/Multi-block.git>). Other features will be
613 added to the GUI with the development of new methods. All the data analysis presented in
614 this work can be replicated with the supplied data. The GUI supports data format of .csv,
615 .xlsx and .mat. The users are welcome to notify the authors if they find any bug or problem
616 related to the use of the toolbox, so that the toolbox can be continuously improved along
617 time. The app can be directly installed in MATLAB or can be used as stand by installing the
618 MATLAB 2018b run time compiler tool at
619 (<https://nl.mathworks.com/products/compiler/matlab-runtime.html>). The password to the start
620 the toolbox is “welovedata” without double colon.

621 Validation

622

623 Dr. Raffaele Vitale

- 624 1) U. Lille, CNRS, LASIRE, Laboratoire de Spectroscopie pour les Interactions, la Réactivité
625 et l'Environnement, Cité Scientifique, F-59000, Lille, France
- 626 2) Molecular Imaging and Photonics Unit, Department of Chemistry, KU Leuven,
627 Celestijnenlaan 200F, B-3001, Leuven, Belgium

628 Email: rvitale86@gmail.com

629

630

631 The MBA-GUI Toolbox v. 1.0 was successfully executed on three different versions of
632 MATLAB (R2016b, R2017b and R2019a) without any issue to report. The Graphical User
633 Interface (GUI) for data loading and visualisation allows a maximum number of 4 distinct
634 numerical arrays to be imported, represented and pretreated by means of a significant amount
635 of computational tools for smoothing, derivation and other types of corrections (*e.g.*,
636 Standard Normal Variate – SNV – Multiplicative Scatter Correction – MSC – *etc.*). It is
637 intuitive and very well-conceived.

638 The imported datasets can afterwards be processed by a remarkable collection of multi-block
639 latent variable-based dimensionality reduction methodologies permitting tasks of various
640 nature to be carried out:

- 641 1. Variable selection (Sequential Orthogonalised Covariance Selection and Sequential
642 Orthogonalised Covariance Selection-Linear Discriminant Analysis);
- 643 2. data exploration/visualisation (Common Dimension – ComDim – Principal
644 Component Analysis, ComDim Canonical Correlation Analysis, ComDim
645 Independent Component Analysis and ComDim Partial Least Squares regression);
- 646 3. multivariate regression (Sequential Orthogonalised Partial Least Squares regression
647 and ComDim Regression); and
- 648 4. multivariate classification (Sequential Orthogonalised Covariance Selection-Linear
649 Discriminant Analysis, ComDim Independent Component Analysis-Linear
650 Discriminant Analysis, ComDim Principal Component Analysis-Linear Discriminant
651 Analysis and Sequential Orthogonalised Partial Least Squares Discriminant Analysis).

652

653 All the toolbox menus are extremely easy to browse and their design is capable of guiding
654 even non-expert users through the sequential steps of multi-block data analysis. They enable
655 not only model training/calibration, but also (when it holds) model optimisation (by two
656 different approaches for cross-validation) and external testing/validation. Furthermore, it is
657 worth mentioning that the presented GUI offers as a valuable add-on an implementation of a
658 recently proposed strategy named SPORT (Sequential Preprocessing through
659 ORThogonalisation) for both the selection of the best pretreatment technique and the
660 extraction of complementary information from the outcomes of multiple preprocessing
661 operations.

662 In its ensemble, the developed toolbox constitutes a comprehensive software suite to address
663 multi-block data analysis problems, which shows a great potential for attracting practitioners
664 by *making their life easier* in scenarios that might exhibit particularly high complexities.

665

666 Disclaimer

667

668 The MBA-GUI is free to use for public as it also involves some algorithms from public
669 sources. Great care has been taken while developing the MBA-GUI, however, the authors do
670 not accept any responsibility or liability.

671 Acknowledgments

672 PM and AN acknowledge the funding received from the European Union's Horizon 2020
673 research and innovation program named MODLIFE (Advancing Modelling for Process-
674 Product Innovation, Optimization, Monitoring and Control in Life Science Industries) under
675 the Marie Skłodowska-Curie grant agreement number 675251.

676

677 References

- 678 [1] L.L. Simon, H. Pataki, G. Marosi, F. Meemken, K. Hungerbühler, A. Baiker, S. Tummala, B.
679 Glennon, M. Kuentz, G. Steele, H.J.M. Kramer, J.W. Ryzak, Z. Chen, J. Morris, F. Kjell, R. Singh, R.
680 Gani, K.V. Gernaey, M. Louhi-Kultanen, J. O'Reilly, N. Sandler, O. Antikainen, J. Yliruusi, P. Froberg,
681 J. Ulrich, R.D. Braatz, T. Leyssens, M. von Stosch, R. Oliveira, R.B.H. Tan, H. Wu, M. Khan, D. O'Grady,
682 A. Pandey, R. Westra, E. Delle-Case, D. Pape, D. Angelosante, Y. Maret, O. Steiger, M. Lenner, K.
683 Abbou-Oucherif, Z.K. Nagy, J.D. Litster, V.K. Kamaraju, M.-S. Chiu, Assessment of Recent Process
684 Analytical Technology (PAT) Trends: A Multiauthor Review, *Organic Process Research &*
685 *Development*, 19 (2015) 3-62.
- 686 [2] E. Borràs, J. Ferré, R. Boqué, M. Mestres, L. Aceña, O. Busto, Data fusion methodologies for food
687 and beverage authentication and quality assessment – A review, *Analytica Chimica Acta*, 891 (2015)
688 1-14.
- 689 [3] H. Zheng, A. Cai, Q. Zhou, P. Xu, L. Zhao, C. Li, B. Dong, H. Gao, Optimal preprocessing of serum
690 and urine metabolomic data fusion for staging prostate cancer through design of experiment, *Anal*
691 *Chim Acta*, 991 (2017) 68-75.
- 692 [4] T.G. Doeswijk, A.K. Smilde, J.A. Hageman, J.A. Westerhuis, F.A. van Eeuwijk, On the increase of
693 predictive performance with high-level data fusion, *Anal Chim Acta*, 705 (2011) 41-47.
- 694 [5] Q. Ouyang, J. Zhao, Q. Chen, Instrumental intelligent test of food sensory quality as mimic of
695 human panel test combining multiple cross-perception sensors and data fusion, *Anal Chim Acta*, 841
696 (2014) 68-76.
- 697 [6] A. Biancolillo, R. Bucci, A.L. Magri, A.D. Magri, F. Marini, Data-fusion for multiplatform
698 characterization of an italian craft beer aimed at its authentication, *Anal Chim Acta*, 820 (2014) 23-
699 31.
- 700 [7] A.R. Martínez Bilesio, M. Batistelli, A.G. García-Reiriz, Fusing data of different orders for
701 environmental monitoring, *Anal Chim Acta*, 1085 (2019) 48-60.
- 702 [8] T. Meng, X. Jing, Z. Yan, W. Pedrycz, A survey on machine learning for data fusion, *Information*
703 *Fusion*, 57 (2020) 115-129.
- 704 [9] A.K. Smilde, I. Måge, T. Næs, T. Hankemeier, M.A. Lips, H.A.L. Kiers, E. Acar, R. Bro, Common and
705 distinct components in data fusion, *Journal of Chemometrics*, 31 (2017) e2900.
- 706 [10] M. Alinaghi, H.C. Bertram, A. Brunse, A.K. Smilde, J.A. Westerhuis, Common and distinct
707 variation in data fusion of designed experimental data, *Metabolomics*, 16 (2019) 2.
- 708 [11] I. Måge, A.K. Smilde, F.M. van der Kloet, Performance of methods that separate common and
709 distinct variation in multiple data blocks, *Journal of Chemometrics*, 33 (2019) e3085.
- 710 [12] Y. Song, J.A. Westerhuis, A.K. Smilde, Separating common (global and local) and distinct
711 variation in multiple mixed types data sets, *Journal of Chemometrics*, 34 (2020) e3197.
- 712 [13] J. Boccard, D.N. Rutledge, A consensus orthogonal partial least squares discriminant analysis
713 (OPLS-DA) strategy for multi-block Omics data fusion, *Analytica Chimica Acta*, 769 (2013) 30-39.
- 714 [14] J. Boccard, D.N. Rutledge, Iterative weighting of multi-block data in the orthogonal partial least
715 squares framework, *Anal Chim Acta*, 813 (2014) 25-34.
- 716 [15] A. Biancolillo, T. Næs. The Sequential and Orthogonalized PLS Regression for Multi-block
717 Regression: Theory, Examples, and Extensions. In: M. Cocchi (Ed.), *Data Fusion Methodologies and*
718 *Applications, Data Handling in Science and Technology*, vol. 31, Elsevier, Oxford, UK, 2019, pp. 157-
719 177.
- 720 [16] A. Biancolillo, F. Marini, J.-M. Roger, SO-CovSel: A novel method for variable selection in a multi-
721 block framework, *Journal of Chemometrics*, 34 (2020) e3120.
- 722 [17] B. Galindo-Prieto, P. Geladi, J. Trygg, Multi-block variable influence on orthogonal projections
723 (MB-VIOP) for enhanced interpretation of total, global, local and unique variations in OnPLS models,
724 arXiv preprint arXiv:2001.06530, (2020).

- 725 [18] E. Acar, M.A. Rasmussen, F. Savorani, T. Næs, R. Bro, Understanding data fusion within the
726 framework of coupled matrix and tensor factorizations, *Chemometr Intell Lab*, 129 (2013) 53-63.
- 727 [19] E. Acar, E.E. Papalexakis, G. Gürdeniz, M.A. Rasmussen, A.J. Lawaetz, M. Nilsson, R. Bro,
728 Structure-revealing data fusion, *BMC Bioinformatics*, 15 (2014) 239.
- 729 [20] T. Næs, O. Tomic, N.K. Afseth, V. Segtnan, I. Måge, Multi-block regression based on
730 combinations of orthogonalisation, PLS-regression and canonical correlation analysis, *Chemometrics
731 and Intelligent Laboratory Systems*, 124 (2013) 32-42.
- 732 [21] K. De Roover, E. Ceulemans, M.E. Timmerman, How to perform multi-block component analysis
733 in practice, *Behavior Research Methods*, 44 (2012) 41-56.
- 734 [22] V. Cariou, D. Jouan-Rimbaud Bouveresse, E.M. Qannari, D.N. Rutledge, ComDim Methods for
735 the Analysis of Multi-block Data in a Data Fusion Perspective. In: M. Cocchi (Ed.), *Data Fusion
736 Methodologies and Applications, Data Handling in Science and Technology*, vol.31, Elsevier, Oxford,
737 UK, 2019, pp. 179-204.
- 738 [23] J.-M. Roger, J.-C. Boulet, M. Zeaiter, D.N. Rutledge, . In: S.D. Brown, R. Tauler, B. Walczak (Eds.),
739 *Comprehensive Chemometrics*, 2nd Ed., vol. 3, Elsevier, Amsterdam, The Netherlands, 2020, pp. 1-
740 75.
- 741 [24] J. Engel, J. Gerretzen, E. Szymańska, J.J. Jansen, G. Downey, L. Blanchet, L.M.C. Buydens,
742 Breaking with trends in pre-processing?, *TrAC Trends in Analytical Chemistry*, 50 (2013) 96-106.
- 743 [25] Å. Rinnan, F.v.d. Berg, S.B. Engelsen, Review of the most common pre-processing techniques for
744 near-infrared spectra, *TrAC Trends in Analytical Chemistry*, 28 (2009) 1201-1222.
- 745 [26] R.J. Barnes, M.S. Dhanoa, S.J. Lister, Standard Normal Variate Transformation and De-Trending
746 of Near-Infrared Diffuse Reflectance Spectra, *Applied Spectroscopy*, 43 (1989) 772-777.
- 747 [27] T. Isaksson, T. Næs, The Effect of Multiplicative Scatter Correction (MSC) and Linearity
748 Improvement in NIR Spectroscopy, *Applied Spectroscopy*, 42 (1988) 1273-1284.
- 749 [28] P.H.C. Eilers, Parametric time warping, *Analytical Chemistry*, 76 (2004) 404-411.
- 750 [29] G. Rabatel, F. Marini, B. Walczak, J.-M. Roger, VSN: Variable sorting for normalization, *Journal of
751 Chemometrics*, 34 (2020) e3164.
- 752 [30] F. Dieterle, A. Ross, G. Schlotterbeck, H. Senn, Probabilistic Quotient Normalization as Robust
753 Method to Account for Dilution of Complex Biological Mixtures. Application in 1H NMR
754 Metabonomics, *Analytical Chemistry*, 78 (2006) 4281-4290.
- 755 [31] Q. Guo, W. Wu, D.L. Massart, The robust normal variate transform for pattern recognition with
756 near-infrared data, *Anal Chim Acta*, 382 (1999) 87-103.
- 757 [32] R. Bro, A.K. Smilde, Principal component analysis, *Anal Methods-Uk*, 6 (2014) 2812-2831.
- 758 [33] P. Geladi, Chemometrics in spectroscopy. Part 1. Classical chemometrics, *Spectrochimica Acta
759 Part B: Atomic Spectroscopy*, 58 (2003) 767-782.
- 760 [34] E.M. Qannari, I. Wakeling, P. Courcoux, H.J.H. MacFie, Defining the underlying sensory
761 dimensions, *Food Quality and Preference*, 11 (2000) 151-154.
- 762 [35] A. El Ghaziri, V. Cariou, D.N. Rutledge, E.M. Qannari, Analysis of multi-block datasets using
763 ComDim: Overview and extension to the analysis of (K + 1) datasets, *Journal of Chemometrics*, 30
764 (2016) 420-429.
- 765 [36] D. Rutledge, Novel extensions and applications of Common Components Analysis in
766 chemometrics, *Twelfth Winter Symposium on ChemometricsSaratov (Russia)*, 2020.
- 767 [37] A. Biancolillo, T. Næs, R. Bro, I. Måge, Extension of SO-PLS to multi-way arrays: SO-N-PLS,
768 *Chemometr Intell Lab*, 164 (2017) 113-126.
- 769 [38] T. Næs, O. Tomic, B.H. Mevik, H. Martens, Path modelling by sequential PLS regression, *Journal
770 of Chemometrics*, 25 (2011) 28-40.
- 771 [39] A. Biancolillo, I. Måge, T. Næs, Combining SO-PLS and linear discriminant analysis for multi-block
772 classification, *Chemometr Intell Lab*, 141 (2015) 58-67.
- 773 [40] J.M. Roger, B. Palagos, D. Bertrand, E. Fernandez-Ahumada, CovSel: Variable selection for highly
774 multivariate and multi-response calibration Application to IR spectroscopy, *Chemometr Intell Lab*,
775 106 (2011) 216-223.

- 776 [41] J.-M. Roger, A. Biancolillo, F. Marini, Sequential preprocessing through ORThogonalization
777 (SPORT) and its application to near infrared spectroscopy, *Chemometrics and Intelligent Laboratory*
778 *Systems*, 199 (2020) 103975.
- 779 [42] W.B. Zheng, H.P. Shu, H. Tang, H.Q. Zhang, Spectra data classification with kernel extreme
780 learning machine, *Chemometr Intell Lab*, 192 (2019).
- 781 [43] H.S. Tapp, M. Defernez, E.K. Kemsley, FTIR spectroscopy and multivariate analysis can
782 distinguish the geographic origin of extra virgin olive oils, *Journal of Agricultural and Food Chemistry*,
783 51 (2003) 6110-6115.
- 784 [44] X. Sun, P. Subedi, R. Walker, K.B. Walsh, NIRS prediction of dry matter content of single olive
785 fruit with consideration of variable sorting for normalisation pre-treatment, *Postharvest Biology and*
786 *Technology*, 163 (2020) 111140.

787

788

789

790

Declaration of interests

The authors declare that they have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper.

Journal Pre-proof