

Recent trends in multi-block data analysis in chemometrics for multi-source data integration

Puneet Mishra^{1,2}, Jean Michel Roger^{3,4}, Delphine Jouan-Rimbaud-Bouveresse⁵, Alessandra Biancolillo⁶, Federico Marini⁷, Alison Nordon², Douglas N. Rutledge^{8,9}

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, INRAE Montpellier, Institut Agro, University Montpellier, Montpellier, France*

⁴*ChemHouse Research Group, Montpellier, France*

⁵*UMR PNCA. AgroParisTech, INRA. Université Paris-Saclay, 75005 Paris, France*

⁶*Department of Physical and Chemical Sciences, University of L'Aquila, Via Vetoio 67100, Coppito, L'Aquila, Italy*

⁷*Department of Chemistry, University of Rome "La Sapienza", P.le Aldo Moro 5, 00185, Rome, Italy*

⁸*Université Paris-Saclay, INRAE, AgroParisTech, UMR SayFood, 75005, Paris, France*

⁹*National Wine and Grape Industry Centre, Charles Stuart University, Wagga Wagga, Australia*

Corresponding author: puneet.mishra@wur.nl

- 22 [Keywords](#)
- 23 2D: 2 dimensional
- 24 3D: 3 dimensional
- 25 CCSWA: Common component and specific weight analysis
- 26 ComDim: Common Dimensions
- 27 CT: Calibration transfer
- 28 DISCO-SCA: Distinct and common simultaneous component analysis
- 29 GCA: Generalized canonical analysis
- 30 GSVD: Generalized singular value decomposition
- 31 GUI: Graphical user interface
- 32 H-PCA: Hierarchal principal component analysis
- 33 H-PLS: Hierarchal partial least-squares
- 34 JIVE: Joint and individual variances explained
- 35 JUMBA: Joint and unique multi-block analysis
- 36 MB-PCA: Multi-block principal component analysis
- 37 MB-PLS: Multi-block partial least-squares
- 38 MB-VIOP: Multi-block variable importance in projection
- 39 MBA-GUI: Multi-block analysis graphical user interface
- 40 **MCR: Multivariate curve resolution**
- 41 MIR: Mid-infrared

- 42 MOCA: Multiple co-inertia analysis
- 43 **MVP: Multi-block variable partitioning**
- 44 NIR: Near-infrared
- 45 O2PLS: Orthogonal 2-block partial least-squares
- 46 OnPLS: Orthogonal n-block partial least-squares
- 47 P-EASCA: Penalized exponential analysis of variance simultaneous component analysis
- 48 P-ESCA: Penalized exponential simultaneous component analysis
- 49 PARAFAC: Parallel factor analysis
- 50 PAT: Process analytical technologies
- 51 PCA-GCA: Principal component analysis generalized canonical analysis
- 52 PCA: Principal component analysis
- 53 PLS: Partial least-squares
- 54 PO-PLS: Parallel orthogonalized partial-least squares regression
- 55 PORTO: Parallel pre-processing through orthogonalization
- 56 SCA: Simultaneous component analysis
- 57 SCD-PCovR: Sparse common and distinct principal covariate regression
- 58 SLIDE: structured learning and integrative decomposition
- 59 SO-CovSel: Sequential orthogonalized covariate selection
- 60 SO-N-PLS: Sequential orthogonalized n-way partial least-squares
- 61 SO-PLS: Sequential orthogonalized partial-least squares regression

62 SPORT: Sequential pre-processing through orthogonalization

63 SR: Selectivity ratio

64 VIP: Variable importance in projection

65

66 1. Abstract

67 In recent years, multi-modal measurements of process and product properties have become
68 widely popular. Sometimes classical chemometric methods such as principal component
69 analysis (PCA) and partial least squares regression (PLS) are not adequate to analyze this kind
70 of data. In recent years, several multi-block methods have emerged for this purpose; however,
71 their use is largely limited to chemometricians, and non-experts have little experience with
72 such methods. In order to deal with this, the present review provides a brief overview of the
73 multi-block data analysis concept, the various tasks that can be performed with it and the
74 advantages and disadvantages of different techniques. Moreover, basic tasks ranging from
75 multi-block data visualization to advanced innovative applications such as calibration transfer
76 will be briefly highlighted. Finally, a summary of software resources available for multi-block
77 data analysis is provided.

78 *Keywords: pre-processing fusion; incremental learning; data fusion; chemometrics*

79 2. Introduction

80 In analytical chemistry, data obtained by multiple sources is frequently encountered [1-3]. A
81 *multi-block* data set can either come from a multi-platform analysis of the same samples (e.g.,
82 to reach a better understanding of the physico-chemical properties of the analyzed objects
83 which is not possible with a single technique [1, 4]), or by the combination of chemical
84 measurements with non-analytical data generated from sensory or consumer sciences [5]. In

85 both cases, the data is not simply multivariate but is *multi-modal*, i.e., multivariate and multi-
86 source. An example of this would be data generated by two different spectroscopic techniques
87 such as mid-infrared spectroscopy (MIR) and Raman spectroscopy. In this case, spectral
88 profiles are multivariate (as the responses are acquired at several wavenumbers), and the modes
89 are represented by the two different spectroscopic techniques. Furthermore, multi-modal data
90 can also be obtained when working under different conditions, for instance, when multiple
91 batches of an industrial process produce data under different processing conditions [6-9].

92 Chemometrics has been developed to handle multivariate data generated from analytical
93 techniques [10, 11]. The foundation of chemometrics lies on the identification of the underlying
94 latent spaces using bilinear or trilinear multivariate data decomposition techniques. These
95 explorations of latent spaces are specifically targeted to find any structured variation and/or
96 correlation with the key property of interest. Once identified, latent spaces can be used to
97 perform several data processing tasks, such as transforming high-dimensional data to a lower
98 dimensional representation for data visualization purposes [12], or developing regression
99 models for predictive analysis [13] and identification of key variables of interest [14-16].
100 Traditional chemometric methods (single-block chemometric techniques), such as principal
101 component analysis (PCA) [12], partial least squares (PLS) regression [13] and their variants,
102 only work properly when the data is single-mode, i.e., generated by only one source of
103 variability, such as a single analytical technique. In the case of multi-block data, the standard
104 single-block-techniques extract only a limited part of the information present in the data [3,
105 17]. It is only by using multi-block data analysis techniques that it is possible to extract the
106 complementary information from data generated in multiple modes [3, 18].

107 Multi-block data analysis accomplishes similar tasks to single-block chemometric techniques,
108 i.e., enhancing data visualization [3, 19] improving predictive performances [1, 4], and
109 identifying the key variables that influence the models [3, 19-21]. Furthermore, multi-block

110 analysis can achieve an enhanced understanding of the common and the distinct information
111 present in the data coming from different platforms [3, 19-21]. The present review provides a
112 brief overview of the multi-block data analysis concept, the various tasks that can be performed
113 with it and the advantages and disadvantages of different techniques. Moreover, basic tasks
114 ranging from multi-block data visualization to advanced innovative applications such as
115 calibration transfer will be briefly highlighted. Special attention has been paid to simplify the
116 explanation of complex concepts and methods related to multi-block data analysis so that users
117 with minimal experience in chemometrics can understand and use advanced multi-block
118 methods in their daily tasks. And finally, a summary of software resources available for multi-
119 block data analysis is provided.

120

121 3. When is the multi-block data generated and what are its 122 characteristics?

123 An example of an experimental that produces a multi-block data set is presented in Fig. 1,
124 where four different spectroscopic techniques (near-infrared spectroscopy, mid-infrared
125 spectroscopy, Raman spectroscopy and fluorescence spectroscopy) are combined to monitor
126 the chemical process taking place in a glass vessel. The data from all the four techniques are
127 acquired simultaneously as presented in Fig. 2. This is just an illustration, but such acquisitions
128 of multi-source data are becoming popular in analytical chemistry [22-25]. A similar example
129 is the non-destructive quantification of (bio-) chemical components in an aqueous process
130 containing fluorescent compounds.

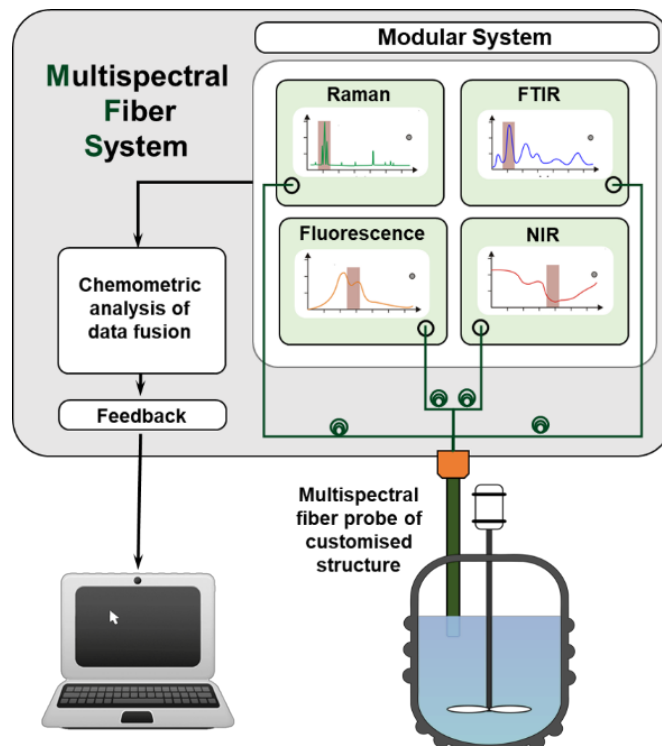
131



132

133 *Fig. 1: Scheme of multispectral fiber system (figure courtesy of Art Photonics GmbH, Germany*
 134 *[26]). Raman system (1); FTIR absorption System (2); NIR reflection System (3); fluorescence*
 135 *detector (4); chemical reactor (5); and fiber optic probes (6).*

136



137

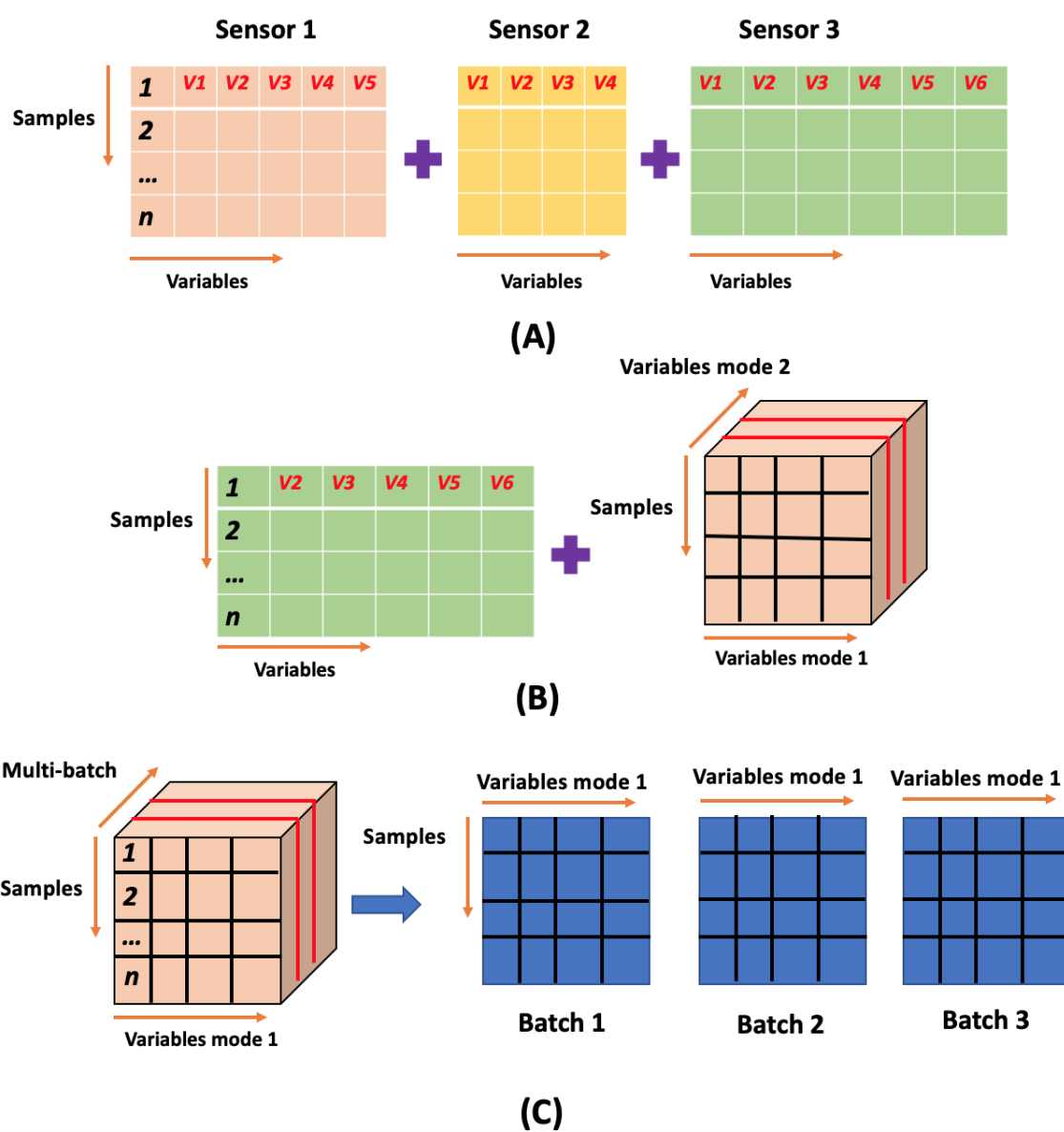
138 *Fig. 2: A schematic of the multiblock data generated in a four-blocks scenario (figure courtesy*

139 *of Art Photonics GmbH, Germany) [26].*

140

141 The main characteristics of multi-block data is that it either consists of multiple matrices
142 corresponding to different analytical platforms generated from measurement on the same
143 sample (Fig. 3A), a combination of matrices and higher order tensors (Fig. 3B), or different
144 independent batch processes (Fig. 3C).

145



146

147 *Fig. 3: A summary of typical multi-block data configurations in analytical chemistry. (A) data*

148 *generated by multiple analytical platforms in the form of 2D matrices, (B) data generated by*
149 *multiple analytical platforms in the form of 2D matrices or higher order tensors, and (C) multi-*
150 *set data from batch processes which can be treated as multi-block data when different batches*
151 *are treated independently.*

152 4. Multi-block data pre-processing

153 Data pre-processing is essential in chemometrics and, like the standard one-block methods,
154 multi-block analysis is also influenced by the pre-processing operation. In the case of multi-
155 block analysis, data pre-treatment can be divided into two stages, i.e., the *inter-block* and the
156 *intra-block* pre-processing [8]. Recently, Campos and Reis provided a comprehensive
157 systemization of multi-block data pre-processing and divided the steps into three levels [8]. In
158 particular, the first level [8] includes the standard chemometric pre-processing operations to
159 correct artefacts and uninteresting variations such as noise, multiplicative effects, scaling,
160 baseline drift, peak shift and variations related to external factors [27]. At the second level, the
161 aim is to equalize the contribution of all variables within each block and this can be achieved
162 by classical methods such as mean-centering and unit variance scaling [8]. The third (and final)
163 level aims at equalizing the inter-block systemic effects such as the differences in the scales,
164 number of variables and the pseudo rank of different blocks [8]. This level of pre-processing
165 is necessary as some multi-block analysis methods tend to favor the blocks with larger
166 variations, leading to model bias. However, with proper scaling or weighting of blocks, it has
167 been proven that model interpretation and predictive accuracy can be increased [28]. The third
168 level (inter-block) pre-processing approaches are mainly scaling, such as block scaling, block
169 variance scaling and block rank scaling [8, 28]. Block scaling and block variance scaling
170 balance the effect of the modelled blocks, to avoid any block dominating the model [8]. More
171 detailed information on multi-block pre-processing can be accessed in a recently published

172 work [8]. Although multi-block pre-processing is important, not all approaches to multi-block
173 data analysis require all levels of pre-processing. For example, the sequential and parallel
174 approaches to partial least-squares regression are less sensitive to the relative scaling of the
175 blocks and can also deal with the differences in the ranks of multi-block data [8], because these
176 methods handle the multi-block data one block at a time, involving orthogonalization steps
177 which do not affect the relative weighting of the blocks [8, 18].

178 In conclusion, the main outcome of all these considerations is that multi-block pre-processing
179 must be carefully planned in accordance with the multi-block analysis to be performed.

180 5. Multi-block exploratory data analysis

181 Exploratory analysis is often performed to obtain low-dimensional representations of high-
182 dimensional multivariate data, to facilitate its interpretation. The key properties of the data can
183 thus be visualized by means of interpretable 2D or 3D plots. In chemometrics, one of the aims
184 is to identify the latent (sub-)spaces capturing key properties of data such as highest
185 variance/closest fit in the case of PCA, or maximum co-variance to the response variables for
186 supervised data decomposition such as PLS. The data decomposition results in a set of scores
187 and loadings, where the scores are the low dimensional representation of the data and the
188 loadings are the latent vectors spanning the relevant sub-space. In standard one block
189 chemometrics, different methods are available for latent space modelling. In fact, the
190 identification of these latent spaces depends on data modes: for a simple 2D data matrix,
191 bilinear decomposition methods such as PCA can be implemented, whereas when the order of
192 the data increases to 3D or more, then higher-order extensions of PCA called Tucker and
193 parallel factor analysis (PARAFAC) can be implemented. However, one block chemometrics
194 methods do not provide a complete solution to deal with multi-block data.

195 To deal with the challenges of visualizing multi-block multivariate data, several extensions of
196 one-block chemometric methods as well as new specific multi-block approaches have emerged
197 in recent decades. A summary of these methods is provided in Table 1. There are different
198 classifications of the multi-block methods; one is based on the separation into two families,
199 depending on how these approaches handle common and distinct information in the blocks.
200 The first family comprises approaches based on identifying the common information among
201 different data blocks and later exploring the contribution of each block to the common
202 components. The second family of methods is based on the simultaneous extraction of the
203 common as well as the distinct information in the different data blocks. The methods belonging
204 to the first family are extensions of PCA. A simple, popular extension of PCA for the multi-
205 block scenario is SUM PCA [17], where multiple data blocks are concatenated in the variable's
206 domain and standard PCA is performed on the joint data, leading to extraction of global
207 principal components. More advanced extensions, called multi-block PCA (MB-PCA) or
208 consensus PCA [29], allow extraction of global components as well as the contribution of the
209 associated blocks. The extraction of subsequent components is performed by deflation of the
210 matrices with respect to the global components. This is done by regressing all the variables in
211 the different blocks with the extracted global component, and the resulting residuals are then
212 used to extract new global components, and so on. The deflation step is performed such that
213 each global component contains unique information. A method like MB-PCA, called multiple
214 co-inertia analysis (MOCA), was also proposed to explore multi-block data [30]. From an
215 algorithmic point of view, MOCA follows the same procedure as MB-PCA to extract the global
216 components, however, in the second step, the block loadings are used for block deflation, not
217 the global scores as used in the case of MB-PCA [31]. The deflation with the block loadings in
218 MOCA allows the blocks to capture orthogonal information. Furthermore, an advanced version
219 of PCA for multi-block analysis is the hierarchal PCA (H-PCA) which provides the global

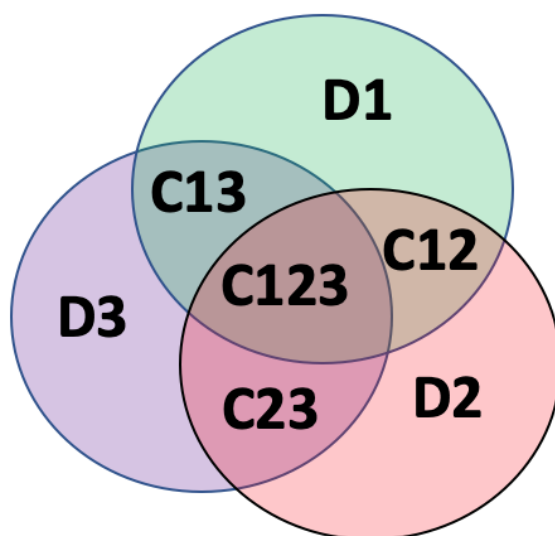
220 components along with the weights of the blocks to reflect the importance of each block in
221 contributing to the extracted global components [29, 32]. Like H-PCA, a method called
222 common components and specific weight analysis (CCSWA or ComDim) has gained attention,
223 as it also allows the extraction of global components and the specific weights for each block to
224 have an insight into each one's importance [33-35]. Ideally, both H-PCA and CCSWA lead to
225 similar solutions but CCSWA is more sophisticated in terms of mathematical explanations with
226 several possibilities of method extensions [35-38].

227 The methods extracting common components (global components) work well when the
228 objective is to globally explore the different blocks, however, they lack the means to highlight
229 which information is unique in each block. A framework for extraction of common and distinct
230 information recently summarized by Smilde et al. [3], is presented in Fig. 4, where the three
231 circles represent three different blocks of data measured on the same samples and the letters D
232 and C indicate the distinct and the common parts of the information, respectively. Several
233 methods can be identified in the framework of common and distinct information extraction,
234 and the evaluation of their performances can be found elsewhere [3, 7, 20, 21]. Some examples
235 of these are distinct and common simultaneous component analysis (DISCO-SCA) [39],
236 principal component analysis - generalized canonical analysis (PCA-GCA) [3], generalized
237 singular value decomposition (GSVD) [40], orthogonal2 PLS (O2PLS) [41], joint and
238 individual variances explained (JIVE) [42], structure revealing data fusion [43, 44] and
239 structured learning and integrative decomposition (SLIDE) [45]. In the case of DISCO-SCA,
240 the first step involves a SCA step to decompose the matrices to a set of scores and loadings
241 matrices. In the second step, the loadings matrices are partitioned and orthogonally rotated to
242 reveal the common and distinct components in the multi-block scenario. In the case of the PCA-
243 GCA, the first step is to perform PCA followed by a GCA. PCA is performed to enhance the
244 stability of the GCA components. The second step is the regression of each block onto its own

245 common components, the residuals of the regression are the distinct subspace which can then
246 be used for subsequent PCA [3]. In the case of GSVD, after performing a preliminary SCA
247 step to filter out the noise from the data, singular value decomposition is jointly applied to the
248 different matrices, under the constraints that the left singular vectors (i.e., normalized scores)
249 be the same for all blocks and that the matrices of singular values \mathbf{D}_b obey $\sum_{b=1}^{N_{blocks}} \mathbf{D}_b^2 = \mathbf{I}$,
250 N_{blocks} being the number of blocks and \mathbf{I} the identity matrix of appropriate dimensionality.
251 Then, identification of a component as common or distinct is made based on the associated
252 singular values for the different blocks [40]. The O2PLS approach can be considered as a multi-
253 block extension of orthogonal PLS (OPLS), with the relevant difference that no asymmetric
254 relation among the blocks is implied, so that the method can be used also for exploratory
255 purposes [41]. At first, the distinct components are extracted from each block, which is the
256 deflated; accordingly, then, a PLS step is carried out to extract the common components from
257 the deflated blocks. In the case of JIVE, the different blocks of data are directly decomposed
258 into a set of common and distinct information by an iterative procedure involving alternating
259 steps of SVD decomposition of the concatenated blocks for the estimation of the common
260 component and SVD decomposition of the individual blocks after deflation of the estimated
261 common components to account for the distinct variation [42]. In the case of structure revealing
262 data fusion, the matrices are jointly factorized and with the help of penalty terms, the common
263 and distinct information is extracted [43,44]. Finally, the SLIDE can be considered as an
264 intermediate model between SUM-PCA and JIVE as it allows components to be partially
265 shared (i.e., common only to some blocks). This is achieved by arranging the loadings in a
266 block-dependent structure and imposing structure sparsity to reveal the common, distinct and
267 the partially shared information [45]. In analytical chemistry, experiments are often organized
268 by means of experimental designs (DoE), and much insight about the samples can be gained in
269 this way. Recently, to deal with this, a new multi-block data visualization tool for exploration

270 of multi-block data generated by designed experiments was proposed. The method is called
271 penalized exponential analysis of variance - simultaneous component analysis (PE-ASCA).
272 PE-ASCA is a combination of penalized exponential - simultaneous component analysis (PE-
273 SCA) [21] and the analysis of variance - simultaneous component analysis (ASCA) [46]. In
274 PE-ASCA, the multi-block data is first partitioned into common and distinct information and
275 later ASCA models are used to incorporate the design information while exploring the data
276 using the SCA. The application of PE-ASCA was recently presented in the domain of
277 metabolomics [19].

278 Another interesting family of methods is the extension of the multivariate curve resolution
279 (MCR) bilinear decomposition technique [47] to the multi-set configuration. MCR operates a
280 self-modeling resolution of mixed profiles into the contribution of the corresponding pure
281 constituents, through a bilinear modeling usually incorporating chemically inspired constraints
282 (e.g., non-negativity, unimodality, mass-balance, selectivity, just to cite a few). The basic trick
283 behind the use of the MCR for dealing with multiple data sources is to first concatenate the
284 data matrices along the common direction and then analyze this augmented data array through
285 MCR, retaining, as a sort of additional constraint, the information related to the presence of
286 different data blocks. In this respect, the method is rather flexible, as it can easily deal with
287 cases where the common direction is represented by the variables (e.g., in multi-batch
288 situations), by the samples (multi-source data integration) or, even by both (e.g., with different
289 sets of samples having all been analyzed by more than one technique). More details on the
290 extension of the MCR for multi-block data analysis can be found elsewhere [47]. Here it should
291 also be stressed that by introducing suitable so-called selectivity constraints, it is possible to
292 guide the model towards the extraction of both common and distinct components. As well,
293 through the use of other constraints (e.g., correlation), MCR can also be employed for
294 predictive purposes.



295

296 *Fig. 4: A framework of common and distinct information extraction from multi-block data.*

297 *Each circle represents data from a different technique. Inside each circle, D is the distinct*

298 *information and C is the common information. The figure is inspired by the multi-block data*

299 *concept presented in [3]).*

300

301 6. Multi-block predictive modelling

302 In recent years, a lot of effort in analytical chemistry has been put into developing spectroscopic

303 methodologies to replace certain complex and highly sophisticated wet chemistry routines for

304 quantitative analysis. In chemometrics, a common method to perform this task is partial least-

305 squares (PLS) regression [13] which decomposes the data into a set of scores and loadings.

306 Later, the scores are used to perform the ordinary least-squares regression. In the case of PLS

307 regression, the scores are extracted to have maximum covariance with the response variable(s).

308 However, PLS cannot be explicitly implemented in the scenario of multi-block data, especially

309 when the aim is to extract common and distinct information. Several approaches to do multi-

310 block predictive analysis have recently gained attention and a summary of methods can be

311 found in Table 1. Unlike standard PLS regression and classification modelling, the aim of
312 multi-block predictive methods is to extract the complementary information from multiple
313 sources to improve the quality of the models (prediction accuracy and/or interpretability). In
314 chemometrics, most of the methods for multi-block predictive analysis are extensions of
315 standard PLS regression to the multi-block scenario. In this regard, one of the first methods
316 developed was multi-block partial least-squares (MB-PLS) regression. This approach, in the
317 formulation proposed by Qin et al., 2001 allows the extraction of global scores by maximizing
318 the covariance with the response variable(s) [29, 48, 49]. The extracted global scores are used
319 in ordinary least squares regression to obtain predictive models. Hierarchical PLS (H-PLS) is a
320 more sophisticated method that allows the extraction of global and block components, giving
321 the possibility of understanding the relative contribution of each block to the global model [29,
322 32, 50]. A similar method, called P-ComDim, or ComDim (k+1), extracts global scores that
323 capture the maximum covariance with the response variable(s) [37]. This is done by
324 maximizing the covariances between the local scores of each block and the scores of the
325 response block. In a procedure like that of CCSWA, the loadings of the variables in each block
326 and the weight (*salience*) of each block can be calculated for each common component.
327 However, the MB-PLS and the ComDim (k+1) approaches do not provide a clear extraction of
328 the common and distinct information from the different data blocks. To deal with this,
329 orthogonal n-PLS (OnPLS) was proposed [51], which is the extension of the two-block O2PLS
330 to the multi-block scenario. As O2PLS, OnPLS does not introduce a priori any asymmetry
331 among the blocks, so that it could in principle be used as an exploratory technique; however,
332 if one block contains the response(s) to be predicted, by suitably combining the scores extracted
333 from all the other matrices, a global regression model can be calculated.

334 More recently, two other methods, i.e., sequential and orthogonalized-PLS (SO-PLS) and
335 parallel and orthogonalized-PLS (PO-PLS), were also proposed as extensions of standard PLS

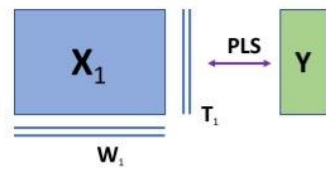
336 [52]. The SO-PLS approach involves a series of standard PLS regression and matrix
337 orthogonalization operations to extract sequentially the complementary information from
338 different data blocks; a generic schema of the algorithm is presented in Fig. 5. As mentioned,
339 in SO-PLS, the extraction of information is sequential, meaning that the aim is to incorporate
340 blocks of data one at a time and to assess their incremental contribution. A PLS regression
341 model is calculated between the first block \mathbf{X}_1 and \mathbf{Y} , yielding scores \mathbf{T}_1 . Then, all the
342 remaining blocks $\mathbf{X}_2, \dots, \mathbf{X}_k$ and \mathbf{Y} are orthogonalized with regards to \mathbf{T}_1 . The process is repeated
343 on the second block, and so on for all the blocks, taking care to orthogonalize all the following
344 blocks with respect to the previously modelled matrices. The major advantages of SO-PLS are
345 linked to the orthogonalization, which removes redundant information, and to its sequential
346 nature, which allows the interpretation of the incremental contributions provided by each data
347 block. The SO-PLS approach is particularly advantageous when the aim is to identify possible
348 extra benefits from the inclusion of each block of information into the model [18]. On the other
349 hand, the PO-PLS approach involves a combination of PLS regression, generalized canonical
350 correlation analysis (GCA) and multiple orthogonalization steps [5]. PO-PLS, unlike SO-PLS,
351 does not explore the blocks sequentially, but aims at identifying the common and the distinct
352 information in different blocks to have a better understanding of how the combinations of
353 blocks contribute to the improved predictive performances.

354 Multi-block variance partitioning (MVP), originally proposed by Skov et al. in 2008 [53],
355 presents some similarities with both SO-PLS and PO-PLS. It was one of the first methods to
356 specifically focus attention on identifying the unique part and the common part of the \mathbf{Y} -related
357 variation in the predictor blocks; this is accomplished by using PLS models between predictor
358 blocks and a common response. For each predictor block \mathbf{X}_k , the total variance of the responses
359 \mathbf{Y} is partitioned into a unique part (that ascribable only to that particular predictor block), a
360 common part (the one shared also with the other independent matrices) and an uninformative

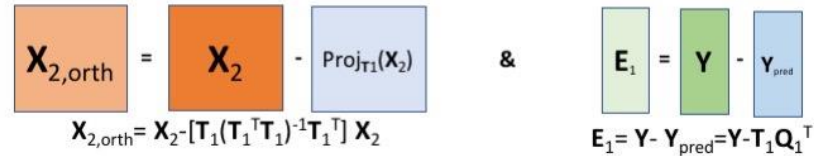
361 part (which is the portion of Y-variance not relevant for by that particular regression model).
362 Operationally, individual PLS models are calculated between each \mathbf{X}_k and \mathbf{Y} . For each predictor
363 block, the uninformative variance is associated to the residuals of that regression, while the
364 unique contribution is calculated after orthogonalizing the predicted responses (variable-wise)
365 with respect to the corresponding predicted responses based on all the other predictor blocks.
366 The common variation is obtained by subtracting the contribution of the unique and
367 uninformative parts from the total variance.

368

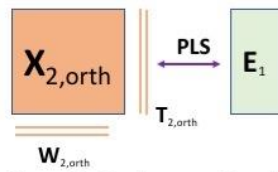
Step 1: First PLS model



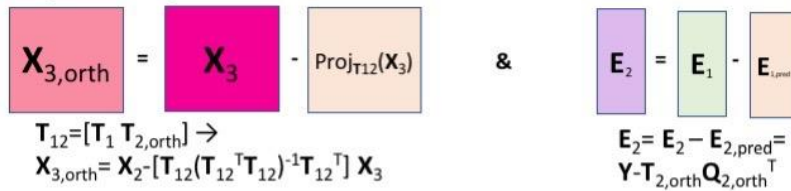
Step 2: Orthogonalization of second block



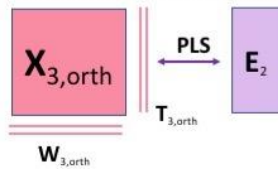
Step 3: Second PLS model



Step 4: Orthogonalization of third block



Step 5: Third PLS model



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

369

370 *Fig. 5: A schema presenting the sequential orthogonalized partial least squares (SO-PLS)*
 371 *regression method [26].*

372

373 **7. Multi-block analysis for variable selection**

374 Data generated from several analytical techniques, such as optical spectroscopy, mass
 375 spectrometry, nuclear magnetic resonance, and chromatography, consists of many variables, of
 376 which only a subset is informative. In fact, most of the variables are often correlated or related
 377 to some background phenomenon which is not of interest to explain the response variable(s).

378 Therefore, in chemometrics, variable selection is often performed to identify the most useful
379 variables for the characterization of the response [14, 15]. To deal with the multi-block
380 scenario, several predictive data fusion methods have also been extended to incorporate
381 variable selection [54, 55]. A summary of the multi-block variable selection methods is
382 presented in Table 1. The main benefit of these approaches is that they allow the selection of
383 complementary variables from multiple sources which jointly achieve good model
384 performance. Variable selection approaches can be divided into three main categories, i.e.,
385 filter, wrapper and embedded methods. Techniques belonging to the first family are based on
386 ranking the predictors according to some model-based criterion, e.g., variable importance in
387 projection (VIP) or selectivity ratio (SR), and retaining only those variables for which specific
388 parameters exceed given thresholds. In a multi-block scenario, some PLS-inspired variable
389 selection methods have been discussed; for instance, in the context of SO-PLS [54] or OnPLS
390 [56, 57]. Second, wrapper approaches directly calculate multi-block models with different
391 combinations of subsets of variables and select the one that gives the best results (usually
392 determined by internal validation, e.g., cross-validation). Finally, embedded methods carry out
393 the variable selection while building the model. In this context, an interesting possibility is the
394 recently proposed multi-block method, called sequential and orthogonalized covariance
395 selection (SO-CovSel) [55]. As the name implies, this approach is strongly related to SO-PLS,
396 and shares some of its advantages. Nevertheless, SO-CovSel is especially suited to variable
397 selection and the interpretation of the system under study because it directly provides
398 information about which variables drive the model most. However, a key thing that SO-CovSel
399 [55] and other multiblock variables selection methods [54, 56, 57] lack is the clear explanation
400 of the variables that contribute to the common and distinct variability of the different data
401 blocks. To deal with this, two new methods recently emerged in the chemometrics domain, the
402 first is the multiblock variable influence on orthogonal projections (MB-VIOP) [57] and the

403 second is the sparse common and distinct covariate regression (SCD-CovR) [58]. The MB-
404 VIOP utilises the variable importance in projection (VIP) approach to sort the variables based
405 on their importance for the simplification and interpretation of the OnPLS model [57]. VIP is
406 performed on the global, common and distinct components extracted by OnPLS, thus reflecting
407 the key variables in the global, common and distinct fusion of multi-block data. The SCD-
408 CovR approach on the other hand combines the sparse principal covariate regression with the
409 simultaneous component analysis to extract the variables explaining common and distinct
410 variations in the multi-block data [58].

411 8. Multi-block analysis for higher order data fusion

412 Higher order data in analytical chemistry is commonly encountered [59]. Let us consider a
413 situation where N samples are analyzed by NIR spectroscopy, giving rise to spectra comprising
414 M variables at T time points. The resulting data structure is a cube of dimensions $N \times M \times T$. If
415 the same objects are analyzed by another platform, at only one time point, this leads to a data
416 matrix of dimensions $N \times K$, where K is the number of variables measured by another platform.
417 The resulting data set is a multi-block one, but the data structures present diverse
418 dimensionalities. To handle such data sets, the most straightforward solution would be to
419 unfold the cube into a matrix and to apply the traditional data fusion approaches. Nevertheless,
420 it has been demonstrated that, when modelling multi-way structures, it is better to leave their
421 natural dimensionality untouched, exploiting suitable methods for their analysis, rather than
422 unfold them out into two-dimensional arrays. In the light of this, multi-block methods for the
423 combination of arrays presenting different number of modes have been proposed. These multi-
424 block methods can be used for both exploratory and predictive purposes. Currently, multi-block
425 methods for unsupervised fusion of higher-order data are mostly based on coupled tensor and
426 matrix factorization approaches [43, 44]. In the context of predictive analysis, both approaches

427 based on multi-block multi-way covariate analysis [60] and the more recently proposed
428 extension of sequential and orthogonalized PLS regression to multi-way arrays (SO-N-PLS)
429 [59] are available. This latter approach resembles SO-PLS presented above, with the main
430 difference being that multi-way blocks are handled by means of N-PLS rather than PLS, to
431 maintain their multi-way structure.

432

433 9. Innovative uses of multi-block analysis

434 Apart from standard chemometric tasks such as exploratory analysis, regression, classification
435 and variable selection, other innovative applications of multi-block methods are emerging. Two
436 such applications are pre-processing selection and fusion, and calibration transfer. Pre-
437 processing selection is a key step in chemometric modeling, and it is largely debated since it is
438 difficult to define an optimal strategy. Often users struggle between different pre-processing
439 techniques to identify the best pre-processing or the best combination of pre-processing
440 techniques. A novel application of multi-block data analysis is to perform the fusion of multiple
441 pre-processing techniques [61-63], where the same data after pre-processing with different
442 methods can be considered as a multi-block dataset and can then be processed by multi-block
443 regression and classification. Recently, a technique called sequential pre-processing through
444 orthogonalization (SPORT) was proposed [64]. SPORT is based on the SO-PLS approach to
445 data fusion where the model learns in an incremental way the complementary information
446 present in different data blocks. Recent applications of SPORT can be found relating to
447 selection of pre-processing [64] and complementary fusion of scatter correction techniques
448 [61] in NIR spectroscopy. Since the SPORT approach is sequential, it is necessary to define
449 the order of pre-processing. The order can be decided upon, based on the complexity of pre-
450 processing techniques so that all easy, fast, and model-free techniques are used at the start and

451 the complex, slow, model-based techniques are reserved for the end. However, to deal with the
452 decision about application, a new pre-processing fusion approach called parallel pre-processing
453 through orthogonalization (PORTO) was proposed [65]. PORTO is based on the PO-PLS
454 procedure of predictive multi-block analysis and allows different pre-processing options and
455 their combinations to be explored in parallel. The PORTO approach has the advantage over the
456 SPORT approach in that it provides a better insight into the common and distinct information
457 highlighted by different pre-processing techniques. However, it has been reported that both the
458 SPORT and PORTO approaches usually lead to the same predictive performance. **The concept
459 of considering differently pre-processed versions of the same matrix as a multi-block data set
460 had already been considered in the framework of the MVP method [53]. In that context, the
461 use of MVP was advocated in order to get a deeper insight into which pre-processings could
462 carry similar information and which ones could possibly add a unique contribution,**

463 The second innovative application of multi-block data analysis is related to calibration transfer
464 (CT). CT is a widely explored task in chemometrics when the aim is to use a model developed
465 using one sensor, on another similar sensor. The aim of calibration transfer is to remove the
466 differences between the two instruments so that the model developed on one instrument can be
467 transferred and used with the other sensor. Recently, methods based on multi-block techniques
468 have emerged for calibration transfer [66]. A recent method called joint and unique multi-block
469 analysis (JUMBA) was proposed for the calibration transfer of NIR models [66]. The method
470 relies on the assumption that the two instruments have a major part of information in common
471 and a minor part that is distinct. Once the common information is identified by the multi-block
472 methods, the model developed on one sensor can be applied on the other.

473

474

475 *Table 1: A summary of multi-block methods available for multi-source data integration in*
 476 *chemometrics.*

Tasks	Data order	Methods	Background	Key features	References
Exploratory data analysis	2 nd order data	Consensus principal component analysis (CPCA)	<ul style="list-style-type: none"> • Global components are extracted maximizing the variance • Individual blocks are later regressed on the global components to extract the weights for individual blocks to have an insight on the contribution of each block to the global component 	<ul style="list-style-type: none"> • Weights of individual blocks provides importance for each block in the final model • Superweights normalized to length = 1 	[50]
		Extensions of multivariate curve resolution (MCR)	<ul style="list-style-type: none"> • The data blocks or matrices are concatenated along the common direction (rows, columns, both) • MCR is applied to the augmented data array 	<ul style="list-style-type: none"> • By suitable definition of the selectivity constraint, can extract both common and distinct components • Can deal with “incomplete” multi-sets (some matrices sharing the row-dimension and some others the column one) • Through the use of suitable constraints, can be used also for predictive modeling 	[47]
		Hierarchical principal component analysis (H-PCA)	<ul style="list-style-type: none"> • Like CPCA but relies on different normalization (superscores normalized to length = 1). • 	<ul style="list-style-type: none"> • Objective function is not clear • May provide different solutions depending on initialization. 	[29, 32]
		Common component and specific weight analysis (CCSWA)	<ul style="list-style-type: none"> • Global components are extracted sequentially by maximizing the variance of the weighted sum of cross-product matrices. • Individual blocks are later deflated on the global components, and the whole procedure is 	<ul style="list-style-type: none"> • On each dimension, weights of individual blocks indicate the importance of each block in the construction of the corresponding component • Both global and local components for 	[34]

			repeated on the deflated blocks.	each block can be obtained, as well as loadings for each block.	
		Multiple co-inertia analysis (MOCA)	<ul style="list-style-type: none"> • Data are preliminary transformed • Orthogonal components are extracted so to maximize the sum of the squared covariance with the scores of each block 	<ul style="list-style-type: none"> • Provides a simultaneous ordination of measurements and variables of multiple blocks 	[30]
		Orthogonal 2 partial least-squares (O2PLS)	<ul style="list-style-type: none"> • Preliminary estimation of common subspace by $\text{svd}(X_2^T X_1)$ • Orthogonalization of the blocks with respect to common subspace • Distinct component extracted from the orthogonalized blocks • After deflation of the distinct components, the final common components are extracted by a PLS-like step between the blocks 	<ul style="list-style-type: none"> • Common components are different between the blocks • No asymmetric relation between the blocks is assumed 	[41]
		Distinct and common simultaneous component analysis (DISCO-SCA)	<ul style="list-style-type: none"> • The joint subspace is extracted by SCA. • Target rotation of the block loadings is used to identify common and distinct components 	<ul style="list-style-type: none"> • Does not allow the extraction of partially shared components 	[39]
		Joint and individual variances explained (JIVE)	<ul style="list-style-type: none"> • Iterative extraction of common and distinct components • SVD on the concatenated data matrices to estimate the common components • Deflation of each block with respect to the common components • SVD on the deflated blocks to estimate the 	<ul style="list-style-type: none"> • The ranks of common and distinctive matrices are determined by permutation tests. 	[42]

			distinct components		
		Principal component analysis Generalized canonical analysis (PCA-GCA)	<ul style="list-style-type: none"> • Preliminary PCA on individual blocks to filter out noise. • Finds linear combination of the blocks which best fit to a set of orthogonal common components 	<ul style="list-style-type: none"> • Focuses on common components • Distinctive components are obtained by PCA on the residual matrices after regressing each block on the common components. 	[3]
		Generalized singular value decomposition (GSVD)	<ul style="list-style-type: none"> • Preliminary SCA to filter out noise. • Joint SVD of the different data matrices • Identification of common and distinct components based on the singular values 	<ul style="list-style-type: none"> • Originally proposed for multi-set data sharing the variable dimensions. 	[40]
		Structured learning and integrative decomposition (SLIDE)	<ul style="list-style-type: none"> • Loadings are organized in a block-dependent structure • Structure sparsity is imposed to reveal the common, distinct and the partially shared information 	<ul style="list-style-type: none"> • Can be considered as an intermediate model between SUM-PCA and JIVE • Components common only to some blocks can be extracted 	[45]
		Penalized exponential simultaneous component analysis (P-ESCA)	<ul style="list-style-type: none"> • Penalties are incorporated in the simultaneous component analysis for separating common and distinct information in the multi-block data 	<ul style="list-style-type: none"> • Common and distinct variation in the data in be explored separately 	[21]
		Penalized exponential analysis of variance simultaneous component analysis (P-EASCA)	<ul style="list-style-type: none"> • Combines the penalized exponential simultaneous component analysis with the analysis of variance simultaneous component analysis • The multi-block data is decomposed into common and distinct part and later the ASCA is used for 	<ul style="list-style-type: none"> • Only multi-block technique available for exploration of designed experimental data 	[19]

			exploratory analysis of common and distinct variation		
	Higher order data	Combined matrix and tensor factorization	<ul style="list-style-type: none"> Combines the matrix factorization and tensor factorization 	<ul style="list-style-type: none"> Data of multiple order such as 2D, 3D etc. can be jointly explored 	[43, 44]
Predictive analysis	2 nd order data	Multi-block partial least-squares regression	<ul style="list-style-type: none"> Global components are extracted maximizing the covariance with the response variable(s) Individual blocks are later regressed on the global components to extract the weights for individual blocks to have an insight into the contribution of each block to the global component 	<ul style="list-style-type: none"> Weights of individual blocks provides importance for each block in the final model Block weights and superweights are normalized to unit length 	[48]
		Hierarchical or consensus partial least-squares regression	<ul style="list-style-type: none"> A CPCA cycle is performed on the multiple X blocks. A PLS cycle is done between the superscores and the response(s) 	<ul style="list-style-type: none"> Superscores are normalized to unit length 	[29, 50]
		Orthogonal n partial least-squares (OnPLS) regression	<ul style="list-style-type: none"> Extension of O2PLS to the multi-block scenario A global regression model is calculated between the block containing the responses to be predicted, and the scores extracted from all the other matrices 	<ul style="list-style-type: none"> Can be also used for exploratory analysis, since no asymmetric relations between the blocks are assumed a priori 	[51]
		ComDim (k+1), or P-ComDim	<ul style="list-style-type: none"> Two sets of global components (for the predictor blocks and for the response blocks) are extracted sequentially by maximizing the variance of the sum of cross-product matrices involving both predictor and response blocks. Individual blocks are later deflated on the global components, and the whole procedure is 	<ul style="list-style-type: none"> The weights (salience) of each block on each dimension, indicates its importance in the determination of that common component 	[37]

			repeated on the deflated blocks.		
		Sequential orthogonal partial least-squares (SO-PLS) regression	<ul style="list-style-type: none"> Includes a combination of partial least-squares regression and sequential orthogonalization step to extract complementary latent variables from multi-block data 	<ul style="list-style-type: none"> Complementary unique information is extracted 	[52]
		Parallel orthogonal partial least-squares (PO-PLS) regression	<ul style="list-style-type: none"> Includes a combination of partial least-squares regression, canonical correlation analysis and orthogonalization step to extract common and distinct latent variables from multi-block data 	<ul style="list-style-type: none"> Common and distinct information can be extracted Good for the cases when order of block is not important or all blocks are of equal importance 	[52]
		Multi-block variance partitioning (MVP)	<ul style="list-style-type: none"> Individual PLS models between each predictor block X_k and Y For each block, unique Y-related variation is obtained by orthogonalizing the predicted responses with respect to the corresponding responses predicted using the other blocks Common variation is obtained by subtracting the unique part and the residuals from the total variance 	<ul style="list-style-type: none"> Extracts common and distinct information Scale invariant Can be extended to evaluate the performances of preprocessing methods 	[53]
Higher order data	Multi-way multi-block covariates regression	<ul style="list-style-type: none"> Extension of principal covariate regression Scores are extracted so to explain the variation in their associated block and convey similarities between the blocks 	<ul style="list-style-type: none"> A different number of scores can be extracted from each block The extent to which inter- and intra-block variation is accounted for is regulated by a metaparameter. 	[60]	

		Sequential orthogonal N-way partial least-squares regression	<ul style="list-style-type: none"> Includes a combination of PLS regression or N-PLS regression (depending on the nature of the block to model) to sequentially extract complementary latent variables from multi-way multi-block data 	<ul style="list-style-type: none"> Complementary information is extracted 	[59]
Variable selection	2 nd order data	Variable importance in projection (VIP) + SO-PLS	<ul style="list-style-type: none"> The method is based on estimating the variable importance on the components extracted by the sequential orthogonalized partial least-squares (SO-PLS) regression 	<ul style="list-style-type: none"> The variables can be extracted with simultaneous SO-PLS modelling 	[54]
		Sequential orthogonalized covariate selection (SO-CovSel)	<ul style="list-style-type: none"> Based on the sequential covariance maximization and orthogonalization step to select variables across multiple blocks of data 	<ul style="list-style-type: none"> The approach is sequential so data blocks based on their importance can be arranged by user Discrete variables are selected which can be used for developing cheap sensors 	[55]
		Multi-block variable important in orthogonal projections (MB-VIOP)	<ul style="list-style-type: none"> The method is based on estimating the variable importance on the common and distinct components extracted by the orthogonal n partial least-squares (OnPLS) regression 	<ul style="list-style-type: none"> Allows exploration of common and distinct variables amount different blocks of data 	[57]
		Sparse common and distinct covariate regression (SCD-PCovR)	<ul style="list-style-type: none"> Combined the sparse principal covariate regression with the simultaneous component analysis to extract the common and distinct variables in multi-block data 	<ul style="list-style-type: none"> Allows exploration of common and distinct variables amount different blocks of data Sparsity parameter can be tuned by user to command the 	[58]

				variable selection	
Calibration transfer	2 nd order data	The joint and unique multi-block analysis (JUMBA) for calibration transfer	<ul style="list-style-type: none"> Based on the framework of Orthogonal n partial least-squares (OnPLS) to identify common and distinct information in multi-block data 	<ul style="list-style-type: none"> Allows multi-instrument calibration transfer Enhanced understanding about intrinsic differences of instruments can be gained within the framework of joint and unique information in multi-block data 	[66]
Pre-processing optimization and fusion	2 nd order data	Sequential pre-processing through orthogonalization (SPORT)	<ul style="list-style-type: none"> Based on sequential orthogonalized partial least-squares regression Includes a combination of partial least-squares regression and sequential orthogonalization step to extract complementary latent variables from multi-block data 	<ul style="list-style-type: none"> Allows a sequential fusion of pre-processing techniques Complementary information is modelled Pre-processing techniques selection can be performed Any pre-processing techniques having zero LVs can be discarded, thus, leading to pre-processing selection 	[64]
		Parallel pre-processing through orthogonalization (PORTO)	<ul style="list-style-type: none"> Based on parallel orthogonalized partial least-squares regression Includes a combination of partial least-squares regression, canonical correlation analysis and orthogonalization step to extract common and distinct latent variables from multi-block data 	<ul style="list-style-type: none"> Allows parallel fusion of pre-processing technique without the need of defining the order Allows insight to the common and distinct information present in different pre-processing techniques which help in selecting a subset of pre-processing techniques 	[65]

477

478 10. Free software resources available for multi-block data 479 analysis

480 Multi-block data analysis is a relatively new domain in chemometrics and software for
481 performing the multi-block analysis are scarce. However, several groups around the world have
482 published freely available codes so that the community can benefit from them. The first
483 publicly available MATLAB-based toolbox is from the University of Copenhagen, Denmark
484 (<http://www.models.life.ku.dk/~courses/MBtoolbox/mbtmain.htm>), which, having been last
485 revised in 2001, focuses only on the two data fusion approaches that were most popular at that
486 time, i.e., multi-block principal component analysis and multi-block partial least squares
487 regression. The second toolbox is that by NOFIMA for multi-block regression by parallel and
488 sequential partial least-squares regression [52]. Both toolboxes provide command line
489 functionalities (within the MATLAB environment) and consist of a limited number of tools.
490 There is also a basic graphical user interface (GUI) available for performing multi-block
491 component analysis in the domain of behavioural research [67]. However, this GUI only
492 proposes principal component analysis on each data block separately, simultaneous component
493 analysis, and cluster-wise simultaneous component analysis for data exploration. Recently, a
494 new graphical user interface, the MBA-GUI (freely available at:
495 <https://github.com/puneetmishra2/Multi-block.git>) has been made available which integrates
496 several advanced multi-block techniques related to data exploration, regression, variable
497 selection, pre-processing selection and fusion [26]. A python library called ‘mbpls 1.0.4’ was
498 also recently developed for performing multi-block PLS and is available at
499 <https://pypi.org/project/mbpls/>.

11. Some comparative examples for predictive multi-block modelling

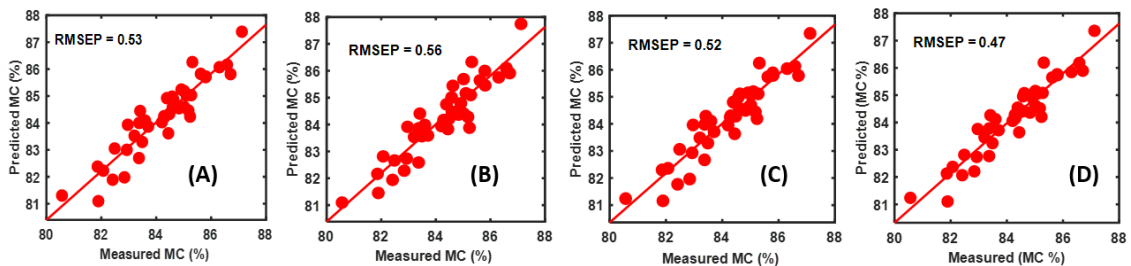
Thanks to the chemometrics developments in recent years, multi-block techniques are now available to perform both exploratory and predictive data modelling. However, there are so many new tools, as well as extensions of standard single-block chemometric techniques (such as, MCR, PLS), that it is becoming difficult to find the best solutions to start with when dealing with a new problem/application. However, in chemometrics, as in other scientific fields, it is difficult if not impossible to univocally define what could be the best technique in an absolute sense. First of all, as the term “best”, itself, can assume several meanings depending on the specific application (e.g., more robust, less impacted by interferences, more accurate, and so on). A well, different techniques have advantages and disadvantages for different data type. For this reason, these different tools could even be ensembled to get a better understanding of the data and solve the background challenges. Recently, several works have tried to compare the performances of different multi-block methods to achieve a deeper understanding of their characteristics, similarities and dissimilarities, in the light of the practical use of those techniques [3, 7, 20, 21, 35]. Since most of these works focus on exploratory approaches (i.e., to symmetric data fusion), while to the authors’ knowledge, so far, no research literature provides a comparative overview of the predictive approaches, hence, three such practical comparisons are provided in the following section. The comparisons are based on the pear data set where a total of 231 pear fruit were measured with two complementary near-infrared (NIR) spectral sensors covering the spectral ranges of ~700-1050 nm and ~1050 to 1600 nm, respectively. The reference property was the moisture content (MC) measured with hot-air oven drying technique [68]. The samples were divided into a calibration set and an independent test set of 190 and 41 individuals, respectively. The data set used in these examples has already been published and it is used here for demonstration purposes only. More details on sampling

526 and reference analysis can be accessed in the original publication related to this data [68]. Out
527 of the three comparisons, the first is the comparative overview of PLS and two multiblock PLS
528 methods (SO-PLS and PO-PLS) [52], the second example is the comparison of two multi-block
529 pre-processing fusion approaches, i.e., SPORT [64] and PORTO [65], and the third is an
530 example of multi-block variable selection with SO-CovSel [55].

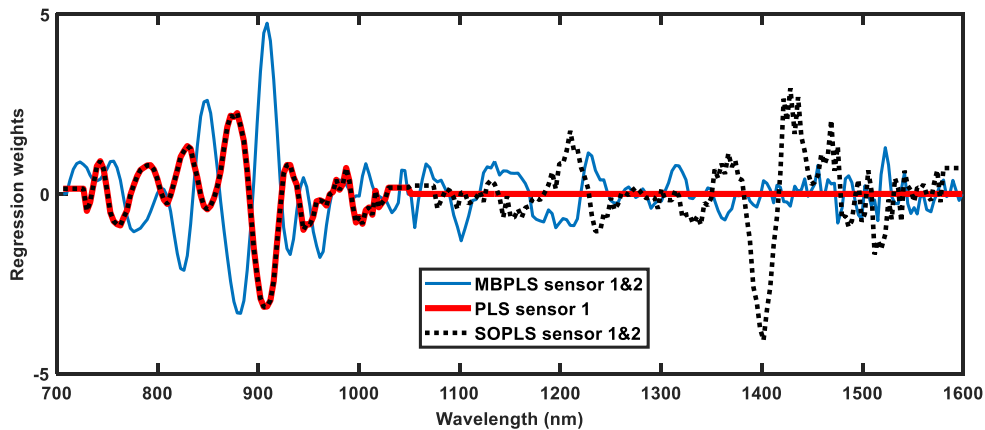
531 11.1. Comparison of PLS, MB-PLS, SO-PLS and PO-PLS

532 A summary of the performances of the different approaches to model the spectral data from
533 multiple complementary NIR sensors is shown in Fig. 6. As a baseline, the results of PLS
534 modeling of the spectra from sensor 1 only (the best model on the individual matrices) was
535 added to show that this block alone is not sufficient to achieve an error as low as the one
536 obtained through the fusion of the complementary information in the multiple sensors. The
537 PLS model built on sensor 1 data only with 7 LVs (optimized using 5-fold cross-validation)
538 (Fig. 6A) attained a root mean square error (RMSEP) of 0.53 % to predict MC. Furthermore,
539 when the data from two sensors was concatenated along the variable direction and a new PLS
540 model was developed (Fig. 6B), the RMSEP was slightly increased, from 0.53 to 0.56 %, SO-
541 PLS also showed a similar RMSEP to that of PLS on concatenated spectral data (MB-PLS),
542 but by extracting complementary information from two sensors. On the other hand, in this case,
543 PO-PLS resulted in the lowest RMSEP (0.47 %). Furthermore, PO-PLS obtained this superior
544 performance by partitioning the common and unique information in the two sensors. However,
545 a key point to note is that the performance of MB-PLS was poorer than that of the SO-PLS
546 approach (i.e., a RMSEP of 0.52 %). Methods like SO-PLS and PO-PLS allow efficient
547 modelling of different data matrices and can lead to more accurate predictions than MB-PLS.
548 Advanced multi-block methods also bring added values such as a better understanding of
549 background chemistry, which can also be noted in Fig. 7, where the regression coefficients
550 corresponding to MB-PLS (dotted blue line) and SO-PLS (dashed black line) are presented. It

551 can be noted that calculating MB-PLS resulted in a model with higher absolute coefficients for
 552 sensor 1 data and several key features in the spectral range of sensor 2 are poorly modelled,
 553 e.g. at 1400 nm, which corresponds to the H₂O overtones directly related to the moisture [69].
 554 A main challenge with MB-PLS is the need to perform proper scaling of the data, but that is
 555 not the case with methods like SO-PLS as they treat each data block sequentially [8], thus
 556 avoiding any negative effect of different data scales.



557
 558 *Fig. 6: Pear data set – Comparison of model performances for the prediction of moisture*
 559 *content (MC). Predicted vs observed MC values (test set) for: (A) the best PLS model on*
 560 *individual blocks; (B) MB-PLS; (C) SO PLS; and (D) PO-PLS.*



561
 562 *Fig. 7.: Pear data set – Comparison of models for the prediction of moisture content (MC).*
 563 *Regression coefficients for the PLS model built on sensor 1 data only (continuous thick red*
 564 *line), MB-PLS on concatenated data from sensors 1 and 2 (continuous thin blue line), and*
 565 *SOPLS on sensor 1 and 2 data (dotted black line; the coefficient vectors for the two separate*
 566 *PLS regressions involved have been concatenated, for a better visualization).*

567

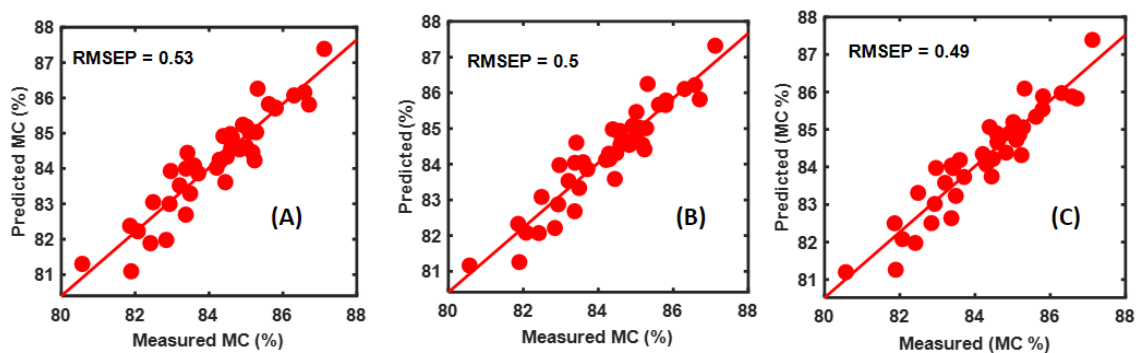
568 11.2. Comparison of pre-processing fusion approaches

569 Pre-processing selection in chemometrics is a challenging task, where a lot of time and
570 resources are usually spent with the aim of achieving optimal pre-processing combinations
571 [70]. However, such an approach can be considered old-fashioned due to emergence of new
572 ensemble pre-processing fusion approaches [70] and especially the multi-block data analysis
573 inspired methods such as SPORT [64] and PORTO [65]. A comparison of the use of SPORT
574 and PORTO on the pear data set already described is shown in Fig. 8, as an example. The data
575 used in this analysis are only those from sensor 1, although SPORT and PORTO can both deal
576 with simultaneous multi-sensor multi-processing. Furthermore, only two data blocks, i.e. raw
577 data and data pre-processed by 2nd derivative, are used for the demonstration. The main thing
578 to note is that using 2nd derivative only (best individual pre-processing) results in a higher value
579 of the RMSEP (0.53 %) (Fig. 8A), whereas modelling with both the raw and 2nd derivative pre-
580 processed data gave a lower RMSEP, 0.50 % for SPORT (Fig. 8B) and 0.49 % with PORTO
581 (Fig. 8C).

582 The good performance from the combined used of raw and 2nd derivative pre-processed data is
583 not a surprise from a fruit property modelling point of view. This is because the NIR spectra
584 of fresh fruit are a mixture of absorption and scattering profiles, the absorption can usually be
585 related to broad peaks in the NIR data, whereas the scattering properties are expressed as the
586 additive and multiplicative effects [71]. The chemical and physical properties of fruits are
587 correlated to both the effect of scattering due to fruit cellular structure (which differs with the
588 ripening stage of fruit) and absorption present in NIR data. Hence, doing a 2nd derivative
589 estimation may eliminate the global intensity differences related to scattering, and therefore,
590 may remove some useful information related to fruit properties. The multi-block approaches
591 compensate this loss of information by first modeling the 2nd derivative and then modelling the

592 remaining variation within the raw data. An example of the complementary modelling
593 performed by SO-PLS is shown in Fig. 9, where the regression vectors for the 2nd derivative
594 pre-processed (solid blue line) and raw data (dashed red line) are shown. It can be seen that the
595 main features of the 2nd derivative are the peaks related to overtones of -OH, -CH and -NH
596 [69], whereas the main information captured from the raw data is the global shape of the
597 spectrum which is an indication of the scattering information.

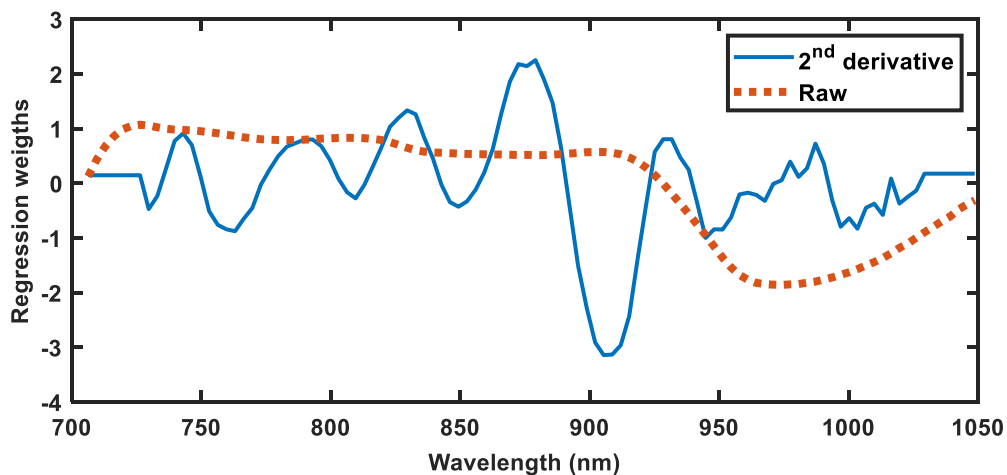
598



599

600 *Fig. 8: Pear data set – Comparison of model performances for the prediction of moisture*
601 *content (MC). Predicted vs observed MC values (test set) for: (A) PLS (on data pre-processed*
602 *with 2nd derivative only), (B) SPORT, and (C) PORTO models.*

603



604

605 *Fig. 9: Pear data set – SPORT model. Regression coefficients for the two blocks of data (raw*

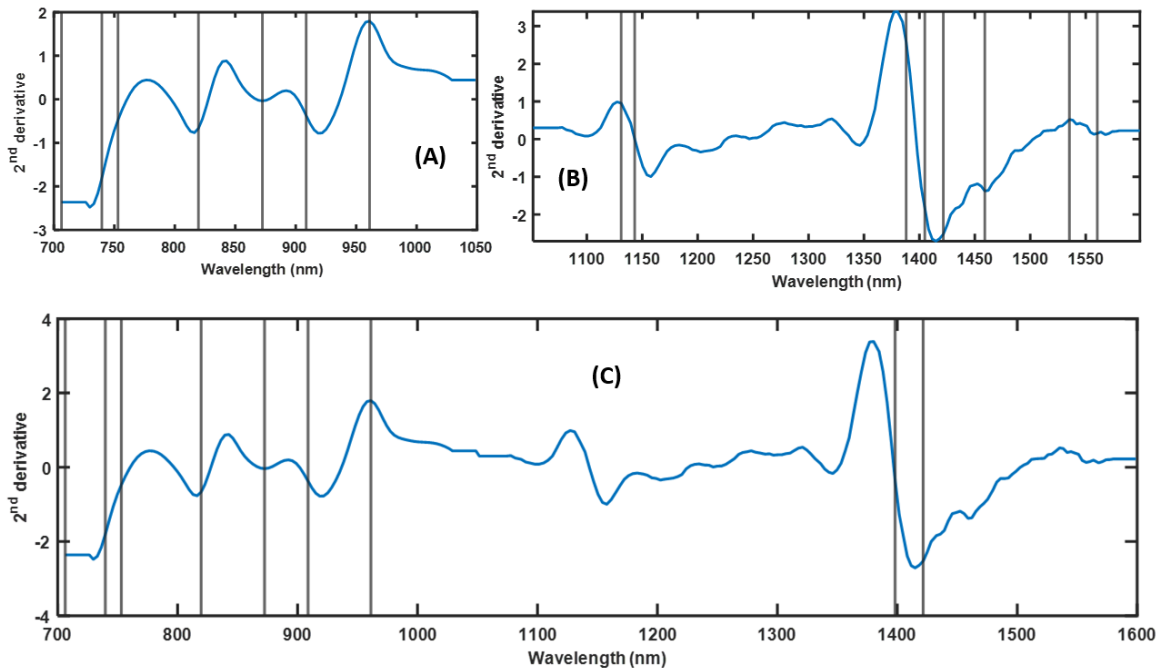
606 *and 2nd derivative pre-processed*).

607

608 11.3. Selecting variables in multi-block scenario

609 Variable selection is useful in chemometrics and is even challenging when the data is multi-
610 block. Particularly, the challenge arises when the redundant information is present in multiple
611 data blocks and the aim is to just use the complementary information that improves the
612 predictive performances of the model. In such a case, new multi-block methods such as SO-
613 CovSel can be used efficiently. Fig. 10 C shows the results of performing SO-CovSel on the
614 two-sensors pear data set. For the sake of comparison, CovSel analysis on individual blocks is
615 also presented and the selected variable for sensor 1 and sensor 2 are shown in Fig. 10 A and
616 B, respectively. Separate CovSel analyses on sensor 1 and sensor 2 data selected 7 and 8
617 wavelengths, respectively. However, most of the wavelengths are related to overtones of
618 similar chemical bonds and, therefore, carry redundant information. In the case of SO-CovSel,
619 due to such redundancy, only 2 bands are selected from sensor 2 and, nevertheless, the RMSEP
620 is reduced from 0.58 % (best individual model) to 0.55 % (Fig. 11).

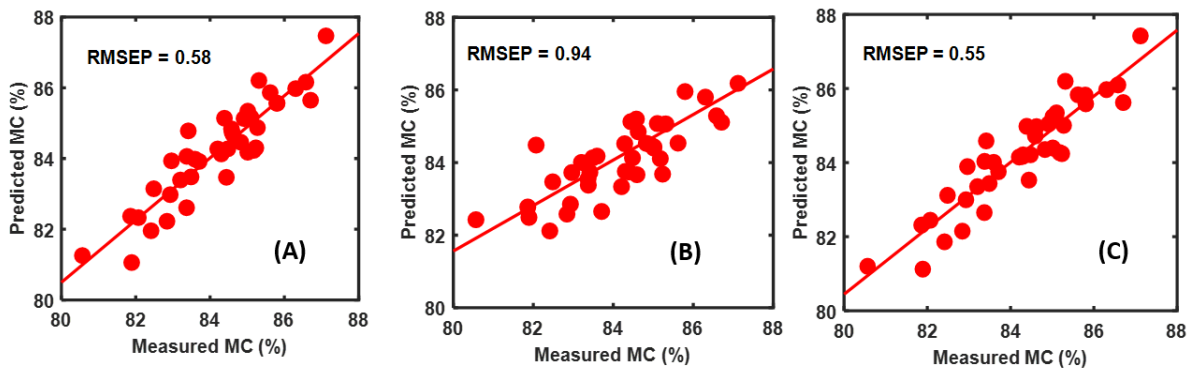
621



622

623 *Fig. 10: Pear data set - Comparison of the results of variable selection on the individual blocks*
 624 *and in the multi-block scenario. (A) Variables selected from sensor 1 data through single-block*
 625 *CovSel analysis, (B) Variables selected from sensor 2 data through single-block CovSel*
 626 *analysis, and (C) Variables jointly selected from sensor 1 and sensor 2 data through the multi-*
 627 *block SO-CovSel approach.*

628



629

630 *Fig. 11: Pear data set - Comparison of the results of variable selection on the individual blocks*
 631 *and in the multi-block scenario for the prediction of moisture prediction in pear fruit. Predicted*
 632 *vs measured values of moisture content based on (A) Single-block CovSel analysis on sensor 1*
 633 *data, (B) Single-block CovSel analysis on sensor 2 data, and (C) Multi-block SO-CovSel*

634 *approach.*

635

636 12. Concluding remarks

637 Multi-block data analysis in chemometrics is gaining increasing attention and the development
638 of new methods in recent years has been rapid. Analytical chemistry can directly benefit from
639 these new techniques to explore and combine data from multiple sources. Due to advances in
640 sensor and computing technologies, multi-source data are now frequently encountered. Multi-
641 block methods are available for diverse tasks such as exploratory data analysis, predictive
642 modelling, variable selection, pre-processing optimization, and calibration transfer. There are
643 also methods available to explore the multi-block data generated by designed experiments,
644 which is often the case with lab-based classical analytical chemistry experiments. The main
645 benefit of multi-block data analysis compared to the standard chemometric methods is that they
646 allow a detailed understanding of common and distinct information present in different data-
647 blocks or data generated from multiple sources. Recently, free software tools such as the MBA-
648 GUI have been made available to the scientific community to explore the possibilities of multi-
649 block data analysis. It can be expected that the future trend will be an exponential increase in
650 the applications of multi-block data analysis methods in analytical chemistry to combine in an
651 optimal way multiple sources of data. *Another important direction that can be foreseen is the
652 development of interactive data visualization tools [72] dedicated to multi-block data analysis,
653 which will allow even non-experts to have a better comprehension of their data.*

654

655 Declaration of Interest

656 None

- 658 [1] A. Biancolillo, R. Bucci, A.L. Magrì, A.D. Magrì, F. Marini, Data-fusion for multiplatform
659 characterization of an Italian craft beer aimed at its authentication, *Analytica Chimica Acta*, 820 (2014)
660 23-31.
- 661 [2] L. Zhou, C. Zhang, Z. Qiu, Y. He, Information fusion of emerging non-destructive analytical
662 techniques for food quality authentication: A survey, *TrAC Trends in Analytical Chemistry*, 127 (2020)
663 115901.
- 664 [3] A.K. Smilde, I. Måge, T. Næs, T. Hankemeier, M.A. Lips, H.A.L. Kiers, E. Acar, R. Bro, Common and
665 distinct components in data fusion, *Journal of Chemometrics*, 31 (2017) e2900.
- 666 [4] A. Biancolillo, I. Måge, T. Næs, Combining SO-PLS and linear discriminant analysis for multi-block
667 classification, *Chemometrics and Intelligent Laboratory Systems*, 141 (2015) 58-67.
- 668 [5] I. Måge, E. Menichelli, T. Næs, Preference mapping by PO-PLS: Separating common and unique
669 information in several data blocks, *Food Quality and Preference*, 24 (2012) 8-16.
- 670 [6] M. Ramos-Barberán, M.V. Hinojosa-Ramos, J. Ascencio-Moreno, F. Vera, O. Ruiz-Barzola, M.P.
671 Galindo-Villardón, Batch process control and monitoring: a Dual STATIS and Parallel Coordinates (DS-
672 PC) approach, *Production & Manufacturing Research*, 6 (2018) 470-493.
- 673 [7] R. Vitale, O.E. de Noord, J.A. Westerhuis, A.K. Smilde, A. Ferrer, Divide et impera: How disentangling
674 common and distinctive variability in multiset data analysis can aid industrial process troubleshooting
675 and understanding, *Journal of Chemometrics*, n/a (2020) e3266.
- 676 [8] M.P. Campos, M.S. Reis, Data preprocessing for multiblock modelling – A systematization with new
677 methods, *Chemometrics and Intelligent Laboratory Systems*, 199 (2020) 103959.
- 678 [9] Z. Ge, Review on data-driven modeling and monitoring for plant-wide industrial processes,
679 *Chemometrics and Intelligent Laboratory Systems*, 171 (2017) 16-25.
- 680 [10] R.G. Brereton, J. Jansen, J. Lopes, F. Marini, A. Pomerantsev, O. Rodionova, J.M. Roger, B. Walczak,
681 R. Tauler, Chemometrics in analytical chemistry—part I: history, experimental design and data analysis
682 tools, *Analytical and Bioanalytical Chemistry*, 409 (2017) 5891-5899.
- 683 [11] R.G. Brereton, J. Jansen, J. Lopes, F. Marini, A. Pomerantsev, O. Rodionova, J.M. Roger, B. Walczak,
684 R. Tauler, Chemometrics in analytical chemistry—part II: modeling, validation, and applications,
685 *Analytical and Bioanalytical Chemistry*, 410 (2018) 6691-6704.
- 686 [12] R. Bro, A.K. Smilde, Principal component analysis, *Analytical Methods*, 6 (2014) 2812-2831.
- 687 [13] P. Geladi, B.R. Kowalski, Partial least-squares regression: a tutorial, *Analytica Chimica Acta*, 185
688 (1986) 1-17.
- 689 [14] T. Mehmood, K.H. Liland, L. Snipen, S. Sæbø, A review of variable selection methods in Partial
690 Least Squares Regression, *Chemometrics and Intelligent Laboratory Systems*, 118 (2012) 62-69.
- 691 [15] T. Mehmood, S. Sæbø, K.H. Liland, Comparison of variable selection methods in partial least
692 squares regression, *Journal of Chemometrics*, n/a (2020) e3226.
- 693 [16] J.M. Roger, B. Palagos, D. Bertrand, E. Fernandez-Ahumada, CovSel: Variable selection for highly
694 multivariate and multi-response calibration: Application to IR spectroscopy, *Chemometrics and
695 Intelligent Laboratory Systems*, 106 (2011) 216-223.
- 696 [17] A.K. Smilde, J.A. Westerhuis, S. de Jong, A framework for sequential multiblock component
697 methods, *Journal of Chemometrics*, 17 (2003) 323-337.
- 698 [18] A. Biancolillo, T. Næs, The Sequential and Orthogonalized PLS Regression for Multiblock
699 Regression: Theory, Examples, and Extensions. In M. Cocchi (Ed.), *Data Fusion Methodology and
700 Applications*, Elsevier, 2019, pp. 157-177.
- 701 [19] M. Alinaghi, H.C. Bertram, A. Brunse, A.K. Smilde, J.A. Westerhuis, Common and distinct variation
702 in data fusion of designed experimental data, *Metabolomics*, 16 (2019) 2.
- 703 [20] I. Måge, A.K. Smilde, F.M. van der Kloet, Performance of methods that separate common and
704 distinct variation in multiple data blocks, *Journal of Chemometrics*, 33 (2019) e3085.
- 705 [21] Y. Song, J.A. Westerhuis, A.K. Smilde, Separating common (global and local) and distinct variation
706 in multiple mixed types data sets, *Journal of Chemometrics*, 34 (2020) e3197.

707 [22] S. Zhu, Z. Song, S. Shi, M. Wang, G. Jin, Fusion of Near-Infrared and Raman Spectroscopy for In-
708 Line Measurement of Component Content of Molten Polymer Blends, *Sensors (Basel, Switzerland)*, 19
709 (2019) 3463.

710 [23] S.E. Barnes, E.C. Brown, M.G. Sibley, H.G.M. Edwards, I.J. Scowen, P.D. Coates, Vibrational
711 Spectroscopic and Ultrasound Analysis for In-Process Characterization of High-Density
712 Polyethylene/Polypropylene Blends during Melt Extrusion, *Applied Spectroscopy*, 59 (2005) 611-619.

713 [24] K. Haroon, A. Arafeh, S. Cunliffe, P. Martin, T. Rodgers, C. Mendoza, M. Baker, Comparison of
714 Individual and Integrated Inline Raman, Near-Infrared, and Mid-Infrared Spectroscopic Models to
715 Predict the Viscosity of Micellar Liquids, *Applied Spectroscopy*, 74 (2020) 819-831.

716 [25] C. Assis, H.V. Pereira, V.S. Amador, R. Augusti, L.S. de Oliveira, M.M. Sena, Combining mid infrared
717 spectroscopy and paper spray mass spectrometry in a data fusion model to predict the composition
718 of coffee blends, *Food Chemistry*, 281 (2019) 71-77.

719 [26] P. Mishra, J.M. Roger, D.N. Rutledge, A. Biancolillo, F. Marini, A. Nordon, D. Jouan-Rimbaud-
720 Bouveresse, MBA-GUI: A chemometric graphical user interface for multi-block data visualisation,
721 regression, classification, variable selection and automated pre-processing, *Chemometrics and
722 Intelligent Laboratory Systems*, (2020) 104139.

723 [27] J.-M. Roger, J.-C. Boulet, M. Zeaiter, D.N. Rutledge, Pre-processing Methods. n: S.D. Brown, R.
724 Tauler, B. Walczak (Eds.), *Comprehensive Chemometrics (Second Edition)*, vol. 3, Elsevier, Oxford,
725 2020, pp.1-75.

726 [28] M.P. Campos, R. Sousa, A.C. Pereira, M.S. Reis, Advanced predictive methods for wine age
727 prediction: Part II – A comparison study of multiblock regression approaches, *Talanta*, 171 (2017) 132-
728 142.

729 [29] J.A. Westerhuis, T. Kourti, J.F. MacGregor, Analysis of multiblock and hierarchical PCA and PLS
730 models, *Journal of Chemometrics*, 12 (1998) 301-321.

731 [30] M. Hanafi, A. Kohler, E.-M. Qannari, Connections between multiple co-inertia analysis and
732 consensus principal component analysis, *Chemometrics and Intelligent Laboratory Systems*, 106
733 (2011) 37-40.

734 [31] M. Hanafi, E.M. Qannari, B. Jaillais, Multi-Block and Three-Way Data Analysis. In: S.D. Brown, R.
735 Tauler, B. Walczak (Eds.), *Comprehensive Chemometrics (Second Edition)*, vol. 3, Elsevier, Oxford,
736 2020, pp. 341-358.

737 [32] S. Wold, N. Kettaneh, K. Tjessem, Hierarchical multiblock PLS and PC models for easier model
738 interpretation and as an alternative to variable selection, *Journal of Chemometrics*, 10 (1996) 463-
739 482.

740 [33] E.M. Qannari, I. Wakeling, P. Courcoux, H.J.H. MacFie, Defining the underlying sensory
741 dimensions, *Food Quality and Preference*, 11 (2000) 151-154.

742 [34] M. Hanafi, G. Mazerolles, E. Dufour, E.M. Qannari, Common components and specific weight
743 analysis and multiple co-inertia analysis applied to the coupling of several measurement techniques,
744 *Journal of Chemometrics*, 20 (2006) 172-183.

745 [35] V. Cariou, D. Jouan-Rimbaud Bouveresse, E.M. Qannari, D.N. Rutledge, ComDim Methods for the
746 Analysis of Multiblock Data in a Data Fusion Perspective. In M. Cocchi (Ed.), *Data Fusion Methodology
747 and Applications*, Elsevier, 2019, pp. 179-204.

748 [36] D. Jouan-Rimbaud Bouveresse, R.C. Pinto, L.M. Schmidtke, N. Locquet, D.N. Rutledge,
749 Identification of significant factors by an extension of ANOVA-PCA based on multi-block analysis,
750 *Chemometrics and Intelligent Laboratory Systems*, 106 (2011) 173-182.

751 [37] A. El Ghaziri, V. Cariou, D.N. Rutledge, E.M. Qannari, Analysis of multiblock datasets using
752 ComDim: Overview and extension to the analysis of $(K + 1)$ datasets, *Journal of Chemometrics*, 30
753 (2016) 420-429.

754 [38] V. Cariou, E.M. Qannari, D.N. Rutledge, E. Vigneau, ComDim: From multiblock data analysis to
755 path modeling, *Food Quality and Preference*, 67 (2018) 27-34.

756 [39] M. Schouteden, K. Van Deun, S. Pattyn, I. Van Mechelen, SCA with rotation to distinguish common
757 and distinctive information in linked data, *Behavior Research Methods*, 45 (2013) 822-833.

758 [40] K. Van Deun, I. Van Mechelen, L. Thorrez, M. Schouteden, B. De Moor, M.J. van der Werf, L. De
759 Lathauwer, A.K. Smilde, H.A.L. Kiers, DISCO-SCA and Properly Applied GSVD as Swinging Methods to
760 Find Common and Distinctive Processes, PLOS ONE, 7 (2012) e37840.

761 [41] J. Trygg, O2-PLS for qualitative and quantitative analysis in multivariate calibration, Journal of
762 Chemometrics, 16 (2002) 283-293.

763 [42] E.F. Lock, K.A. Hoadley, J.S. Marron, A.B. Nobel, JOINT AND INDIVIDUAL VARIATION EXPLAINED
764 (JIVE) FOR INTEGRATED ANALYSIS OF MULTIPLE DATA TYPES, The annals of applied statistics, 7 (2013)
765 523-542.

766 [43] E. Acar, M.A. Rasmussen, F. Savorani, T. Næs, R. Bro, Understanding data fusion within the
767 framework of coupled matrix and tensor factorizations, Chemometrics and Intelligent Laboratory
768 Systems, 129 (2013) 53-63.

769 [44] E. Acar, E.E. Papalexakis, G. Gürdeniz, M.A. Rasmussen, A.J. Lawaetz, M. Nilsson, R. Bro, Structure-
770 revealing data fusion, BMC Bioinformatics, 15 (2014) 239.

771 [45] I. Gaynanova, G. Li, Structural learning and integrative decomposition of multi-view data,
772 Biometrics, 75 (2019) 1121-1132.

773 [46] A.K. Smilde, J.J. Jansen, H.C.J. Hoefsloot, R.-J.A.N. Lamers, J. van der Greef, M.E. Timmerman,
774 ANOVA-simultaneous component analysis (ASCA): a new tool for analyzing designed metabolomics
775 data, Bioinformatics, 21 (2005) 3043-3048.

776 [47] R. Tauler, M. Maeder, A. de Juan, Multiset Data Analysis: Extended Multivariate Curve Resolution.
777 In: S.D. Brown, R. Tauler, B. Walczak (Eds.), Comprehensive Chemometrics (Second Edition), vol.2,
778 Elsevier, Oxford, 2020, pp. 305-336.

779 [48] L.E. Wangen, B.R. Kowalski, A multiblock partial least squares algorithm for investigating complex
780 chemical systems, Journal of Chemometrics, 3 (1989) 3-20.

781 [49] S.J. Qin, S. Valle, M.J. Piovoso, On unifying multiblock analysis with application to decentralized
782 process monitoring, Journal of Chemometrics, 15 (2001) 715-742.

783 [50] S. Wold, PLS Modeling with Latent Variables in Two Or More Dimensions, 1987.

784 [51] T. Löfstedt, J. Trygg, OnPLS—a novel multiblock method for the modelling of predictive and
785 orthogonal variation, Journal of Chemometrics, 25 (2011) 441-455.

786 [52] T. Næs, O. Tomic, N.K. Afseth, V. Segtnan, I. Måge, Multi-block regression based on combinations
787 of orthogonalisation, PLS-regression and canonical correlation analysis, Chemometrics and Intelligent
788 Laboratory Systems, 124 (2013) 32-42.

789 [53] T. Skov, D. Ballabio, R. Bro, Multiblock variance partitioning: A new approach for comparing
790 variation in multiple data blocks, Analytica Chimica Acta, 615 (2008) 18-29.

791 [54] A. Biancolillo, K.H. Liland, I. Måge, T. Næs, R. Bro, Variable selection in multi-block regression,
792 Chemometrics and Intelligent Laboratory Systems, 156 (2016) 89-101.

793 [55] A. Biancolillo, F. Marini, J.-M. Roger, SO-CovSel: A novel method for variable selection in a
794 multiblock framework, Journal of Chemometrics, 34 (2020) e3120.

795 [56] B. Galindo-Prieto, J. Trygg, P. Geladi, A new approach for variable influence on projection (VIP) in
796 O2PLS models, Chemometrics and Intelligent Laboratory Systems, 160 (2017) 110-124.

797 [57] B. Galindo-Prieto, P. Geladi, J. Trygg, Multiblock variable influence on orthogonal projections (MB-
798 VIOP) for enhanced interpretation of total, global, local and unique variations in OnPLS models, arXiv
799 preprint arXiv:2001.06530, (2020).

800 [58] S. Park, E. Ceulemans, K. Van Deun, Sparse common and distinctive covariates regression, Journal
801 of Chemometrics, n/a (2020) e3270.

802 [59] A. Biancolillo, T. Næs, R. Bro, I. Måge, Extension of SO-PLS to multi-way arrays: SO-N-PLS,
803 Chemometrics and Intelligent Laboratory Systems, 164 (2017) 113-126.

804 [60] A.K. Smilde, J.A. Westerhuis, R. Boqué, Multiway multiblock component and covariates regression
805 models, Journal of Chemometrics, 14 (2000) 301-331.

806 [61] P. Mishra, J.M. Roger, D.N. Rutledge, E. Woltering, SPORT pre-processing can improve near-
807 infrared quality prediction models for fresh fruits and agro-materials, Postharvest Biology and
808 Technology, 168 (2020) 111271.

- 809 [62] P. Mishra, F. Marini, A. Biancolillo, J.-M. Roger, Improved prediction of fuel properties with near-
810 infrared spectroscopy using a complementary sequential fusion of scatter correction techniques,
811 *Talanta*, (2020) 121693.
- 812 [63] P. Mishra, A. Nordon, J.-M. Roger, Improved prediction of tablet properties with near-infrared
813 spectroscopy by a fusion of scatter correction techniques, *Journal of Pharmaceutical and Biomedical*
814 *Analysis*, (2020) 113684.
- 815 [64] J.-M. Roger, A. Biancolillo, F. Marini, Sequential preprocessing through ORThogonalization
816 (SPORT) and its application to near infrared spectroscopy, *Chemometrics and Intelligent Laboratory*
817 *Systems*, 199 (2020) 103975.
- 818 [65] P. Mishra, J.M. Roger, F. Marini, A. Biancolillo, D.N. Rutledge, Parallel pre-processing through
819 orthogonalization (PORTO) and its application to near-infrared spectroscopy, *Chemometrics and*
820 *Intelligent Laboratory Systems*, (2020) 104190.
- 821 [66] T. Skotare, D. Nilsson, S. Xiong, P. Geladi, J. Trygg, Joint and Unique Multiblock Analysis for
822 Integration and Calibration Transfer of NIR Instruments, *Analytical Chemistry*, 91 (2019) 3516-3524.
- 823 [67] K. De Roover, E. Ceulemans, M.E. Timmerman, How to perform multiblock component analysis in
824 practice, *Behavior Research Methods*, 44 (2012) 41-56.
- 825 [68] P. Mishra, F. Marini, B. Brouwer, J.M. Roger, A. Biancolillo, E. Woltering, E.H.-v. Echtelt, Sequential
826 fusion of information from two portable spectrometers for improved prediction of moisture and
827 soluble solids content in pear fruit, *Talanta*, 223 (2021) 121733.
- 828 [69] K.B. Walsh, J. Blasco, M. Zude-Sasse, X. Sun, Visible-NIR 'point' spectroscopy in postharvest fruit
829 and vegetable assessment: The science behind three decades of commercial use, *Postharvest Biology*
830 *and Technology*, 168 (2020) 111246.
- 831 [70] P. Mishra, A. Biancolillo, J.M. Roger, F. Marini, D.N. Rutledge, New data preprocessing trends
832 based on ensemble of multiple preprocessing techniques, *TrAC Trends in Analytical Chemistry*, (2020)
833 116045.
- 834 [71] R. Lu, R. Van Beers, W. Saeys, C. Li, H. Cen, Measurement of optical properties of fruits and
835 vegetables: A review, *Postharvest Biology and Technology*, 159 (2020) 111003.
- 836 [72] T. Skotare, R. Sjögren, I. Surowiec, D. Nilsson, J. Trygg, Visualization of descriptive multiblock
837 analysis, *Journal of Chemometrics*, 34 (2020) e3071.

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.