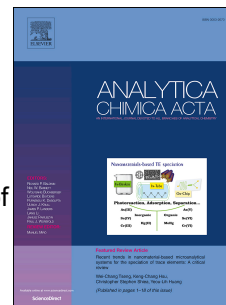# Journal Pre-proof

Multi-level Data Fusion Strategies for Modeling Three-way Electrophoresis Capillary and Fluorescence Arrays Enhancing Geographical and Grape variety Classification of Wines

Rocío Ríos-Reina, Silvana M. Azcarate, José M. Camiña, Héctor C. Goicoechea

Please cite this article as: R. Ríos-Reina, S.M. Azcarate, J.M. Camiña, H.C. Goicoechea, Multi-level Data Fusion Strategies for Modeling Three-way Electrophoresis Capillary and Fluorescence Arrays Enhancing Geographical and Grape variety Classification of Wines, *Analytica Chimica Acta*, https://doi.org/10.1016/j.aca.2020.06.014.

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.
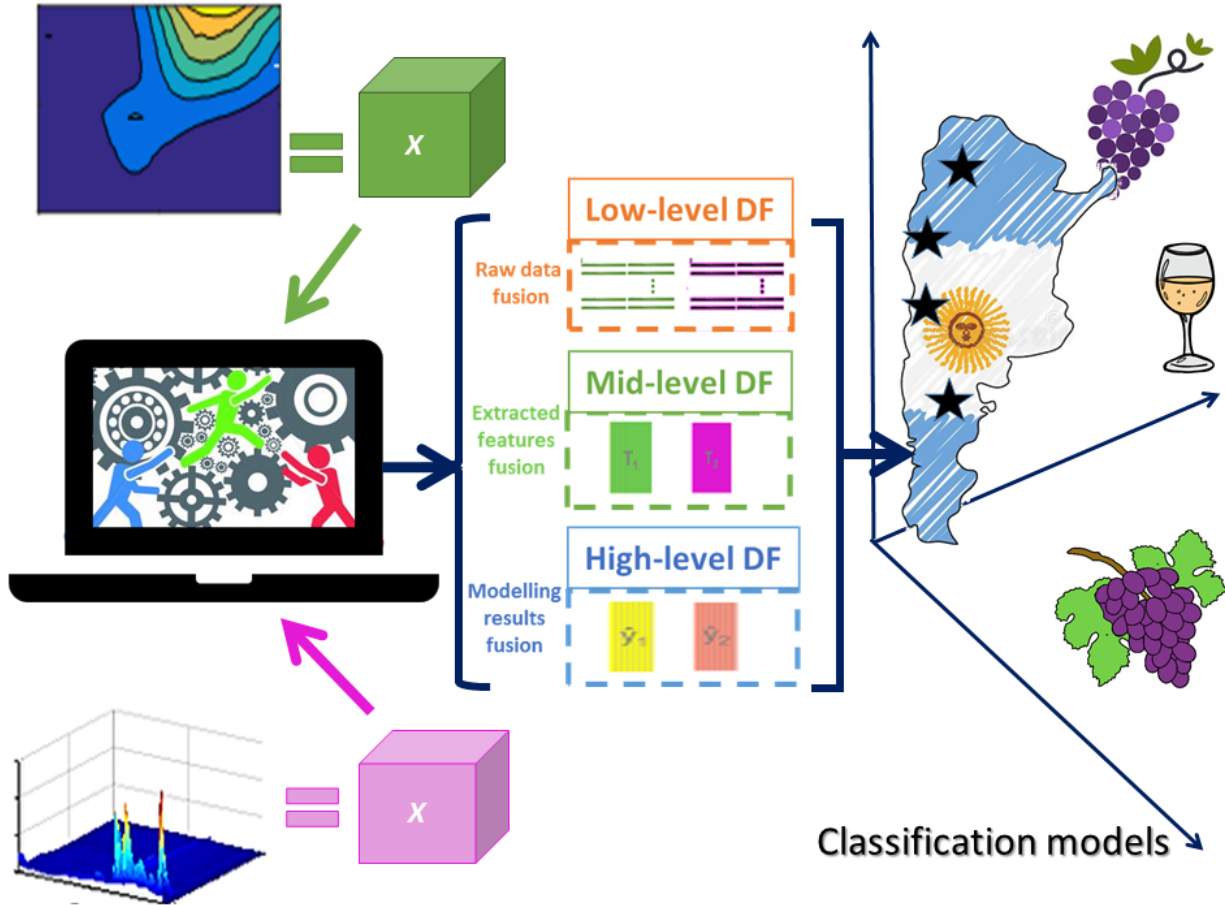
# CRediT author statement

**Silvana M. Azcarate:** Conceptualization, Methodology, Investigation, Writing - Review & Editing.

**Rocío Ríos Reina:** Software, Formal analysis, Visualization.

**José M. Camiña:** Supervision.

**Héctor C. Goicoechea:** Conceptualization, Investigation, Resources, Supervision, Funding acquisition.

# Multi-level Data Fusion Strategies for Modeling Three-way

# Electrophoresis Capillary and Fluorescence Arrays

# Enhancing Geographical and Grape variety Classification of

# Wines

Rocío Ríos-Reina [a], Silvana M. Azcarate [b,*], José M. Camiña [b], Héctor C. Goicoechea [c]

**Affiliations**

[a] Área de Nutrición y Bromatología, Fac. Farmacia, Univ. Sevilla, C/P. García González no. 2, E-41012 Sevilla, Spain.

[b] Facultad de Ciencias Exactas y Naturales, Universidad Nacional de La Pampa-CONICET, and Instituto de Ciencias de la Tierra y Ambientales de La Pampa (INCITAP), Av. Uruguay 151 (6300) Santa Rosa, La Pampa, Argentina.

[c] Laboratorio de Desarrollo Analítico y Quimiometría (LADAQ), Cátedra de Química Analítica I, Facultad de Bioquímica y Ciencias Biológicas, Universidad Nacional del Litoral-CONICET, Ciudad Universitaria, Santa Fe (S3000ZAA), Argentina

*To whom correspondence should be addressed. E-mail silvanaazcarate@gmail.com

18

19 **Abstract**

20 Capillary electrophoresis with diode array detection (CE-DAD) and multidimensional

21 fluorescence spectroscopy (EEM) second-order data were fused and chemometrically

22 processed for geographical and grape variety classification of wines. Multi-levels data

23 fusion strategies on three-way data were evaluated and compared revealing their

24 advantages/disadvantages in the classification context. Straightforward approaches

25 based on a series of data preprocessing and feature extraction steps were developed for

26 each studied level. Partial least square discriminant analysis (PLS-DA) and its multi-

27 way extension (NPLS-DA) were applied to CE-DAD, EEM and fused data matrices

28 structured as two-way and three-way arrays, respectively. Classification results

29 achieved on each model were evaluated through global indices such as average

30 sensitivity non-error rate and average precision. Different degrees of improvement were

31 observed comparing the fused matrix results with those obtained using a single one,

32 clear benefits have been demonstrated when level of data fusion increases, achieving

33 with the high-level strategy the best classification results.

34

37

38

39

40

**1. Introduction**

In multivariate classification setting, second-order data handling is producing a significant impact on the development of analytical strategies, especially for determining characteristic patterns of the analytes of interest in highly complex matrices [1, 2]. In particular, food characterization tasks are in an innovative progressive change going from the development of dedicated methods for quantification of specific compounds to the fingerprint acquisition by rapid, non-destructive and non-selective instrumental techniques [3–5].

Hence, second-order data, especially three-way arrays modeling applications, could supply interesting improvements as regards results attained when extremely complex systems should be classified. Thereby, it has been demonstrated that data analysis can be more effective when modeling second- or higher -data with multiway algorithms compared to unfolding procedures [1, 6]. Otherwise, it has been recently revealed that the ability in terms of discrimination power can be improved by using second-order data arranged in a three-way structure instead of first-order data [2].

Over the years, analytical methods and data analysis tools commonly used in food quality and process control had to be re-evaluated and modified to fit these new tasks [1]. In this progression of gathering more and better information, the multivariate statistical analysis of fused data has become a powerful tool for enhancing the reliability of the results. Being the key point how the information sources can be combined to provide the joint classification prediction of the samples, three levels of data fusion (DF) have been reported [7, 8].

Firstly, low-level DF (1-DF) implies a simple concatenation of the individual matrices to build a single array that is then used for calculating a single model for final classification. In the food authentication and quality control field, it has been the most

2

66  used fusion approach to improve the results since is a common, conceptually simple,

67  first attempt with outputs from different sources providing first-order data [9–11].

68  Nevertheless, out of the food analysis scope, three-way arrays have been concatenated

69  using low-level DF strategy in a recent original report [12].

70      Secondly, mid-level DF (2-DF) first extracts some relevant features from each

71  data source separately and then these outputs are concatenated to build a single array to

72  be then processed by the desired chemometric technique. This approach has been

73  probably the one that has purposed more challenges chiefly for second-order data

74  analysis, where witty strategies for data compression, extraction or reduction have been

75  explored for improving outcomes. Mainly, the attempt has been addressed to evaluate

76  the combination of first- and second order data provided by multiple platforms since

77  data are very different in structure, size or scale [13–18]. Otherwise, the performance of

78  mid-level applications has commonly been compared to low-level fusion as well as to

79  single models [19–22].

80      Lastly, the high-level DF (3-DF) builds separate models for the different blocks,

81  and the individual results are then integrated into a single final response. This strategy

82  has been lesser explored than the two mentioned above. Although several

83  methodologies for final identity declaration by modelling the individual matrices

84  independently have been reported [8, 23], only few of them have been inquired in food

85  classification context. High-level DF has been mainly implemented for the comparison

86  with the other two DF levels [24, 25].

87      The aim of this work was to develop multiple strategies to assess the three DF

88  levels on two second-order arrays, with different data complexity, in order to know the

89  correlation and analogy between both information sources for twofold classification

90  purposes. The focus was put on the development of models able to distinguish among

91  white wines of three different grape varieties with geographical indication (GI) from the

92  four main wine production regions of Argentina [26]. For that, fluorescence excitation–

93  emission matrix spectroscopy (EEM), and capillary electrophoresis with diode array

94  detector (CE-DAD) were applied as non-target analysis in order to acquire a fingerprint

95  to characterize the wines.

96      To our best knowledge, it is the first time that the multi-levels of DF strategies

97  on second-order data are evaluated and compared revealing their

98  advantages/disadvantages in the classification context. Thus, we developed multiple and

99  straightforward approaches based on a series of data preprocessing and feature

100 extraction steps, which constitutes a significant improvement in the DF analysis, and it

101 offers a wide range of possibilities when second-order data of different nature are

102 assessed. Finally, the challenge consisted in finding the optimal combination of data

103 preprocessing, fused data and data modeling that would provide the best results.

104

105 **2. Materials and methods**

106 *2.1. Samples*

107      Thirty-nine samples of commercial white wine from four wine-producing

108 origins, all belonging to provinces of Argentina (Mendoza-M, San Juan-SJ, Salta-S, and

109 Rio Negro-RN) and three different grape varieties (Chardonnay-CH, Sauvignon Blanc-

110 SB and Torrontés-TO), were included in this study: 14 Chardonnay wines (10 from

111 Mendoza and 4 from San Juan), 13 Sauvignon Blanc wines (10 from Mendoza, 1 from

112 San Juan, and 2 from Río Negro) and 12 wines from grapes of the variety Torrontés (4

113 from Mendoza, 1 from San Juan, 5 from Salta, and 2 from Río Negro). Wine samples

114 were selected from the 2011 to 2013 vintages and bought from a local market. The

115   alcoholic content ranged from 12.2 to 13.8% v/v of ethanol. These samples were

116   analyzed in triplicate by the two techniques described below.

117

118   *2.2. Data acquisition*

119   2.2.1. EEM data

120       All spectrofluorimetric measures were acquired according to the method

121   reported by Azcarate *et al.* [26] using a Cary Eclipse Fluorescence Spectrophotometer

122   (Agilent Technologies, Waldbronn, Germany) with a 1×1 cm quartz fluorescence cell,

123   xenon flash lamp. CaryEcplise software package was used to control the instrument,

124   data acquisition and data analysis. Fluorescence excitation spectra were recorded by

125   varying the wavelengths between 245 and 341 nm (increment 5 nm), and by recording

126   the emission spectra from 300 to 500 nm (spaced by 0.5 nm interval). Fluorescence

127   measurements were done in triplicate for each sample.

128

129   2.2.2. CE-DAD data

130       The electrophoretic run conditions are treated in detail in our previous work [27]

131   and here only main analysis steps will be recalled. All electropherograms were acquired

132   on a CE system (Agilent Technologies, Waldbronn, Germany) equipped with a DAD

133   and an uncoated fused silica capillary of 40 cm total length (31.5 cm effective length)

134   and 75 µm inner diameter (MicroSolv Technology Corporation, Eatontown, NJ, USA).

135   Separation was performed by applying a voltage of 24 kV and with a typical current of

136   approximately 80 µA. The hydrodynamic injection was performed in the positive

137   electrode of the capillary by applying a pressure of 40 mbar for 8 s. The cartridge was

138   maintained at 25.0°C. The electropherograms were recorded during 10 min at 0.3 s

139    steps and recording UV spectra between 189 and 401 nm each 2 nm and samples were

140    analyzed by triplicate.

141

142    *2.3. Data analysis*

143    In order to extract and/or merge the information presented in each data set

144    obtained by the two different instrumental analysis of each sample, different

145    chemometric algorithms were employed. As both techniques produced out-puts with the

146    same data structure (i.e. three-way arrays), they were analyzed by similar algorithms in

147    order to decompose and compress the data.

148    The data analysis workflow developed in this study is schematized in **Fig. 1.** In

149    general terms, it includes: 1) building separate classification models on data obtained

150    from the individual analytical techniques by applying 3 different approaches; and 2)

151    building classification multiplatform models by applying different DF strategies: l-DF,

152    2-DF and 3-DF (assessing different approaches). Then, all the classification models

153    obtained were assessed and compared.

154

155    **Insert here Fig. 1**

156

157    2.3.1. Data set and preprocessing

158    In order to validate the classification models, the dataset corresponding to each

159    technique, containing 39 samples, was split into a training set of 24 samples (12 CH, 12

160    SB and 12 TO or 15 M, 4 SJ, 2 RN and 2 S) and a test set of 15 samples (6 CH, 5 SB

161    and 4 TO or 9 M, 2 SJ, 2 S, 2 RN) by using the Duplex algorithm [28], keeping the

162    triplicates in the same set (i.e. the training set contained 72 analysis and the test set 45

163    analysis). The split between training and test sets was done by keeping the ratio of

6

164 samples of each class like in the original set, balancing the representation of each

165 category and keeping the replicates together. Moreover, after checking by exploratory

166 analysis that both sets spanned the whole variability domain, the same split was

167 maintained for all the data sets (the individual and the fused data sets).

168 To find the optimal classification results for each class studied (variety or

169 origin), different preprocessing options were considered in each model: both mean

170 centering and autoscaling were used depending on the nature of the data, as well as none

171 preprocessing.

172

173 2.3.2. Decomposition and compression methods

174 As can be observed in **Fig. 1**, different decomposition and compression methods

175 (i.e. exploratory and reduction data analysis) were applied. Then, the features obtained

176 were used for the DF models. Thus, on the one hand, the original EEM and CE-DAD

177 three-way arrays ($117 \times 49 \times 41$ and $117 \times 676 \times 107$, respectively, considering the samples

178 by triplicate) were unfolded in a multiset structure via row-wise augmentation and then

179 these new matrices (a matrix of $117 \times 2009$ for EEM unfolded data and of $117 \times 72332$ for

180 CE-DAD unfolded data), as well as the fusion of both, were used for the classification

181 or compressed by principal component analysis (PCA).

182 On the other hand, the EEM three-way array was decomposed by parallel factor

183 analysis (PARAFAC) [29] into trilinear components, related to the main fluorophores

184 present in the samples, whose scores (first mode loadings) were used as features for the

185 classification, or to build a fused dataset previous to the classification process**.** A three-

186 factor model, constrained with non-negativity in all modes, was obtained as the

187 optimum model according to the CORe CONsistency DIAgnostic test (COR-CONDIA)

188 [18, 30], the explained variance, the visual inspection of the profiles and residuals [26].

189    Therefore, the three-way array matrix was decomposed by PARAFAC to three new

190    matrices containing the PARAFAC scores of the three fluorophores, as well as their

191    excitation and emission loadings.

192        Finally, CE-DAD array decomposition was carried out by Tucker3 and the

193    resulted **A** matrix (Tucker3 output with the concentrations) was directly used, or fused

194    in a new data set, for the classification purpose. In this study, Tucker3 was selected due

195    to the high complexity of the CE-DAD data, which require different number of factors

196    in each mode. The number of factors selected was 18, 18 and 6 for each recorder mode,

197    obtaining a model with a 95% of total explained variance. Non-negativity was imposed

198    as unique constraint in all modes in agreement with a previous work [27] and three

199    matrices **A**, **B** and **C** were obtained containing the concentration, electropherogram and

200    spectra profiles, respectively, together with a *G* core ($18\times18\times6$) corresponding to the

201    magnitude of the interaction among factors in different modes.

202

203    2.3.3. Classification methods

204        In this work, two classification techniques derived from the regression algorithm

205    partial least squares (PLS) were used: the PLS-DA (DA for discriminant analysis) [31]

206    for first-order data and its multi-way or multilinear extension (NPLS-DA) [32] for

207    second-order data (three-way arrays). In order to select the proper number of latent

208    variables (LVs), i.e. the dimensionality of the model, the minimum classification error

209    rate in cross-validation (venetian blind) was considered. In discriminant analysis, the

210    dependent variable, **Y**, holds the class information (as many *y*-variables as number of

211    classes). The raw predictions from a PLS-DA model is a value of nominally zero or one.

212    A value closer to zero indicates the new sample is not in the modeled class; a value of

213    one indicates a sample is in the modeled class [31].

214

215 *2.4. Individual and data fusion strategies*

216       In this study, four modeling strategies were tested in order to obtain the best

217 classification of the wine samples according to their origin and/or grape variety:

218 classification models of individual techniques and classification models by low-, mid-

219 and high-level data fusion approaches (**Fig. 1**).

220       Before the data fusion, three classification models were obtained with the

221 individual data matrices EEM and CE-DAD: a NPLS-DA model obtained by each

222 original three-way array, and two PLS-DA models, one with the unfolded matrices and

223 other with the decomposed matrices by PARAFAC or Tucker3, respectively. The

224 strategy followed is described at the top of **Fig. 1.** Each matrix was split into a

225 training/test set (72/45) before building the classification models.

226

227 2.4.1. Low-level data fusion (1-DF)

228       In the 1-DF approach, the data matrices are directly concatenated to provide

229 sample classification [8]. In this study, the unfolded data from CE-DAD and EEM

230 matrices were concatenated before any model calculation. These single blocks were

231 joined in a single matrix providing an overall data set with 74341 variables (2009

232 variables from the unfolded EEM matrix plus 72332 variables from the unfolded CE-

233 DAD matrix). After that, the data fused matrix was split into training and test sets.

234 Then, two different 1-DF options were tested (**Fig.1.**).

235       Despite there are many options to be carried out in this 1-DF approach (i.e.

236 applying PCA on the concatenated matrix and then LDA on the scores, or the direct

237 application of other classification algorithms), the selected option was to develop a

238 PLS-DA model, with 6 LVs, and built directly with the concatenated data matrix, after

239   mean centering data preprocessing, by means of the previously described validation

240   protocol (named as Low- level DF: opt1 in **Fig.1**) in order to perform the same

241   classification method in all the strategies of the study.

242

243   2.4.2. Mid-level data fusion (2-DF)

244       In general, in 2-DF strategies, the analytical data are merged at the features level

245   [25]**.** This means that relevant features are independently extracted from each analytical

246   data matrix, which are then concatenated into a single global matrix that is used as input

247   to perform a classification model[8]. In comparison with 1-DF strategies, this method

248   allows to guarantee a more balanced representation of each source of information, in the

249   case of each analytical data matrix has a huge difference in the number of variables

250   [25]**.** However, in this 2-DF the main issues to control are the features to retain from

251   each model and the method to extract them as well as the preprocessing method to

252   adopt. In this study, two different 2-DF strategies were tested, differing from the feature

253   extraction method used.

254       In the first 2-DF option (named as Mid-level DF: opt-1 in **Fig. 1**)**,** the relevant

255   feature extraction was performed by the development of a PCA model for each data

256   block. The number of principal components (PCs) chosen for each PCA model was

257   again selected in order to give more than 90% of cumulative variance in both blocks.

258   Thus, 7 and 4 PCs were selected for the unfolded both CE-DAD and EEM matrices

259   (previously pre-processed by mean-centering), respectively. Then, the PCA-scores

260   associated to the first 7 and 4 PCs for each data block (Str) were considered as extracted

261   features and were then fused in a new matrix (with 11 variables). This fused matrix was

262   pre-processed by auto-scaling or none-preprocessed and then modeled by means of

10

263   PLS-DA, obtaining PLS-DA models with 4 and 5 LVs according to the preprocessing

264   method applied.

265         In the second 2-DF option tested in this study (named as Mid-level DF: opt-2 in

266   **Fig. 1**)**,** the relevant features of each data block were obtained by the development of a

267   PARAFAC and a Tucker3 model for EEM and CE-DAD matrices, respectively. These

268   models are similar to those described above for the individual data modeling (Section

269   2.3.2)**.** Then, the scores associated to the 3 PARAFAC factors extracted from EEM's

270   array were concatenated with the 18 Tucker3 scores extracted from CE-DAD's array,

271   forming a fused matrix with 21 variables, which was autoscaled and used for building of

272   a PLS-DA model of 6 LVs.

273         In all these strategies, PLS-DA models were applied to the fused score-matrices

274   starting from the training-test set split procedure.

275

276   2.4.3. High-level data fusion (3-DF)

277         In 3-DF strategies, the classification of the samples is performed independently

278   on each analytical data block, and then the predictions provided by the models

279   calibrated on the single blocks are combined together [8]**.** In other words, the

280   information in the different data matrices is joined at the level of the prediction obtained

281   by each individual model into a unique solution [33]**.**

282         In this study, a PLS-DA and N-PLS-DA models were first independently fit for

283   EEM and CE-DAD data matrices (data unfolded and decomposed, and original three-

284   way arrays, respectively) and then the decisions/prediction obtained by each single-

285   block model were fused by two different 3-DF strategies proposed in the literature [34]**:**

286   Majority voting and Bayesian consensus with discrete probability distributions.

287    On the one hand, Majority voting was carried out by directly merging the

288    predictions of the single PLS-DA or NPLS-DA models (**Fig. 1**). This 3-DF method is

289    based on a democratic (weighted) process that combines the predictions provided by the

290    individual classification models and classifies a sample into a class according to the

291    most frequent class assignment. Within this method, there are three criteria deriving by

292    applying specific limits or thresholds. The 'loose' criterion is the simplest and most

293    intuitive, in which a sample is assigned based on the most frequent class assignment,

294    and a sample is not classified in case of ties (frequency of assignments to a class >50%).

295    The "intermediate" and "strict" majority voting criteria classify a sample if the

296    agreement of predictions is higher than or equal to 75% and 100% (full prediction

297    agreement of all the considered models), respectively [34]. In this study, as only two

298    analytical methods are fused, only a sample is classified into a class when both

299    techniques classify it into the same class, so the criterion used was the 'strict' (100% of

300    frequency assignments).

301    On the other hand, from the confusion matrix, the Bayesian consensus estimates

302    the probability that a sample belongs to a specific class on the basis of each analytical

303    data block and then combines these probabilities into a joint probability used for the

304    final assignation [34]. As Bayesian results are affected by the model sequence followed

305    in the iteration process, all combinations of analytical sources were considered in our

306    study (i.e. both blocks were selected as the initial block), and according to the

307    classification results, the best order was to start with the EEM's dataset and then with

308    the CE-DAD's dataset.

309    In a first step, the prior probability has been estimated as equal probability. Considering

310    three classes according to variety and four classes according to origin, the used prior

311    probabilities were 0.33 and 0.25, respectively. Then, likelihood conditional probabilities

312     were estimated from the confusion matrix of each classification model being calculated

313     by dividing the number of classes correct and incorrect predicted by the total of samples

314     of each class. Then, once the posterior probabilities have been calculated for the first

315     analytical block, the fusion approach proceeds iteratively, that is, the posterior

316     probabilities of the first model were used as new prior probabilities in the second model.

317     For that, the class predicted by the first-block classification model (with EEM array,

318     unfolded or decomposed) was initially considered. Then, the posterior probabilities of

319     the first model were used as new prior probabilities in the second model, where the

320     class predicted by the second block model (i.e. those obtained from the CE-DAD array,

321     unfolded or decomposed) was the new evidence.

322        Finally, this last probability obtained (i.e. the consensus probability derived from

323     the combination of the information of both data blocks) was the one used to predict the

324     class according to the maximum posterior probability obtained. Hence, this last

325     posterior probability was used to accept or reject the predicted class depending on a

326     predefined probability threshold, that in this study was defined as >50%. The

327     corresponding equations and further details of this method can be found in the literature

328     [34, 35].

329

330     *2.4. Evaluating models*

331        The classification models were internally validated by using venetian blind

332     cross-validation (CV) and the final models' performance was confirmed by a test set

333     validation (TV). For that purpose, as it was mentioned above, the dataset was split into a

334     training and test set, with 61.5% and 38.5% of the samples, respectively, by keeping the

335     ratio of samples of each class like in the original set and the triplicates together.

336     The quality of the models was assessed from the classification and prediction

337  abilities. The optimal conditions were decided by means of primary measures related to

338  single classes as sensitivity (Sens.), specificity (Spec.) and precision (Prec.) of the

339  calibration and prediction, which were calculated on each class separately encoding

340  different aspects of the classification [34]. This information can be found in the

341  **Supporting Information** as **Table I** and **Table II** for grape variety and geographical

342  classification results, respectively. Additionally, to provide an overall evaluation of the

343  classification quality, the global indices derived from primary class measures such as

344  average sensitivity (non-error rate –NER-) and average precision (PREC.) were also

345  calculated according the recommendation of Ballabio *et al* [34].

346

347  *2.5. Software*

348     Spectra preprocessing, and low-level, mid-level and high-level DF strategies

349  were carried out by means of hand-made routines written in Microsoft Excel v. 2016

350  (Microsoft Corporation, USA) and Matlab R2014b (The Mathworks, Natick, MA,

351  USA). Decomposition and compression methods (PARAFAC, Tucker3, PCA) and PLS-

352  DA classification models were calculated with the PLS_Toolbox 7.9.5 (Eigenvector

353  Research Inc., Wenatchee, WA) working under MATLAB environment.

354

355  **3. Results and discussion**

356  *3.1. Data visualization*

357     The fluorescence and CE-DAD landscapes of several samples belonging to the

358  three grape varieties of the four geographical origins are shown in **Fig. 2**. On the one

359  hand, it could be observed that the shape of the EEM spectra varied within the same

360  origin among varieties, as well as within the same grape variety among origins. Thus,

14

361　the visual assessment of all the fluorescence features of the grape varieties pointed out a

362　general trend for the spectral maxima to be shifted towards 450 and 350 nm of em/ex.

363　Furthermore, similar fluorescence trend was observed for the different origins but

364　maintaining the characteristic shape of its variety.

365　　　　On the other hand, as can be seen in the CE-DAD landscapes (**Fig. 2**), they

366　showed many overlapping peaks corresponding to the complex mixture of chemical

367　compounds that are present in the white wines. Moreover, as was shown in a previous

368　report [27] a remarkable peak misalignment and shape deformation in electrophoretic

369　mode was produced. It could be also observed some differences between geographical

370　origins and grape varieties. Thus, all the samples showed marked peaks around 3 and 5

371　min but with strong differences among varieties and origins. From these observations,

372　all these differences could make possible the classification of the samples according to

373　origins and varieties.

374

375　　　　　　　　　　　　　　　**Insert here Fig. 2**

376

377　*3.2. Individual models*

378　3.2.1. General considerations

379　　　　In the first stage, excitation-emission matrices (EEM) and capillary

380　electrophoresis (CE-DAD) data were treated separately to build the classification

381　models. The data matrices were organized and analyzed as two- and three-dimensional

382　arrays. Thus, three different classification approaches considering the data structures

383　and modeling were performed: (a) three-way data by PLS-DA using factors obtained

384　from a resolution method (PARAFAC or Tucker3); (b) three-way data by N-PLS-DA;

385　and (c) full unfolded data using PLS-DA (schematized in **Fig. 1**).

15

386    It should be noted that both datasets differ in the complexity of the data

387    structure. EEMs represents a well-known illustration of bilinear data fulfilling with the

388    so-called low rank trilinearity condition, which can be decomposed into the excitation

389    and the emission spectra for a given fluorescent component [2]. In return, the three-

390    dimensional array built from CE-DAD data is non-trilinear. Moreover, it presents

391    certain drawbacks like remarkable peak misalignment and shape deformation in a data

392    mode associated with deficiency rank in the other one.  In those cases, in which multi-

393    way data are involved for classification issues, the choice of the appropriate multi-way

394    approach will be decisive in the validity of the solution found.

395    For building models, the latent variables were selected considering the lowest

396    CV classification error rate (data not shown). The best preprocess and region

397    (variables), together with the optimal configuration of each model, such as the number

398    of latent variables retained, were selected as those leading to the lowest inaccuracy and

399    highest sensitivity and specificity obtained with the prediction set.

400

401    3.2.2. Classification models for EEM data

402    For the approach (a) of Fig. 1, EEM dataset was arranged in a three-way data

403    array (72 training samples, 41 emission wavelengths and 49 excitation wavelengths)

404    and then, it was analyzed by means of PARAFAC. A three-factor model was chosen

405    representing the best compromise between explained variance (99.5%) and core

406    consistency (71%). The obtained model presented results that were in good agreement

407    with works presented in literature [26] where the loadings for second (emission) and,

408    third (excitation) modes and PCA scores have already been reported.

409    On the other hand, for the approach (c) of Fig 1, the EEM data array ($72 \times 49 \times$

410    41) was unfolded into a two-dimensional array ($72 \times 2009$). To check the repeatability

411     of the measurements, detect outliers and recognize patterns in the samples' distribution,

412     PCA analyses on PARAFAC factors and unfolded matrix were performed (**Fig. IA.** and

413     **IB. SM**). By analyzing the scores plots, it could not be observed a clear differentiation

414     of the samples by means of geographical origin on both two- and three-way data

415     structures, showing in both cases a similar overlapping, mainly between samples of M

416     and SJ. The same situation was observed when the differentiation among samples was

417     assessed by means of the three grape varieties, as it was shown in a previous report [26].

418             A PLS-DA model was performed on the PARAFAC factors and the unfolded

419     matrix, which were prior preprocessed by autoscaling and mean centering, respectively.

420     On the other side, a classification model based on NPLS-DA on the three-way data

421     matrix was built (approach (b) of Fig 1). The obtained classification results of these data

422     sets are reported in **Table 1** and **Table 2,** when grape variety and geographical origin

423     were used as classifier as well as in **Fig. II** and **III. SM**, respectively. Then separate

424     models were evaluated comparing the number of latent variables retained and the

425     indices derived from confusion matrix (Sens., Spec., Prec.; and PREC. and NER).

426             All the built models for sample classification according to grape variety showed

427     a similar performance that the obtained in a previous work [26]. However, in the present

428     study, the best individual model for grape variety classification seemed to be the NPLS-

429     DA model reaching the highest NER and PREC values in prediction stage being 81.1%

430     and 82.1%, respectively (**Table 1**).

431             Furthermore, suitable classification results were attained according to

432     geographical origin by the three built models. In this case, the model acquired from

433     PLSDA on the PARAFAC factors was the optimal reaching the highest rate of well-

434     classified samples in the prediction set and displaying the highest NER values (**Table**

435     **2**).

17

436

437                                **Insert here Table 1**

438                                **Insert here Table 2**

439

440     3.2.3. Classification models for CE-DAD data

441          For the (a) and (b) approaches (Fig 1), CE-DAD dataset was arranged in a cube

442     structure (72 training samples, 676 times and 107 wavelengths) and it was unfolded in a

443     matrix of size (72 × 72332) for the strategy (c). In order to avoid drawbacks, three-way

444     array was modeled by means of Tucker3 that allows using a different number of factors

445     in each mode. Thus, the number of eigenvalues explaining 95% of the variance of the

446     data were 18, 18 and 6, for modes 1, 2 and 3, respectively. After that, PCA analysis on

447     unfolded matrix and Tucker3 factors were performed. The scatter plots of these PCA

448     analyses are shown in **Fig. IC. SM** and **Fig. ID SM**, respectively. It can be seen that the

449     CE-DAD data showed higher variability than EEM data when the reproducibility was

450     assessed.

451          In the same way, PLS-DA was performed on the Tucker3 factors and the

452     unfolded matrix, applying autoscaling and mean centering as preprocessing,

453     respectively. On the other side, a classification model based on NPLS-DA on the three-

454     way data array was built (approach b). It is relevant to highlight that both (b) and (c)

455     approaches were able to deal rank deficiency [14].

456          The three models for grape variety classification showed similar performances;

457     however, for prediction set, the NPLS-DA model attained higher indices of 66.1 and

458     70.8 for NER and PREC, respectively (**Table 1**). Concerning geographical origin

459     classification, despite the obtained models with the CE-DAD data achieved promising

460     results in the calibration stages, they were not able to predict the samples correctly.

461 Thus, they showed in all cases NER values lower than 56.9 % (**Tabla 1**). These

462 classification results could be also observed by looking the scores and loadings plots of

463 the PLS-DA models reported in **Fig. II** and **III. SM.**

464 Thus, the application of fusion approaches was expected to increase the overall

465 classification ability according to variety and origin classification, by integrating these

466 different behaviors of single analytical sources.

467

468 *3.3. Data fusion models*

469 With the goal of improving the classification results according to geographical

470 origin and grape variety, different strategies were assessed in the three data fusion

471 levels.  Thus, in the case of having second-order data, several are the strategies that

472 could be adopted.

473

474 3.3.1. Low-level data fusion

475 PLS-DA was carried out in the first 1-DF option, directly on the concatenated

476 unfolded data matrix (**Fig. 1**). The validation results on the test samples are reported in

477 **Table 1 and Table 2.**

478 The models based on 1-DF achieved similar classification performances than

479 models from individual blocks when grape variety classification was evaluated. Despite

480 the two 1-DF options were able to correctly discriminate the same number of samples,

481 they did not improve the prediction results respect to NPLS-DA model from EEM data,

482 which obtained the best performance of all individual models. However, the option 1

483 reached better results in the calibration stage.

484 On the other side, the results obtained for geographical origin classification

485 showed relevant improvements of both 1-DF approaches in comparison with the

486    individual models, inasmuch as a total of 39 samples pf the training set were correctly

487    predicted. In this case, this strategy seemed to provide significantly better classification

488    results reaching 94.4% and 87.5% for NER and PREC, respectively. It is important to

489    highlight than these results exhibited an increment of more than 5% for NER and a great

490    improvement in PREC of more than 16%, displaying the sterling model ability to avoid

491    wrong predictions in the classes.

492         As a second option, a PCA can be also applied as a necessary step to compress

493    the information when algorithms as LDA are applied. Thus, it is important to be careful

494    inasmuch as, as in this case, the concatenated data from 1-DF consists on an extremely

495    large matrix where the number of irrelevant variables becomes larger than the really

496    meaningful ones and therefore, the selection of the more relevant variables could result

497    difficult [2, 40-41].

498         The scores plot of the best 1-DF models for grape variety and geographical

499    origin classifications are shown in **Fig. 3A** and **3C**, respectively. By comparing these

500    models to those for PLS-DA, obtained for the individual data matrices (**Fig. II** and **III**

501    **SM**), the improvement in the separation of classes by the DF models is clearer.

502

503                         **Insert here Fig. 3**

504

505         Data fusion showing better discrimination ability than individual spectroscopies

506    have been also reported for fist-order data [36]. Indeed, most of the researches founded

507    in the literature carried out 1-DF on first-order data [10, 37–39] due to the ease of

508    performance together with the satisfactory results. Basically, 1-DF involves the

509    straightforward concatenation by combining variables of the data blocks. Thus, direct

510    first-order data concatenation is easier than second order-data concatenation since, in

511 these last, a prior step, e.g. data unfolding, could be needed when data arrays differ in

512 structure and complexity. In this case, fewer studies could be found in the literature

513 [18].

514

515 3.3.2. Mid-level data fusion

516      As described above, two different approaches for 2-DF were evaluated. Option 1

517 was based on a first extraction of the relevant features by the development of a PCA

518 model for each unfolded data block, and then, the fusion of the PCA-scores matrices

519 obtained, being this matrix used in the development of the PLS-DA classification

520 models. Within this option, autoscaling or none-preprocessing were tested. The option 2

521 was based on the feature extraction by PARAFAC and Tucker3 of the EEM and CE-

522 DAD matrices, respectively, and then the fusion of the scores associated to the 3-

523 PARAFAC factors from EEM array and the 18-Tucker3 scores from CE-DAD array,

524 being the matrix used for building the PLS-DA classification models.

525      On the one hand, with regards to the two 2-DF options, by observing both grape

526 variety (**Table 1**) and geographical origin (**Table 2**) classification results for calibration

527 steps, the option 2 showed better classification results, calibrating correctly almost the

528 total of samples of the training set. On the other hand, by assessing the prediction rate of

529 grape variety classification (**Table 1**), it was again difficult to select the best option due

530 to once again, the two options were able to discriminate correctly the same amount of

531 samples for prediction set, but none of them was able to improve the performance

532 respect to the individual models.

533      Otherwise, the geographical origin classification rate of the option 2 was better

534 than the option 1 achieving 91.7 and 95.5 for NER and 79.2 and 87.5 for PREC,

535    respectively (**Table 2**). These results can be also observable by looking the scores plots

536    obtained by these PLS-DA models (bottom of **Fig. 3B and Fig. 3D**).

537         By making a general comparison of these 2-DF results with the individual

538    classification models, they also showed relevant improvements in the prediction rates in

539    comparison to the individual classification results (mainly for geographical origin

540    classification). However, similar classification results were achieved in comparison with

541    the 1-DF approaches, for both grape variety and geographical origin classification.

542    Nevertheless, better results in classification have been reported when 2-DF approach

543    was compared with the analysis of individual datasets or with 1-DF [14, 18, 22, 25, 42].

544         Despite both low- and mid- levels improved the individual classification results,

545    being similar between them, there are different advantages and disadvantages that could

546    be considered. For both cases, the data block obtained is then processed by the desired

547    chemometric technique. However, on the one hand, 1-DF only implies that the matrices

548    describing the individual blocks, after proper preprocessing, are concatenated to build a

549    single array, being easier even more if the data has a first-order structure. However, a

550    disadvantage of 1-DF is that typically data sets are obtained in which the number of

551    observations is much smaller than the number of variables, which prevents to apply

552    many multivariate data analysis techniques that are not directly applicable. The most

553    popular way of trying to solve the problem of many variables is to reduce the

554    dimensionality of each data matrix separately, before attempting to link them by means

555    of DF and this is how 2-DF works.

556         Otherwise, there are multiple possibilities that can be applied to carry out a 2-DF

557    strategy. The most remarkable techniques reported in the bibliography for multiway

558    data sets have been sequential and orthogonalized partial least squares (SO-PLS) [44]

559    and coupled matrix and tensor factorization (CMTF) [45]. Other approaches based on

560 multiblock analysis less used but also suitable in data fusion context, is the Common

561 Components and Specific Weights Analysis (CCSWA, also so-called ComDim) [18,

562 46].

563

564 3.3.3. High-level data fusion

565      In this level of DF, two different approaches were also developed and assessed:

566 Majority Voting and Bayesian consensus methods. These approaches were implemented

567 by using the classification results of the 3 individual models (NPLS-DA, and PLS-DA

568 from the unfolded matrix and the extracted features matrix) of both analytical methods.

569 The classification results obtained by both approaches are shown in **Table 1** and **Table**

570 **2** for grape variety and geographical origin classification, respectively.

571      Concerning the prediction results of Bayesian consensus DF, the model

572 performed using the outputs of the individual NPLS-DA models for grape variety

573 classification provided the best results. Although this model could only match the

574 amount of samples correctly classified in the prediction set with the NPLS-DA

575 individual model from EEM data (NER = 81.1), it was able to improve the calibration

576 stage getting over in more than 15% in both indices NER and PREC.

577      Otherwise, for geographical origin classification, the Bayesian consensus model

578 developed from the PLS-DA results obtained with PARAFAC and Tucker3 scores

579 provided the best predictions. These results agreed with the classification results

580 obtained by the individual classification models discussed in section 3.2.2 and 3.2.3.

581 Therefore, in general terms, discrimination performances based on Bayesian high-level

582 fusion approaches resulted to be better than those obtained on single analytical sources,

583 as occurs in another work founded in the literature [25]. Hence, this improvement was

584 better observable in the case of geographical origin classification, for which the

23

585   Bayesian consensus fusion obtained the best prediction results achieving 97.2% and

586   91.7% for NER and PREC. This improvement could indicate that the reliability and

587   confidence of the final outcome are increased by the integration of heterogeneous

588   predictions. Moreover, classification performances have been previously reported by

589   means of Bayesian consensus 3- DF fusion achieving slightly better than those obtained

590   in the mid- level approach [25].

591       However, the classification results obtained by the Majority Voting approach

592   were worse than the results obtained by the PLS-DA models made with the individual

593   data blocks. The main reason of that could be the fact that, considering only two

594   analytical techniques, the criteria applied was the "strict", which means that only the

595   samples that were perfectly classified in both techniques could be classified into a

596   specific class. Hence, in the present study, there were many cases in which one

597   technique classified a sample to one class and the other technique classified the same

598   sample as another class, making that the final decision was "not-classified". For this

599   reason, Majority Voting as 3-DF strategy should be preferably applied when three or

600   more techniques are studied simultaneously.

601       In comparison to the other two DF strategies (i.e. 1-DF and 2-DF), 3-DF has not

602   the problem of needing to adjust an adequate scaling due to each model is fitted

603   independently with its best scaling. However, a disadvantage of 3-DF is that the order

604   of combining the obtained predictions affects the final decisions.

605       As a final remark, it is important to highlight the evident improvement of the

606   classification results as level of data fusion increases. In fact, the improvement of the

607   DF prediction models can be linked in the level advance. Despite the combination of

608   multiple analytical sources increases the complexity of data treatment, this is

609   compensated by significantly better classification ability.

610

## 4. Conclusions

611

612     The proposed multi-level fusion strategies provide a useful and reliable way of

613     improving the analytical quality of the results in second-order data for classification

614     outcomes. The benefit of fusion is highlighted in prediction stage when samples cannot

615     be classified from individual sources. In particular, these advantages were more evident

616     when geographical origin classification was assessed, especially taking in account the

617     complexity of the system presenting unbalanced classes. In addition, multi-level data

618     fusion from multi-via modeling accomplished the best classification models. Thus, it is

619     noteworthy that the benefits of data fusion at different levels are added to the second-

620     order data advantage, furnishing a synergistic effect on the classification results.

621     Although both techniques provided good classification results separately, data

622     fusion approaches improved the classification results and provided a larger description

623     of the sample. Hence, the statistic mathematical integration of the information from the

624     different analytical sources can be helpful because it leads to the minimization of the

625     overall uncertainty due to a compensation effect among the single experimental

626     uncertainties. This finally translates into increased reliability of the outcome, and,

627     therefore, it can be concluded that high-level strategies are suitable approaches to obtain

628     greater confidence on the combined (fused) analytical predictions. Notwithstanding the

629     practical application seems to be more cumbersome insomuch as first, the independent

630     models for each platform must be fit. However, model outputs combination is easy to

631     implement and analyse, and it does not require higher effort to be performed.

632

633     **Acknowledgements**

638

639     **References**

640     [1]    E. Salvatore, M. Bevilacqua, R. Bro, F. Marini, M. Cocchi, Classification

641            Methods of Multiway Arrays as a Basic Tool for Food PDO Authentication,

642            Compr. Anal. Chem. 60  (2013) 339-382.

643     [2]    S.M. Azcarate, Araújo, A. Muñoz de la Peña, H.C. Goicoechea, Modeling

644            Second-Order Data for Classification Issues : Data Characteristics, Algorithms,

645            Processing Procedures and Applications, TRAC-Trend Anal. Chem. 107 (2018)

646            151-168.

647     [3]    I. Reinholds, V. Bartkevics, I.C.J. Silvis, S.M. Ruth, S. Esslinger, Analytical

648            Techniques Combined with Chemometrics for Authentication and Determination

649            of Contaminants in Condiments : A Review, J. Food Compos. Anal. 44 (2015)

650            56-72.

651     [4]    M.P. Callao, I. Ruisánchez, An Overview of Multivariate Qualitative Methods

652            for Food Fraud Detection, Food Control 86 (2018) 283–293.

653     [5]    A.M. Gómez-Caravaca,  R.M. Maggio, L.Cerretani, Chemometric Applications

654            to Assess Quality and Critical Parameters of Virgin and Extra-Virgin Olive Oil .

655            A Review, Anal. Chim. Acta  913 (2016) 1-21.

656     [6]    J.A. Arancibia, C.E. Boschetti, A.C. Olivieri, G.M. Escandar, Screening of Oil

657            Samples    on    the    Basis    of    Excitation-Emission    Room-Temperature

658        Phosphorescence Data and Multiway Chemometric Techniques . Introducing the

659        Second-Order Advantage in a Classification Study, Anal. Chem. 80 (2008) 2789–

660        2798.

661    [7]    A. Biancolillo, R. Boqué, M. Cocchi, F. Marini, Data Fusion Strategies in Food

662        Analysis. Data Handl. Sci. Techn. 31 (2019) 271-310.

663    [8]    E. Borràs, J. Ferré, R. Boqué, M. Mestres, L. Aceña, O. Busto, Data Fusion

664        Methodologies for Food and Beverage Authentication and Quality Assessment -

665        A Review, Anal. Chim. Acta 891 (2015) 1–14.

666    [9]    M.S. Godinho, M.R. Blanco, F.F. Gambarra Neto, L.M. Lião, M.M. Sena, R.

667        Tauler, A.E. De Oliveira, Evaluation of Transformer Insulating Oil Quality Using

668        NIR, Fluorescence, and NMR Spectroscopic Data Fusion, Talanta 129 (2014)

669        143–149.

670   [10]    S. Schwolow, N. Gerhardt, S. Rohn, P. Weller, Data Fusion of GC-IMS Data and

671        FT-MIR Spectra for the Authentication of Olive Oils and Honeys—Is It Worth to

672        Go the Extra Mile?, Anal. Bioanal. Chem. 411 (2019) 6005–6019.

673   [11]    C. Pizarro, S. Rodríguez-Tecedor, N. Pérez-Del-Notario, I. Esteban-Díez, J.M.

674        González-Sáiz, Classification of Spanish Extra Virgin Olive Oils by Data Fusion

675        of Visible Spectroscopic Fingerprints and Chemical Descriptors, Food Chem.

676        138 (2013) 915–922.

677   [12]    J. Guillemant, A. Berlioz-Barbier, F. Albrieux, L.P. de Oliveira, M. Lacoue-

678        Nègre, J.F. Joly, L. Duponchel, Low-Level Fusion of FT-ICR MS Data Sets for

679        the Characterization of Nitrogen and Sulfur Compounds in Vacuum Gas Oils

680        Low-Level Fusion of FT-ICR MS Data Sets for the Characterization of Nitrogen

681        and Sulfur Compounds in Vacuum Gas Oils, Anal. Chem. 92 (2020) 2815-2823.

682    [13]    L. Mandrile, L. Barbosa-Pereira, K.M. Sorensen, A.M. Giovannozzi, G. Zeppa,

683            S.B. Engelsen, A.M. Rossi, Authentication of Cocoa Bean Shells by Near- and

684            Mid-Infrared Spectroscopy and Inductively Coupled Plasma-Optical Emission

685            Spectroscopy, Food Chem. 292 (2019) 47–57.

686    [14]    M. Silvestri, A. Elia, D. Bertelli, E. Salvatore, C. Durante, M. Li Vigni, A.

687            Marchetti, M. Cocchi, A Mid Level Data Fusion Strategy for the Varietal

688            Classification of Lambrusco PDO Wines, Chemom. Intell. Lab. Syst. 137 (2014)

689            181–189.

690    [15]    A. Biancolillo, R. Bucci, A.L. Magrì, A.D. Magrì, F. Marini, Data-Fusion for

691            Multiplatform Characterization of an Italian Craft Beer Aimed at Its

692            Authentication, Anal. Chim. Acta 820 (2014) 23–31.

693    [16]    H.M. Santos, J.P. Coutinho, F.A.C. Amorim, I.P. Lôbo, L.S. Moreira, M.M.

694            Nascimento, R.M. de Jesus, Microwave-Assisted Digestion Using Diluted $HNO_3$

695            and $H_2O_2$ for Macro and Microelements Determination in Guarana Samples by

696            ICP OES, Food Chem. 27 (2019) 159–165.

697    [17]    C.R. Carneiro, C.S. Silva, M.A. De Carvalho, M.F. Pimentel, M. Talhavini, I.T.

698            Weber, Identification of Luminescent Markers for Gunshot Residues:

699            Fluorescence, Raman Spectroscopy, and Chemometrics, Anal. Chem. 91 (2019)

700            12444–12452.

701    [18]    R. Ríos-Reina, R.M. Callejón, F. Savorani, J.M. Amigo, M. Cocchi, Data Fusion

702            Approaches in Spectroscopic Characterization and Classification of PDO Wine

703            Vinegars, Talanta  198 (2019) 560–572.

704    [19]    A. Bajoub, S. Medina-Rodríguez, M. Gómez-Romero, E.A. Ajal, M.G. Bagur-

705            González, A. Fernández-Gutiérrez, A. Carrasco-Pancorbo, Assessing the Varietal

706    Origin of Extra-Virgin Olive Oil Using Liquid Chromatography Fingerprints of

707    Phenolic Compound, Data Fusion and Chemometrics, Food Chem. 215 (2017)

708    245–255.

709    [20]    E. Borràs, J. Ferré, R. Boqué, M. Mestres, L. Aceña, A. Calvo, O. Busto, Olive

710    Oil Sensory Defects Classification with Data Fusion of Instrumental Techniques

711    and Multivariate Analysis [PLS-DA], Food Chem. 203 (2016) 314–322.

712    [21]    K.M. Nunes, M.V.O. Andrade, A.M.P. Santos Filho, M.C. Lasmar, M.M. Sena,

713    Detection and Characterisation of Frauds in Bovine Meat in Natura by Non-Meat

714    Ingredient Additions Using Data Fusion of Chemical Parameters and ATR-FTIR

715    Spectroscopy, Food Chem. 205 (2016) 14–22.

716    [22]    Y. Li, Y. Wang, Synergistic Strategy for the Geographical Traceability of Wild

717    Boletus Tomentipes by Means of Data Fusion Analysis, Microchem. J. 140

718    (2018) 38–46.

719    [23]    T.G. Doeswijk, A.K. Smilde, J.A. Hageman, J.A. Westerhuis, F.A. Van Eeuwijk,

720    On the Increase of Predictive Performance with High-Level Data Fusion, Anal.

721    Chim. Acta 705 (2011) 41–47.

722    [24]    C. Márquez, M.I. López, I. Ruisánchez, M.P. Callao, FT-Raman and NIR

723    Spectroscopy Data Fusion Strategy for Multivariate Qualitative Analysis of Food

724    Fraud, Talanta 161 (2016) 80–86.

725    [25]    D. Ballabio, E. Robotti, F. Grisoni, F. Quasso, M. Bobba, S. Vercelli, F. Gosetti,

726    G. Calabrese, E. Sangiorgi, M. Orlandi, Chemical Profiling and Multivariate

727    Data Fusion Methods for the Identification of the Botanical Origin of Honey,

728    Food Chem. 266 (2018) 79–89.

729    [26]    S.M. Azcarate, A. Araújo Gomes, M.R. Alcaraz, M.C.U. Araújo, J.M. Camiña,

730            H.C. Goicoechea, Modeling Excitation – Emission Fluorescence Matrices with

731            Pattern Recognition Algorithms for Classification of Argentine White Wines

732            According Grape Variety, Food Chem. 184 (2015) 214–219.

733    [27]    S.M. Azcarate, A. Araújo Gomes, L. Vera-Candioti, M.C.U. Aráujo, J.M.

734            Camiña, H.C. Goicoechea, Second-Order Capillary Electrophoresis Diode Array

735            Detector Data Modeled with the Tucker3 Algorithm: A Novel Strategy for

736            Argentinean White Wine Discrimination Respect to Grape Variety,

737            Electrophoresis 37 (2016) 1902–1908.

738    [28]    R.D. Snee, Validation of Regression Models : Methods and Examples,

739            Technometrics 19 (1977) 415–428.

740    [29]    R. Bro, Multi-Way Analysis in the Food Industry. Models, Algorithms, and

741            Applications, Thesis manuscript, 1998.

742    [30]    R. Bro, H.A.L. Kiers, A New Eficient Method for Determining the Number of

743            Components in PARAFAC Models, J. Chemometrics 17 (2003) 274–286.

744    [31]    D. Ballabio, V. Consonni, Classification Tools in Chemistry. Part 1: Linear

745            Models. PLS-DA, Anal. Methods 5 (2013) 3790–3798.

746    [32]    A.G.E.K. Smilde, Comments on Multilinear PLS, J. Chemom. 11 (1997) 367–

747            377.

748    [33]    A. Smolinska, J. Engel, E. Szymanska, L. Buydens, L. Blanchet, General

749            Framing of Low-, Mid-, and High-Level Data Fusion With Examples in the Life

750            Sciences, Data Handl. Sci. Techn. 31 (2019) 51-79.

751    [34]    D. Ballabio, R. Todeschini, V. Consonni, Recent Advances in High-Level Fusion

752     Methods to Classify Multiple Analytical Chemical Data, Data Handl. Sci. Techn.

753     31 (2019) 129-155.

754   [35]  A. Fernández, A. Lombardo, R. Rallo, A. Roncaglioni, F. Giralt, E. Benfenati,

755     Quantitative Consensus of Bioaccumulation Models for Integrated Testing

756     Strategies, Environ. Int. 45 (2012) 51–58.

757   [36]  A. Dankowska, W. Kowalewski, Tea Types Classification with Data Fusion of

758     UV–Vis, Synchronous Fluorescence and NIR Spectroscopies and Chemometric

759     Analysis, Spectrochim. Acta - Part A Mol. Biomol. Spectrosc. 211 (2019) 195–

760     202.

761   [37]  B.P. Geurts, J. Engel, B. Ra, L. Blanchet, A. Suppers, E. Szyma, J.J. Jansen,

762     L.M.C. Buydens, Improving High-Dimensional Data Fusion by Exploiting the

763     Multivariate Advantage, Chemometr. Intell. Lab. 156 (2016) 231–240.

764   [38]  Z. Haddi, H. Alami, N. El Bari, M. Tounsi, H. Barhoumi, A. Maaref, N. Jaffrezic-

765     renault, B. Bouchikhi, Electronic Nose and Tongue Combination for Improved

766     Classification of Moroccan Virgin Olive Oil profiles, Food Res. Int. 54 (2013)

767     1488–1498.

768   [39]  Q.Q. Wang, H.Y. Huang, Y.Z. Wang, Geographical Authentication of

769     Macrohyporia Cocos by a Data Fusion Method Combining Ultra-Fast Liquid

770     Chromatography and Fourier Transform Infrared Spectroscopy, Molecules 24

771     (2019) 1320-1338.

772   [40]  A.K. Smilde, M.J. van der Werf; S. Bijlsma, B.J.C. van der Werff-van der Vat,

773     R.H. Jellema, Fusion of Mass Spectrometry-Based Metabolomics Data, Anal.

774     Chem. 77  (2005) 6729–6736.

775  [41]  S. Roussel, J. Roger, P. Grenier, Authenticating White Grape Must Variety with

776       Classification Models Based on Aroma Sensors , FT-IR and UV Spectrometry, J.

777       Food Eng. 60 (2003) 407–419.

778  [42]  M. Silvestri, L. Bertacchini, C. Durante, A. Marchetti, E. Salvatore, M. Cocchi,

779       Application of Data Fusion Techniques to Direct Geographical Traceability

780       Indicators, Anal. Chim. Acta 769 (2013) 1–9.

781  [43]  M. Cocchi, Introduction: Ways and Means to Deal With Data From Multiple

782       Sources, Data Handl. Sci. Techn. 31 (2019) 1-26.

783  [44]  A. Biancolillo, T. Næs, R. Bro, I. Måge, Extension of SO-PLS to multi-way

784       arrays: SO-N-PLS, Chemom. Intell. Lab. Syst. 164 (2017) 113-126.

785  [45]  E. Acar, R. Bro, A.K. Smilde, Data Fusion in Metabolomics Using Coupled

786       Matrix and Tensor Factorizations, P. IEEE 103 (2015) 1602-1620.

787  [46]  M. Hanafi, G. Mazerolles, E. Dufour, E.M. Qannari, Common components and

788       specific weight analysis and multiple co-inertia analysis applied to the coupling

789       of several measurement techniques, J. Chemometrics 20 (2006) 172–183.

790

791

792

793

794

795

796

797

798

799 **Figure captions**

800

801 **Figure 1.** Schematic representation of data analysis workflow

802 **Figure 2.** Typical landscapes of (A) EEM and (B) CE-DAD data for a wine sample

803 showing within each geographical origin –Mendoza (M), San Juan (SJ), Río Negro

804 (RN) and Salta (S)- each grape variety -Chardonay (CH) Sauvignon blanc (SB) and

805 Torrontés (T)-

806 **Figure 3.** Scores plots for the first three LVs exhibiting the best classification results

807 obtained from (A and C) 1-DF and (B and D) 2-DF models showing the differentiation

808 among wines from (A and B) grape variety and (C and D) geographical origin

809 classifications. 95% confidence ellipses for each class are plotted in 3D in each scores

810 plot.
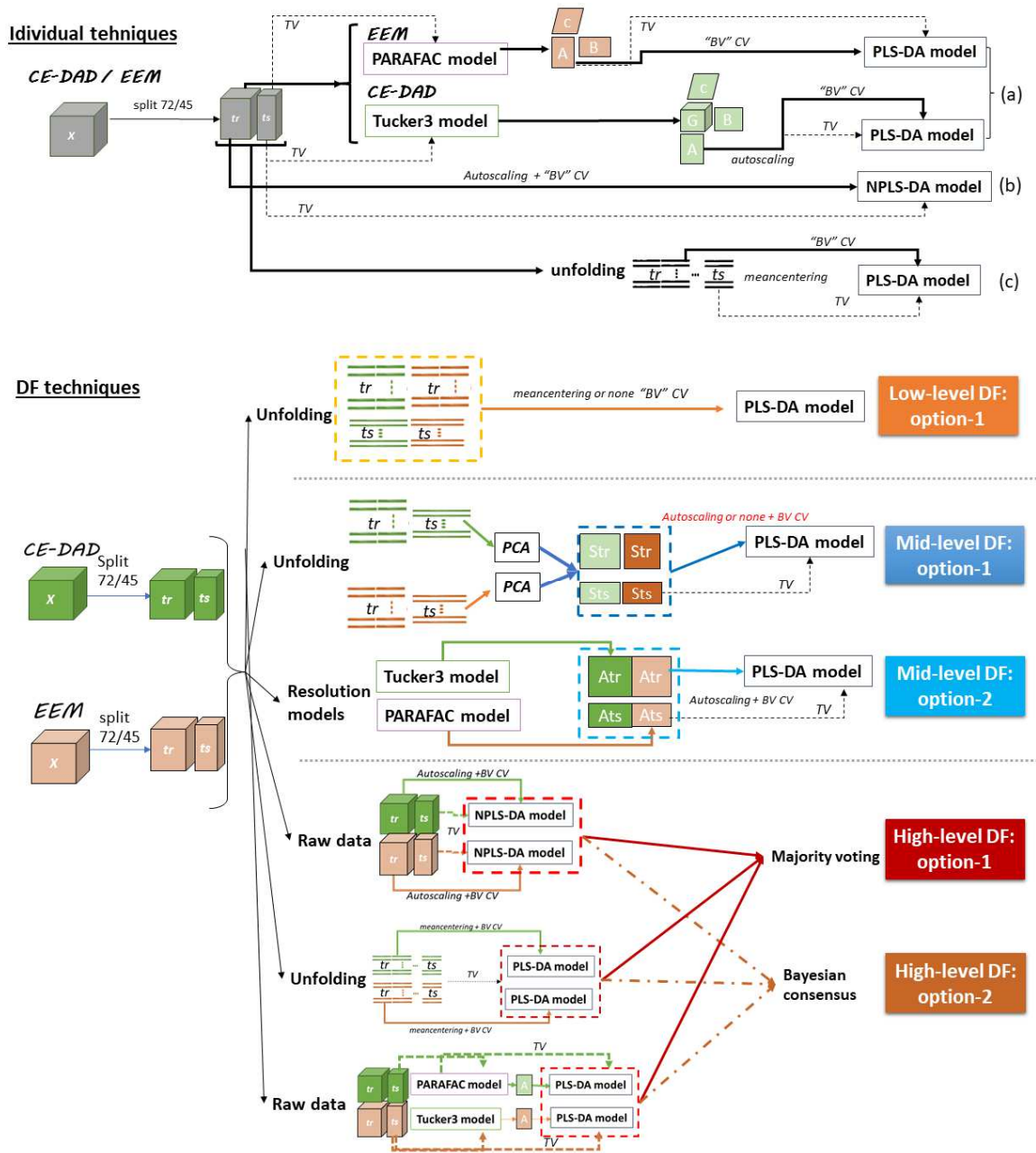
811

812

813

814

815

816

817

818

819

820

821

822

823

33

824    **Figure 1**



825
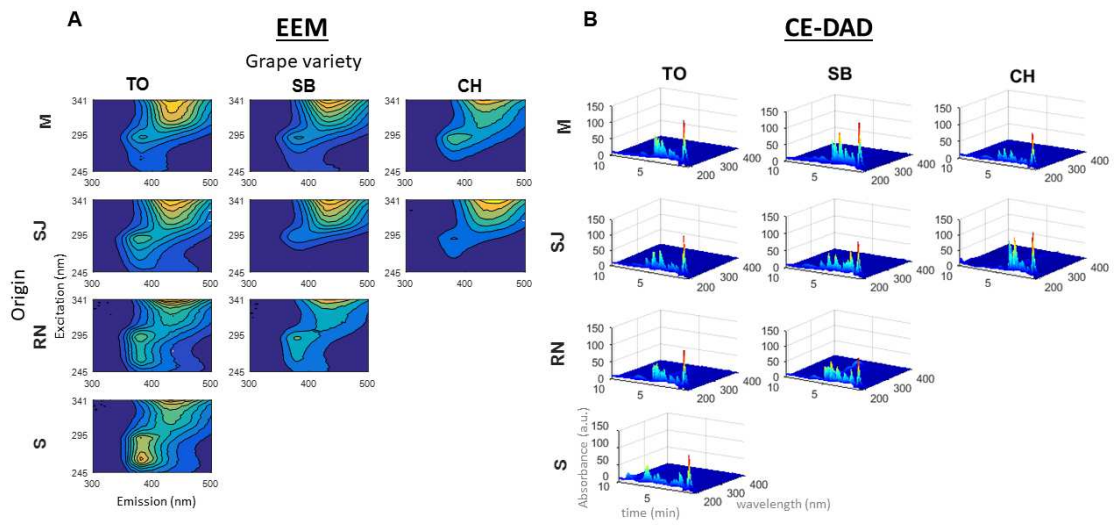
826

827

828

829

830

831

832     **Figure 2**

833



834
835

836

837

838

839
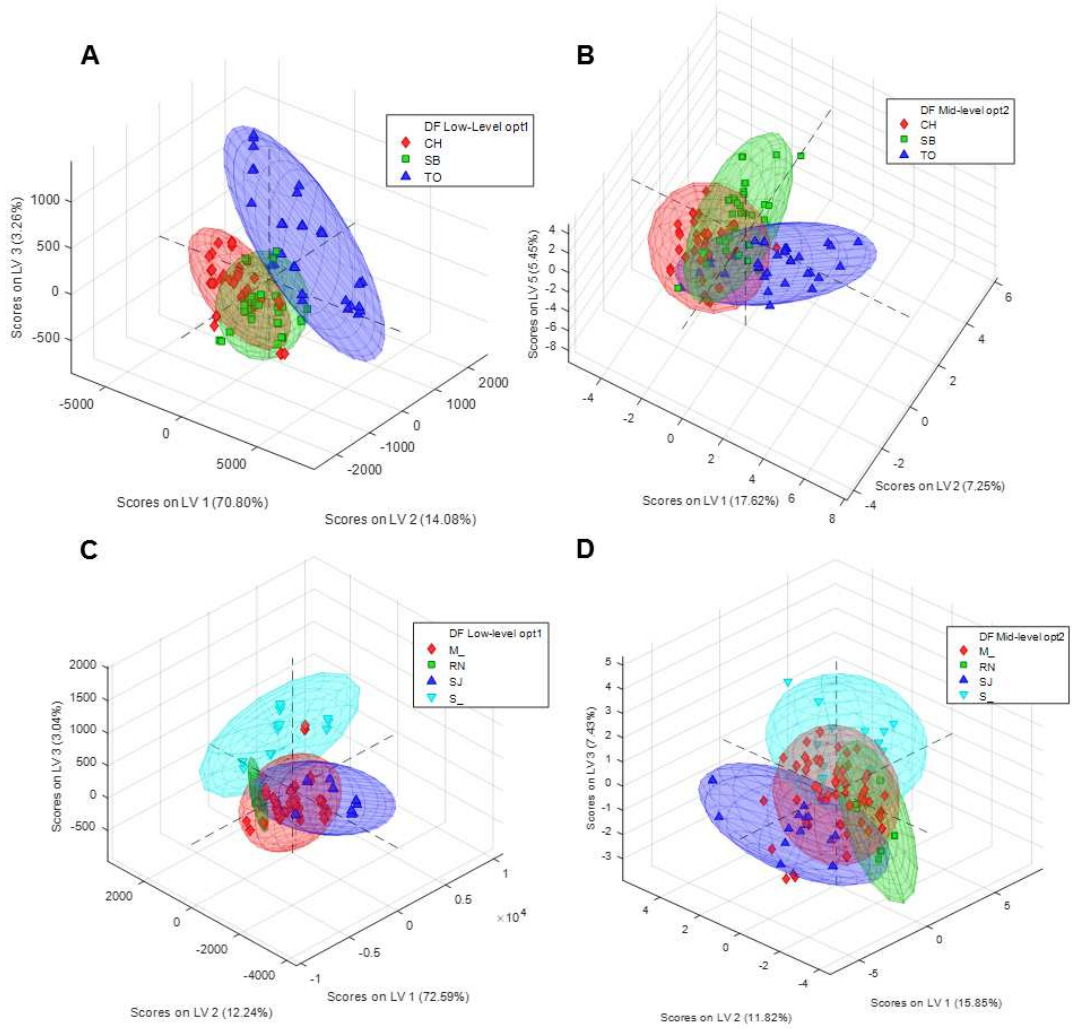
840

841

842

843

844

845

846

847    **Figure 3**

848



849

850

851

852

853

854

855

856 **Tables**

857 **Table 1.** Classification results according to grape variety (Chardonnay -CH-, Sauvignon

858 blanc –SB-, and Torrontés -TO-) obtained in the calibration and prediction stage from

859 individual data blocks (EEM and CE-DAD) and fused data (EEM-CE-DAD) evaluating

860 1-DF, 2-DF and 3-DF. For each model it is displayed the number of samples correctly

861 classified, not error rate (NER) and average precision (PREC) for both calibration

862 (CAL) and prediction (PRED) sets of each evaluated model.

863

| GRAPE VARIETY CLASSIFICATION | DATA STRUCTURE | | Correct classified samples CAL/PRED | | | NER (CAL) | NER (PRED) | PREC (CAL) | PREC (PRED) |
|---|---|---|---|---|---|---|---|---|---|
| | | | CH 24/18 | SB 24/15 | TO 24/12 | | | | |
| EEM | Unfolded data | | 24/9 | 18/9 | 18//6 | 83.3 | 53.3 | 86.1 | 64.3 |
| | Three-way data | | 18/15 | 15/9 | 18/12 | 70.8 | 81.1 | 71.0 | 82.1 |
| | 3-factors PARAFAC scores | | 15/15 | 15/9 | 12/6 | 58.3 | 64.4 | 61.9 | 70.8 |
| CE-DAD | Unfolded data | | 24/15 | 24/3 | 24/9 | 100.0 | 59.4 | 100.0 | 66.7 |
| | Three-way data | | 18/15 | 21/6 | 24/9 | 87.5 | 66.1 | 87.8 | 70.8 |
| | Tucker3 scores | | 24/12 | 24/6 | 24/9 | 100.0 | 60.6 | 100.0 | 66.7 |
| EEM-CE-DAD | LOW-LEVEL | Opt. 1 | 24/15 | 24/6 | 24/9 | 100.0 | 66.1 | 100.0 | 74.1 |
| | MID-LEVEL | Opt. 1 | 24/12 | 18/9 | 24/6 | 91.7 | 58.9 | 93.3 | 69.0 |
| | | Opt. 2 | 24/15 | 24/3 | 24/9 | 100.0 | 59.4 | 100.0 | 66.7 |
| | HIGH-LEVEL | Bayesian consensus | 24/18 | 21/0 | 24/9 | 95.8 | 58.3 | 100.0 | 58.3 |
| | | | 18/15 | 21/9 | 24/12 | 87.5 | 81.1 | 87.8 | 82.1 |
| | | | 24/15 | 24/9 | 24/9 | 100.0 | 72.8 | 100.0 | 77.1 |
| | | Majority voting | 24/6 | 18/0 | 18/3 | 83.3 | 19.4 | 86.1 | 40.7 |
| | | | 18/12 | 12/6 | 18/9 | 66.7 | 60.6 | 67.2 | 65.7 |
| | | | 15/9 | 15/6 | 12/6 | 58.3 | 50.0 | 61.9 | 52.2 |

864

865

866

867

868

869

870

871

872 **Table 2.** Classification results according to geographical origin (Mendoza –M-, Río

873 Negro –RN- San Juan –SJ-, and Salta –S-) obtained in the calibration and prediction

874 stage from individual data blocks (EEM and CE-DAD) and fused data (EEM-CE-DAD)

875 evaluating 1-DF, 2-DF and 3-DF. For each model it is displayed the number of samples

876 correctly classified, not error rate (NER) and average precision (PREC) for both

877 calibration (CAL) and prediction (PRED) sets of each evaluated model.

878

| GEOGRAPHICAL ORIGIN CLASSIFICATION | DATA STRUCTURE | | Correct classified samples CAL/PRED | | | | NER (CAL) | NER (PRED) | PREC (CAL) | PREC (PRED) |
|---|---|---|---|---|---|---|---|---|---|---|
| | | | M 45/27 | RN 6/6 | SJ 12/6 | S 9/6 | | | | |
| EEM | Unfolded data | | 36/18 | 6/3 | 12/3 | 9/3 | 95.0 | 54.2 | 85.0 | 60.4 |
| | Three-way data | | 24/15 | 3/3 | 9/3 | 9/3 | 69.6 | 51.4 | 62.0 | 47.9 |
| | 3-factors PARAFAC scores | | 36/15 | 6/6 | 3/6 | 3/6 | 59.6 | 88.9 | 63.3 | 70.8 |
| CE-DAD | Unfolded data | | 45/24 | 6/0 | 12/3 | 9/0 | 100.0 | 34.7 | 100.0 | 27.9 |
| | Three-way data | | 33/21 | 6/0 | 9/3 | 9/6 | 87.1 | 56.9 | 77.5 | 50.8 |
| | Tucker3 scores | | 33/18 | 6/3 | 9/6 | 6/0 | 78.8 | 54.2 | 66.7 | 51.7 |
| EEM-CE-DAD | LOW-LEVEL | Opt. 1 | 39/21 | 6/6 | 9/6 | 9/6 | 90.4 | 94.4 | 85.7 | 87.5 |
| | MID-LEVEL | Opt. 1 | 33/18 | 6/6 | 6/6 | 6/6 | 72.5 | 91.7 | 65.5 | 79.2 |
| | | Opt. 2 | 36/21 | 6/6 | 12/6 | 9/6 | 95.0 | 95.5 | 85.0 | 87.5 |
| | HIGH-LEVEL | Bayesian consensus | 45/21 | 6/6 | 12/6 | 9/3 | 100.0 | 81.9 | 100.0 | 84.4 |
| | | | 39/15 | 6/6 | 12/6 | 9/6 | 96.7 | 88.9 | 91.7 | 75.0 |
| | | | 39/24 | 6/6 | 12/6 | 9/6 | 96.7 | 97.2 | 86.7 | 91.7 |
| | | Majority voting | 36/18 | 6/0 | 12/0 | 9/0 | 95.0 | 16.7 | 85.0 | 13.6 |
| | | | 21/12 | 3/0 | 6/0 | 9/3 | 61.7 | 23.6 | 55.9 | 25.0 |
| | | | 27/12 | 6/3 | 3/6 | 3/0 | 54.6 | 48.6 | 46.3 | 35.1 |

879
880

**Highlights**

1) Second-order data were fused and chemometrically processed.

2) Multiple strategies for multi-levels data fusion were evaluated.

3) Straightforward approaches for classification purposes are presented.

4) Different degrees of improvement were observed on the results.

5) High-level strategy provided the best classification results.

**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.

☐The authors declare the following financial interests/personal relationships which may be considered as potential competing interests: