This is the peer reviewd version of the followng article:

Assessing feature relevance in NPLS models by VIP / Favilla, Stefania; Durante, Caterina; LI VIGNI, Mario; Cocchi, Marina. - In: CHEMOMETRICS AND INTELLIGENT LABORATORY SYSTEMS. - ISSN 0169-7439. -STAMPA. - 129:(2013), pp. 76-86. [10.1016/j.chemolab.2013.05.013]

Terms of use:

The terms and conditions for the reuse of this version of the manuscript are specified in the publishing policy. For all terms of use and more information see the publisher's website.

03/05/2024 15:19

Chemometrics and Intelligent Laboratory Systems xxx (2013) xxx-xxx

Contents lists available at SciVerse ScienceDirect



CHEMOM-02660; No of Pages 11

### Chemometrics and Intelligent Laboratory Systems

journal homepage: www.elsevier.com/locate/chemolab

#### Assessing feature relevance in NPLS models by VIP **Q4**1

Q1Q2 2 3

> 4 5

### Stefania Favilla<sup>a</sup>, Caterina Durante<sup>b</sup>, Mario Li Vigni<sup>b</sup>, Marina Cocchi<sup>b,\*</sup>

Department of Biomedical Sciences, Metabolic and Neuroscience, University of Modena and Reggio Emilia, Italy

<sup>b</sup> Department of Chemical and Geological Sciences, University of Modena and Reggio Emilia, Italy

#### ARTICLE INFO

#### 6 Article history: 7 Received 9 January 2013 8 9 Received in revised form 17 April 2013 10 Accepted 26 May 2013 11 Available online xxxx 13

- 15Keywords: VIP
- 16 17
- Multi-way data 18 NPLS
- NPLS-DA 19
- Feature selection 20

#### ABSTRACT

Multilinear PLS (NPLS) and its discriminant version (NPLS-DA) are very diffuse tools to model multi-way data 21 arrays. Analysis of NPLS weights and NPLS regression coefficients allows data patterns, feature correlation 22 and covariance structure to be depicted. In this study we propose an extension of the Variable Importance 23 in Projection (VIP) parameter to multi-way arrays in order to highlight the most relevant features to predict 24 the studied dependent properties either for interpretative purposes or to operate feature selection. The VIPs 25 are implemented for each mode of the data array and in the case of multivariate dependent responses con- 26 sidering both the cases of expressing VIP with respect to each single y-variable and of taking into account 27 all y-variables altogether. 28

Three different applications to real data are presented: i) NPLS has been used to model the properties of 29 bread loaves from near infrared spectra of dough, acquired at different leavening times, and corresponding 30 to different flour formulations. VIP values were used to assess the spectral regions mainly involved in deter- 31 mining flour performance; ii) assessing the authenticity of extra virgin olive oils by NPLS-DA elaboration of 32 gas chromatography/mass spectrometry data (GC-MS). VIP values were used to assess both GC and MS dis- 33 criminant features; iii) NPLS analysis of a fMRI-BOLD experiment based on a pain paradigm of acute 34 prolonged pain in healthy volunteers, in order to reproduce efficiently the corresponding psychophysical 35 pain profiles. VIP values were used to identify the brain regions mainly involved in determining the pain in- 36 tensity profile. 37

© 2013 Published by Elsevier B.V. 38

CHEMOMETRICS BAND INTELLIGENT LABORATORY

SYSTEMS

#### 30

03

42 41

#### 43 1. Introduction

Multilinear PLS (NPLS) and its discriminant version (NPLS-DA) are **O5**44 very diffuse tools to model multi-way data arrays. NPLS represents 45the multi-way extension of two-way partial least squares regression 46 (PLS) for multi-way data and was first developed by Bro in 1996 [1] 47 and successively by Bro, Smilde and De Jong [2-4]. 48

It has been demonstrated that multi-way data analysis tools, tak-4950ing into account the multi-way structure of data are much more efficient compared to unfolding procedures, that is re-arranging the 51 multi-way data into a two-way matrix structure and then applying 52bilinear models. Multi-way analysis allows simplifying the interpreta-53 tion of the results and providing more adequate and robust models 54using relatively few parameters [5,6]. While this is true in general, it 55 is worth noticing that when dealing with real-time monitoring, e.g. 56in batch process monitoring, N-way models may not represent a 57real advantage with respect to adopting a proper unfolding/refolding 58procedure as by using Multiway-PCA [7]. 59

60 In particular, the use of NPLS shares all the advantages of latent varjable based regression and discrimination methods, from the point of 61 view of data visualization and interpretation [8-10]. In fact, analysis of 62

0169-7439/\$ - see front matter © 2013 Published by Elsevier B.V. http://dx.doi.org/10.1016/j.chemolab.2013.05.013

NPLS weights and NPLS regression coefficients allows data patterns, fea- 63 ture correlation and covariance structure to be depicted. 64

However, it is often needed to define which are the most relevant 65 features to predict the studied dependent properties either for interpre- 66 tative purposes, e.g. to provide a better understanding of the underlying 67 process that generated the data, or to operate feature selection in order 68 to reduce the noise generated by irrelevant features or to reduce data 69 redundancy. 70

Some of the several variable selection methods applied to two-way 71 data matrices in the context of PLS regression [11], such as interval PLS 72 (iPLS) [12] or genetic algorithms [13], can be as well suited for NPLS if 73 the X-block multi-way data array has only one spectral dimension, e.g. 74 samples  $\times$  spectral profiles  $\times$  time [14]. When the data array has two 75 spectral dimensions, or more generally when a two-dimensional signal 76 map characterizes each sample, as generated by hyphenated analytical 77 techniques, such as emission/excitation fluorescence, chromatography/ 78 mass spectrometry, etc., these variable selection methodologies present 79 significant challenges and it is suggested to apply them after unfolding 80 the data array [15]. However, in this way, the multi-way data structure 81 is not taken into account in the variable selection step, thus loosing the 82 multi-way analysis advantage.

Moreover, a general distinction can be made among tools, which ac- 84 complish feature selection by deleting a set of features and re-assessing 85 the performance of the reduced models, thus requiring extensive model 86

<sup>\*</sup> Corresponding author. Tel.: +39 059 2055029; fax: +39 059 373543. E-mail address: marina.cocchi@unimore.it (M. Cocchi).

S. Favilla et al. / Chemometrics and Intelligent Laboratory Systems xxx (2013) xxx-xxx

In the case of a two-way data matrix, 
$$\mathbf{Y}_{\text{LM}}$$
 is defined by: 140

validation, and those that operate a ranking of the features according to their relevance. Concerning the latter type, congruence loadings [16], VIP (Variable Importance in Projection) [17–20], and selectivity ratio [21] have gained increasing attention as an important measure of each explanatory variable or predictor.

The aim of this work is to extend the VIP method to multi-way ar-92rays and to develop accompanying code. A similar attempt has been 93 reported previously [22] but our formulation is, in our opinion, 94more straightforward and closer to VIP definition for two-way data. 95 In fact, the VIP definition given in Ref. [22] does not reproduce the 96 two-way VIP formulation, which consists, for each X-variable, of a 97 98 sum, over latent variables, of its PLS-weight weighted by the percentage of explained Y variance. The reason is due to the fact that the 99 NPLS mode 1 scores for X-block (T), which are linked to Y by the 100 NPLS inner-relation, are substituted by a different projection of 101 unfolded-X through NPLS weights. Usually, in the case of several de-102pendent responses (multivariate Y) VIP is defined taking into account 103 all y-variables altogether. Here we consider as well, the possibility of 104 expressing VIP with respect to each single y-variable (this is a further 105106 difference with the approach presented in Ref. [22] that does not allow this possibility). This offers higher flexibility to the method 107 and can be particularly useful to interpret discriminant NPLS-DA 108 models, since the VIP for single y-variables corresponds, in this case, 109 to the most discriminant feature for each category. However, it is be-110 yond the aim of this paper, to compare the use of VIPs with other fea-111 ture selection methodologies for multi-way arrays, actually in the 112 applications presented here variable selection is not operated and 113 114 VIPs are used more on an interpretative ground.

### 115 2. Methods

143

116 2.1. Multilinear partial least squares (NPLS)

Multilinear PLS (NPLS) represents the extension of two-way partial
least squares regression (PLS) to data arrays of any order considering
both X and Y-blocks. In the following, the method is described considering the case of a three-way data array, X, but the extension to further dimensions can be simply deduced. As for the dependent variables
Y-block, we will describe here the case of a two-way matrix, but the
method can be easily extended to higher orders in the Y-block [1].

Specifically, PLS regression aims to find a relationship between a set 124 of predictor (independent) data, X, and a set of responses (dependent), 125Y. In the more general case, the arrays of independent, X and dependent 126127Y variables are decomposed in such a way that the score vectors from these models have pair-wise maximal covariance [3,4]. Multilinear 128 PLS was firstly developed as a PARAFAC-like model of X and it was 129shown that the method could be easily extended to any desired order 130 for both **X** and **Y** arrays. This method was further elaborated and lastly 131 improved with respect to residual analyses by introducing a core 132133 array in the model of X [2].

Considering an  $\underline{X}$  array of dimension  $I \times J \times K$ , the NPLS model is obtained by modeling  $\underline{X}$  as in Tucker3 decomposition:

$$\mathbf{X} = \mathbf{T}\mathbf{G}_{\mathbf{X}} \left( \mathbf{W}^{\mathbf{K}} \otimes \mathbf{W}^{\mathbf{J}} \right)^{\mathrm{T}} + \mathbf{E}_{\mathbf{X}}$$
(1)

where **X** is the **X** array unfolded to an  $I \times JK$  matrix, **T** holds the first mode scores (sample mode), **W**<sup>*J*</sup> and **W**<sup>*K*</sup> are the second and the third mode weights, respectively. The symbol  $\otimes$  denotes the Kronecker product [5].

140 **G**<sub>X</sub> is the matricized core array of size  $F \times F \times F$  where *F* is the 141 number of NPLS components (factors) and it is defined by:

$$\mathbf{G}_{\mathbf{X}} = \mathbf{T}^{+} \mathbf{X} \left( \left( \mathbf{W}^{K} \right)^{+} \otimes \left( \mathbf{W}^{J} \right)^{+} \right)^{\mathrm{T}}.$$

144 Here the superscript '+' means that the Moore–Penrose is pseudo 145 inverse.

$$\mathbf{Y} = \mathbf{U}\mathbf{Q} + \mathbf{E}_{\mathbf{Y}} \tag{3}$$

where U holds the Y scores and Q is the loading matrix.143 $\underline{E}_X$  and  $E_Y$  hold  $\underline{X}$  and Y residuals, respectively. In analogy with the149two-way PLS algorithm, the weights are determined such that the150scores obtained from the  $\underline{X}$  decomposition (T) have maximum151covariance with the scores obtained from Y decomposition (inner152relation:  $\mathbf{U} = \mathbf{TB} + \mathbf{E}_U$ ).153

By regressing the data onto their weights vectors, a score vector is 154 found in the  $\underline{X}$ -space providing a least squares model of the  $\underline{X}$  data. 155 Furthermore, by choosing the weights such that the covariance between X and Y is maximized a predictive model is obtained as: 157

$$\mathbf{Y} = \mathbf{TBQ} + \mathbf{E}_{\mathbf{Y}}.$$
 (4)

159

163

100

166

179

Regression coefficients that apply directly to  $X(I \times JK)$  may also be 160 derived [4,22]: 161

$$R = \left[ w_1 \left( I - w_1 w_1^T \right) w_2 \dots \prod_{f=1}^{F-1} \left( I - w_f w_f^T \right) w_f \right]$$
(5)

$$\mathbf{B}_{\mathsf{PLS}} = \mathbf{RB}\mathbf{1} \tag{6}$$

$$\mathbf{Y} = \mathbf{X}\mathbf{B}_{\text{PLS}}.$$
 (7)

The NPLS-DA formulation is the same but the dependent variable 168 block is a matrix **Y** holding the class information, i.e. for each category 169 a y-variable is defined as a dummy variable assuming values one/ 170 minus one to indicate class membership or not (notation one/zero 171 is also used). As the predicted y-values can assume real values and 172 not only minus one and one, classification of the samples is accomplished by assigning the sample to the category corresponding to 174 the highest value of the predicted response, i.e. if the predicted vector 175 of responses for an unknown sample, is:  $[-0.5\ 0.8\ 0.5]$  (in the case of 176 three classes problem), it will be assigned to class two. 177

2.2.1. Two-way case

The variable importance in the projection (VIP) [17,19] represents 180 Q6 the influence of each variable j of the data matrix  $\mathbf{X}_{IJ}$  on the model of 181 the responses matrix  $\mathbf{Y}_{I,M}$  182

$$VIP_{j}^{2} = S_{f}w_{jf}^{2} \cdot SSY_{f} \cdot J/(SSY_{tot.expl.} \cdot F)$$
(8)

where, *F* is the number of latent variables of the PLS model and *J* the **184** number of **X** variables. 185

In the case of mono-dimensional  $\mathbf{y}^{l \times 1}$  holds: 186

$$SSY_f = b_{f,f}^{\ 2} \mathbf{t}_f^{\ T} \mathbf{t}_f \qquad SSY_{\text{tot.expl.}} = \mathbf{b}^2 \mathbf{T}^{\mathsf{T}} \mathbf{T}$$
(9)

where **T** is the **X** score matrix and **b** the PLS inner relation coefficients. Thus a VIP value for each variable is computed in order to quantify its importance by using the PLS weight *w<sub>if</sub>* weighted by how much of **y** is explained in each model dimension (latent variable).

VIP formulation as originally proposed [17] is intended to be a parameter varying in a fixed range since the sum of squared VIP for all 193 variables is the sum to the number of variables. Thus, the variables 194 **Q7** with a VIP value larger than 1 (i.e. larger than the average of square 195 VIP values) have an above average influence on the model and are, 196 therefore, considered the most relevant for explaining **Y**. The choice 197 of the VIP threshold to assess the salient variables is a critical issue, 198 as in any ranking method. The original proposal, that will be adopted 199 here as well, of a threshold of one is acceptable if variable relevance is 200

Please cite this article as: S. Favilla, et al., Assessing feature relevance in NPLS models by VIP, Chemometrics and Intelligent Laboratory Systems (2013), http://dx.doi.org/10.1016/j.chemolab.2013.05.013

(2)

2

87

88

89

90

91

S. Favilla et al. / Chemometrics and Intelligent Laboratory Systems xxx (2013) xxx-xxx

201 discussed but feature selection is not accomplished. In the cases of

202 marker identification and variable selection, resampling methods 203 such as bootstrap are more appropriate [19,20] to assess the signifi-

204 cance of the VIPs.

#### 205 2.2.2. Three-way case

In the case of a two-dimensional  $\mathbf{Y}(I \times M)$  the previous relation applies to each mode, e.g. in the case of a three-way array  $\underline{\mathbf{X}}(I \times J \times K)$ :

$$\operatorname{VIP}_{j}^{2} = \Sigma_{f} w_{jf}^{2} \cdot \operatorname{SSY}_{mf} J / \left( \operatorname{SSY}_{\operatorname{tot.expl.},m} \cdot F \right)$$
(10)

$$\operatorname{VIP}_{k}^{2} = \Sigma_{f} w_{kf}^{2} \cdot \operatorname{SSY}_{mf} \cdot K / \left( \operatorname{SSY}_{\operatorname{tot.expl.},m} \cdot F \right)$$
(11)

where, *F* is the number of total latent variables, *J* the number of  $\underline{X}$  variables in Mode2 and *K* the number of variables in Mode3. For each latent variable *f*:

$$SSY_{tot.expl.,m} = \Sigma_i \Big( \mathbf{T}_{(l \times F)} \mathbf{B}_{(F \times F)} \mathbf{q}^{\mathsf{T}}_{(m,F)} \Big)^2$$
(12)

**214** and each y-variable  $y_m$ :

$$SSY_{mf} = \Sigma_i \left( \mathbf{t}_f b_{ff} \mathbf{q}_{mf} \right)^2 \tag{13}$$

216 where, I is the number of samples in Mode2, T is the Mode1 score ma-218 trix, **B** holds the NPLS inner relation coefficients and **Q** the **Y** loadings. For a given model dimension *f* and each variable *j*, the VIP value is 219 given by the squared weight  $w_{if}^2$  of that parameter (i.e. the weight  $w_{if}$  in-220 dicates the importance of the *j*th variable in the model dimension *f*), mul-221 222 tiplied by the percent of **Y** explained sum of squares by that *f* dimension. 223The variable importance is then normalized so that VIP<sup>2</sup> equals the 224 number of the variables.

While considering all **Y** variables together, Eqs. (12) and (13) are reduced to:

$$SSY_{tot.expl.} = \Sigma_i \left( \mathbf{T}_{(I \times F)} \mathbf{B}_{(F \times F)} \mathbf{Q}^{\mathsf{T}}_{(M,F)} \right)^2$$
(14)

NIR spectra

$$SSY_f = \sum_m \sum_i \left( \mathbf{t}_f \boldsymbol{b}_{f,f} \mathbf{q}^{\mathrm{T}}_{m,f} \right)^2 \tag{15}$$

and Eq. (10) is reduced to:

$$VIP_{j}^{2} = \Sigma_{f} w_{jf}^{2} \cdot SSY_{f} J / \left( SSY_{tot.expl.} \cdot F \right).$$
(16)

Extension to the other **Y** modes can be easily obtained. 233

230

232

234

258

### 3. Data sets and pretreatment

In this study, we present applications of VIP to different three-way 235 Q8 data sets. Two data sets are related to optimization of food processing 236 and authentication issue for products with protected denomination of 237 origin, respectively, and the third one is related to a neuroscience 238 problem. Each data set allows exploring the different situations, pre- 239 dictive and discriminant models, partial and overall VIP contribution 240 with respect to **Y** block together with the different aspects of comple- 241 mentary information that VIP can highlight with respect to e.g. NPLS 242 weights or regression coefficients. 243

Li Vigni and Cocchi [24] presented a multi-way study related to 245 Q9 the influence of flour formulation on bread quality. Ten different 246 flour mixtures were investigated by means of Near Infrared Spectros-247 copy (NIRS) to obtain information on flour performance in a critical 248 phase such as dough leavening. For each mixture, a laboratory-scale 249 bread making experiment was carried out according to a standard-250 ized recipe and the leavening phase of each dough sample was mon-251 itored by means of NIRS at different times. NPLS was applied to model 252 the properties of bread loaves (dimensions, volume, weight, height) 253 from near infrared spectra, acquired at different leavening times, of 254 the dough obtained from different flour formulations. 255 The data are arranged as follows and schematically shown in Fig. 1: 256

- *X*-block: a three-way array  $\mathbf{X}(I \times J \times K)$  (10 flour mixtures × 173 257

NIR wavelengths  $\times$  7 leavening time intervals)

Bread properties



S. Favilla et al. / Chemometrics and Intelligent Laboratory Systems xxx (2013) xxx-xxx



Fig. 2. Data set arrangement for EVOO data. (Left) GC–MS data (X): Mode1, samples; Mode2, Retention times (1514 time points); Mode3, 77 selected m/z fragments. (Right) Dummy class variables (Y): Mode1, samples; Mode2 classes (3).

259 - *Y*-block:  $\mathbf{Y}(I \times M)$  (10 flour mixtures  $\times$  4 bread loaf properties: 260 weight, height, diameter and density).

NIR signals were preprocessed by applying Savitsky–Golay Smooth ing (15 points window, second order polynomial) coupled to Standard
 Normal Variate normalization (SNV) to remove the baseline shift.

VIP values are used to assess the spectral regions mainly involvedin determining flour performance.

266 3.2. Data set 2: Extra virgin olive oil (EVOO)

The data set [25] consists of a set of extra virgin olive oil (EVOO), be-267 longing to different olive cultivars and coming from different Mediter-268ranean areas: Liguria (Northern Italy), Apulia (Southern Italy), Greece, 269270 Tunisia and Spain. The aim is to assess the authenticity of Ligurian 271EVOO that has been designed by protected denomination of origin (PDO) certification and represents one of the most highly esteemed 272EVOOs, of high economic value. The EVOO samples have been charac-273274terized by the analysis of aroma (Head Space Solid Phase Micro Extraction coupled with Gas Chromatography-Mass Spectrometry, i.e. 275HS-SPME/GC-MS), which is well suited for analyzing the volatile 276fraction that is of relevance for the sensory quality of olive oil. The 277278differentiation among classes has been obtained by NPLS-DA, defining 279three classes: Liguria, Apulia and Foreign, which includes the EVOO 280 from Turkey, Spain and Greece.

The data set is arranged as follows (Fig. 2):

281

282 - *X*-block: a three-way array  $\underline{X}(I \times J \times K)$  (73 EVOO samples  $\times$  1514 283 retention time points  $\times$  77 m/z fragments)

- *Y-block*:  $\mathbf{Y}(I \times M)$  (73 EVOO samples × 3 dummy variables holding class memberships).

For each class the data were randomly split in a training and validation 286 (test) set as shown in Fig. 2. The training set was preprocessed by center-287ing across the first mode, block-scaling within the second mode, by defin-288ing four retention time regions in order to allow both major and minor 289 constituents to contribute to the model without up-weighting baseline 290291 contribution [25] and scaled by inverse standard deviation within the third mode (selected mass fragments). The pretreatments were applied 292in the order Mode3, Mode2, and Mode1. 293

VIP values were used to assess both GC and MS discriminant features.

### 3.3. Data set 3: Neuroscience data set

296

This data set derives from a functional magnetic resonance imag- 297 ing (fMRI) experiment where the psychophysical pain profile, corre- 298 sponding to subjective responses to acute prolonged noxious 299 stimulation of one hand, was acquired in healthy volunteers. The ex- 300 periment lasted 20 min (300 time points), the sensory intensity of 301 pain (psychophysical pain profile) and the hemodynamic response 302 (blood-oxygen-level contrast registered by a magnetic resonance 303 pulse sequence, fMRI-BOLD signal) were recorded simultaneously 304 during the experiment. The functional fMRI-BOLD signals (fMRI-BOLD 305 time series) acquired at each brain voxel, as described in Prato et al. 306 [26], were summarized for forty four brain regions of interest (ROIs) 307 by taking the first principal singular vector (1st-SVD) of the data ma-308 trix containing the fMRI-BOLD time series for each voxel in that specific ROI. 310





**Fig. 3.** Data set arrangement for Neuroscience data. (Left) fMRI-BOLD data ( $\underline{X}$ ): Mode1, Times Points (300); Mode2, fMRI-BOLD intensity for the 44 ROIs; Mode3, volunteers (10). (Right) psychophysical responses ( $\underline{Y}$ ): Mode1, Time points (300); Mode2, perceived pain intensity by volunteers (10), in this case the actual  $\underline{Y}(300 \times 1 \times 10)$  has been rotated for illustration purposes.

S. Favilla et al. / Chemometrics and Intelligent Laboratory Systems xxx (2013) xxx-xxx



Fig. 4. Data set for Bread. (Left) Y-loadings plot (Q) for the three NPLS components, F1, F2, and F3; (Right) NPLS-weight plot for Mode2 (spectra), F1 vs. wavelengths.

NPLS was applied to build a model that could express the main
 variation of the ROI time series of different volunteers and obtain a
 fitted model that could reproduce the corresponding psychophysical
 pain profile efficiently.

315 The data arrays have been arranged as (Fig. 3):

<sup>316</sup> - *X*-block: a three-way array  $\underline{\mathbf{X}}(I \times J \times K)$  (300 time points × 44 <sup>317</sup> ROIs × 10 volunteers).

<sup>318</sup> - *Y*-block: the  $\underline{\mathbf{Y}}(I \times 1 \times K)$  array is actually a matrix  $\mathbf{Y}(I \times K)$  com-<sup>319</sup> prised of 300 time points (psychophysical pain profile) × 10 volun-<sup>320</sup> teers, as shown in Fig. 3, and so computationally handled as such.

The choice of defining time as mode one was motivated by the applicability of the model. In fact, for this approach the main scope was to identify those ROI time series strictly connected (i.e., in terms of covariance) with the psychophysical pain profile of each volunteer (see Fig. 3).

The X and Y data were not centered or scaled within any mode. VIP values were used for ranking the ROIs according to their relevance in the NPLS model hence to depict brain region activation profile in response to pain stimulus.

### 330 4. Results and discussion

### 331 4.1. Bread

Near infrared spectra acquired on dough at subsequent leavening 332 times is an efficient way to characterize the leavening process. In partic-333 ular, NPLS was used to study the relationship between the modifications 334 335 recorded by the NIR signal during the leavening time and four properties measured on bread loaves, namely height, weight, volume and den-336 sity. The dough samples correspond to ten different wheat mixtures 337 (combining four distinct wheat varieties), performed according to a 338 G-optimal design, thus bread performance can be linked to best mixture 339 formulation in terms of wheat varieties. 340

341The dimensionality of NPLS model was chosen on the basis of the342best compromise of the minimum values of RMSECV for the four343properties, modeled as a single Y block. Leave One Out cross valida-344tion was chosen due to the limited number of samples, ten.

345A three factor NPLS model explains 75% of the total Y variance in fit,346with acceptable performances in fit for all responses and generally poor347results in cross-validation. As discussed in Ref. [24] these are mainly due348to two mixtures showing extreme behavior in property space. However,349a qualitative identification of the relationship among NIR signal variabil-350ity along with the progression of the leavening step and bread proper-351ties can be obtained by inspecting NPLS model results.

In particular, **Y** loading plot (Fig. 4, left plot) shows that the first NPLS latent variable (F1) mainly models bread loaf height and volume, while latent variables 2 and 3 (F2 and F3) model bread loaf weight. NPLS Mode2 weights indicate which spectral regions mostly influence each factor and hence bread loaf properties. For example, the first factor for Mode2 weights (**w**<sub>i1</sub>, Fig. 4, right plot), which is linked to height/volume properties, shows that the most relevant contributions can be assigned 358 to water and its redistribution across the macro-polymeric components 359 of the dough, such as gluten and starch. In fact, the most relevant fre- 360 quencies are those near 1400, and between 1900 and 1950 nm, for 361 which weights have a negative sign, and above 2100 nm, for which 362 weights have positive sign. These regions correspond, respectively, to ab- 363 sorptions that can be associated to the O–H stretching first overtone, the 364 O–H bending second overtone mode and to overtone and combination 365 mode contributions from the starch, protein and lipid fractions.

Inspection of Mode3 weights (leavening times, Fig. 5) indicates 367 that height and volume are influenced by initial (t0–t10) and last 368 leavening phases (yeast activity, dough strength); in fact these 369 times have relevant weight values on the first NPLS factor (Fig. 5, 370 top plot); weight is mainly influenced by initial (t0) time; hence 371 flour properties are mostly important, as shown by second and third 372 NPLS factor weights (Fig. 5, middle and bottom plots). 373

The regression coefficient maps are not so easily interpretable, see 374 Fig. 6, while in order to assess the most relevant spectral regions and 375 leavening times VIP values, proved to be very useful (Fig. 7A and B). 376 Fig. 7A reports the VIP values for Mode3. It is interesting to notice that, 377 for bread weight, the spectrum at the beginning of the leavening phase 378 (time 0) has VIP values higher than one while height property exhibits 379 significant alternating variation at times 10, 40 and 60. The other two 380 properties, volume and density, show that the times at which the most 381 significant variations in the NIR spectrum are recorded correspond main- 382 ly to the initial conditions (0 and 10 min) and the final one (60 min). 383

The most influential NIR spectrum regions, as highlighted by Mode2 384 VIP values (Fig. 7B), are the same for all properties: at about 1425 nm. 385



Fig. 5. Data set for Bread. Plot of NPLS weight for Mode3 (leavening times).

S. Favilla et al. / Chemometrics and Intelligent Laboratory Systems xxx (2013) xxx-xxx



Fig. 6. Data set for Bread. 2D plot of the NPLS regression coefficient map for the y-property Volume.



Fig. 7. Data set Bread. (A) Squared VIP values for Mode3 for each bread property; (B) squared VIP values for Mode2 for each bread property.

S. Favilla et al. / Chemometrics and Intelligent Laboratory Systems xxx (2013) xxx-xxx



Fig. 8. Data set for Bread. Plot of NPLS regression coefficients at selected leavening times. Gray bar highlights the most relevant NIR spectral regions.

There are also contributions from C–H combination and first O–H stretching overtone, between 1900 and 1950 nm, corresponding to absorptions which can be associated to the second overtone O–H bending mode, and above 2100 nm, where overtone and combination mode contributions from the starch fractions are present.

391 On the basis of VIP analysis it is now possible to plot for each modeled bread property the NPLS regression coefficients corresponding 392 only to the most influential leavening times, so that the spectral contri-393 butions can be discussed in terms of increasing/decreasing of specific 394 absorption bands, taking into account the regression coefficients sign, 395 in an easier way. Fig. 8, as an example, shows the regression coefficient 396 plot for the most relevant leavening times for the bread Volume. The 397 most significant spectral regions discussed above change correlation 398 399 sign at the different times as it may be expected by the dynamic of the leavening process where several rearrangements of the starch and glu-400 ten networks take place [24]. 401

The combination of multi-way methods applied to NIR spectra is here useful to supervise changes of the system according to the

leavening time, and can be used as a reference to evaluate the behav- 404 ior of dough obtained from different wheat flour mixtures, and poten- 405 tially to identify anomalous leavening situations. Also, it has been 406 shown that VIP values give the same information as the joint discus- 407 sion of Y loadings and NPLS weights, but with a more direct highlight 408 of the most influential contribution for each property. Moreover, they 409 can be taken into account as a guide to plot in an interpretable manner the NPLS regression coefficients, which may otherwise result of 411 difficult interpretation. 412

#### 4.2. Extra virgin olive oil (EVOO) 413

The performance of the NPLS-DA classification models has been 414 evaluated by means of percentage of correct classification in cross- 415 validation (CV), by using venetian blind with six cancelation groups. 416 Three NPLS components gave the best performance with 100% correct 417 classification in fit and CV for Liguria and Foreign classes and 92% (fit 418 and CV) for Apulia. The test set for Liguria and Foreign classes was 419



Fig. 9. Data set for EVOO. Plot of the squared VIP. (Left) Mode2 VIPs vs. Retention Times, threshold sets to 5. (Right) Mode3 VIPs vs. m/z fragments, threshold sets to 1.

S. Favilla et al. / Chemometrics and Intelligent Laboratory Systems xxx (2013) xxx-xxx



Fig. 10. Data set for EVOO. Average TIC chromatograms for each class, overlaid with a shift of 0.5 on time axis (x-axis) and of 1 on y-axis. Dots indicate retention time with VIP values higher than one.

420 correctly classified with no error and with one misclassified sample for421 Apulia.

8

The VIP values for Mode2 and Mode3 for each y-variable, hence cat-422 egory, are reported in Fig. 9. The VIP values threshold is drawn at 1 for 423424 Mode3 (m/z fragment), Fig. 9 on the right; while for Mode2 (Retention time) is drawn at 5, Fig. 9 on the left. This choice is motivated by consid-425ering that this is the number of points generally defining a peak, thus a 426 variable to be considered important has to have a contribution at least 427 as one peak in the signal. Taking into account the most important vari-428 ables of both modes some considerations about the main compositional 429difference of each EVOO category can be made. The chromatographic 430peaks, which seem to characterize the volatile fraction of Liguria olive 431 432 oil, are some low molecular weight compounds (Retention time about 16 min; m/z < 40) and C6 linear unsaturated aldehydes (Retention 433 time regions at about 27 min and 32 min; m/z region 61-70) character-434 istics of high quality virgin olive oils. The latter compounds are also 435

present in Apulia and the Foreign class but in a lower amount. The Apulia class is mainly characterized by the retention time region 57–60 min and by higher molecular weights compounds (m/z values higher than 80), such as alpha-copaene (Rt = 57.7 min; m/z 81, 93, 105). The high VIP values (mainly for Liguria and Foreign classes) at the retention time regions 39–40 and 43–45 min highlight compounds that are more 441 related to, i.e. specific for, the Foreign class as it can be seen from Fig. 10, 442 where the average total ion count chromatograms (TIC) for each class are shown.

The NPLS weights for Mode2 and Mode3 for each of the three NPLS 445 components (F1, F2 and F3) are reported in Fig. 11. In order to interpret 446 these plots the Y loadings have to be inspected (figure not shown). These 447 indicate that Liguria (y1) has high positive loadings on components 2 448 and 3 (F2 and F3) and close to zero on component 1 (F1), the opposite 449 holds for Foreign (y3) while Apulia has high positive loadings values 450 on component 1 (F1) and almost zero on the other two components. 451



Fig. 11. Data set EVOO. (a) NPLS weights for Mode2 vs. retention times, for NPLS components F1 (top), F2 (middle) and F3 (bottom). (a) NPLS weights for Mode3 vs. m/z fragments, for NPLS components F1 (top), F2 (middle) and F3 (bottom).



**Fig. 12**. Data set for Neuroscience. Plot of the squared VIP calculated for all **Y** altogether vs. ROIs, the threshold is set to 1. Labels refer to ROIs: Contralateral (*C*) and Ipsilateral (*I*) to the injected hand, respectively; Pre-Central Gyrus (*PreCG*), Middle (*Mid*); Medial Superior Frontal Gyrus/Paracentral Lobule (*MSFG*); Parietal Operculum (*PO*); Posterior Insula (*pINS*); Anterior Insula (*alns*); Mid-Cingulate Cortex (*MCC*); Medial Thalamus (*Med Th*); Lateral Thalamus (*Lat Th*); Caudate Nucleus, head (*hCaud*); Caudate Nucleus, body (*bCaud*); Inferior Parietal lobule (*IPL*).

452 Taking this into account by looking at the Mode2 weights (Fig. 11a) the

453 observations about the characteristics retention time regions for each

454 class, made on the basis of the VIP values, are confirmed, e.g. the negative 455 sign of the weights at retention time 39–40 and 43–45 indicates that this

456 region is characteristic of the Foreign class.

### 457 4.3. Neuroscience data set

The goal was to build a model that, starting from the fMRI-BOLD characterization as expressed by the 1st SVD component of each ROIs fMRI-BOLD time series, could predict efficiently the psychophysical pain intensity for each volunteer. A NPLS model with 4 Latent Variables (LVs) was selected according to Cross Validation results, lowest RMSECV value, i.e. 28.12, with explained **Y** variance around 95%. The VIP values were used for ranking the brain regions (ROIs) main-464 ly involved in pain perception. In this case it is interesting to consider 465 the VIP values calculated for the overall **Y**-responses, to gain a global 466 picture of brain activation common to all subjects. Fig. 12 reports the 467 VIP values versus ROI number as bar plot, highlighting the ROIs showing 468 VIP values greater than one. These brain regions are in agreement with 469 the results reported in Prato et al. [26]. 470

The information gathered by the VIPs is complementary to the one 471 carried by the N-PLS weights. Fig. 13 reports the weights for the first 472 two N-PLS factors for the second mode (ROIs). Extreme positive  $\mathbf{w}_j$  473 values on F1 are those related to ROI with VIP values higher than 474 one. Inspection of the corresponding weights plot for Mode3 (volun-475 teers), i.e.  $\mathbf{w}_k$ , Fig. 14, shows that the different volunteers are distrib-476 uted along the first factor that differentiates e.g. volunteer #9 from 477 volunteer #10.



Fig. 13. Data set for Neuroscience. Scatter plot of NPLS weights for Mode2 (ROIs): first vs. second component.

S. Favilla et al. / Chemometrics and Intelligent Laboratory Systems xxx (2013) xxx-xxx



Fig. 14. Data set for Neuroscience. Scatter plot of NPLS weights for Mode3 (volunteers): first vs. second component.

This can be discussed by considering the different behaviors in per-479ceived pain as underlined by the first mode loadings plots (functional 480fMRI profiles), shown in Fig. 15. The first factor describes the average 481 pain profile, the second seems almost dedicated to a delayed maximum 482 peak and more persistent pain. High positive loading values on factor 483 one for Mode2, Fig. 14 (see for instance volunteer #7) represent an op-484 posite behavior with an anticipated maximum peak in comparison with 485 the average pain profile. The highest negative value (volunteer #9) with 486 the extreme negative position (with respect to the abscissa) identifies a 487 positive shift of the maximum pain perceived in comparison with the 488 mean profile (reference volunteer #3). This may be retrieved by direct 489 490observation of the psychophysical responses for these subjects in comparison with those subjects showing a profile close to average as sub-491 492 jects three, see Fig. 16.

493The second factor can be considered as a component that takes into account a sort of "prolonged activation due to tonic pain input" (see 494 Fig. 15, in gray dashed) and it is particularly dedicated to describe vol-495unteer #2 with its ample bell-shape of the pain perceived. In Fig. 16 496the comparison between the pain profiles of volunteer #2 and the refer-497498ence volunteer #3 is shown. High weight  $w_{k2}$  value (see Fig. 14) for this volunteer seems to be only related to its particular behavior that is also 499 responsible for the separation of the ROI regions in two groups with re-500spect to the weight values on the second mode (Fig. 13). 501

Thus, discussion of the weights plot is useful to recover the detailed information on specific subject behavior while the overall VIP values point to the most relevant ROIs whose activation is involved in pain perception and that are thus capable of reproducing the Y-psycho responses.

#### 5. Conclusions

The recent developments in feature selection methods for two-way 508 data have addressed the problem of increasing the performance of re- 509 gression models, such as PLS. Complex filter, wrapper or embedded 510 methods [11,20] improve predictor performance compared to simpler 511 variable ranking methods, but the improvements are not always signifi- 512 cant, they are computationally costly and in case of a large number of 513 variables, the risk of over-fitting can be not negligible in the process of 514 variable selection. The extension of variable selection methods to the 515 multi-way data arrays, in multilinear regression context (NPLS, NPLS- 516 DA) without recurring to unfolding, has been rather limited. 517

In this work, we introduced an extension of the Variable Importance in 518 Projection (VIP) parameter to multi-way arrays. The proposed method 519 has been tested on three different data sets where VIPs were discussed 520 comparatively with respect to NPLS weight and regression coefficients. 521

VIPs naturally point out the identification of the most relevant varibles related to **Y** in a multi-way array **X**. In particular, VIPs can be calculated for each mode of **X** and both considering the single y-responses or all the **Y** altogether. The former can be useful to assess the relevant variables with respect to each modeled properties, especially in the case of discriminant NPLS-DA to highlight the discriminant features, since each y-variable corresponds to a given class. While the latter offer a useful summary in order to operate feature selection. 529

However, when considering VIP, it is important to remember that, 530 from an interpretative point of view, this metric suffers from the fact 531 that PLS/NPLS components carry with them the unresolved contribution 532 of both Y-related and Y-orthogonal parts of the X-variance. Both contribu-533 tions are fundamental for a correct prediction of Y and the VIP metric rep-534 resents a valid support whenever the interest is aimed to rank variables 535 according to their influence to the whole model, while for contexts and 536 purposes where the interest is mainly focused on assessing the X-part co-537 varying with Y only, other methodologies may represent a valuable solution, such as OPLS [28] and Selectivity Ratio metric [21]. 539

In the studied cases, VIPs provided an easier and complementary 540 way to interpret the variable relevance in NPLS models, especially 541 when examination of regression coefficients was not so straightfor- 542 ward due to the unreadable complex patterns associated [27], as in 543 the case of spectral data (as for the Bread data set) and moreover 544 with two signal dimensions as in the case of hyphenated analytical 545 techniques. In the EVOO application, which is an example of chroma- 546 tography/mass spectrometry data, the joint information from the VIPs 547 on the retention time and m/z directions allows discussion in chemi- 548 cal terms of the most discriminant features.



Fig. 15. Data set for Neuroscience. Mode1 loading (score) plot vs. time points (the temporal mode showing the evolution of functional fMRI-BOLD): (black) first component; (gray dashed) second component.

Please cite this article as: S. Favilla, et al., Assessing feature relevance in NPLS models by VIP, Chemometrics and Intelligent Laboratory Systems (2013), http://dx.doi.org/10.1016/j.chemolab.2013.05.013

507

S. Favilla et al. / Chemometrics and Intelligent Laboratory Systems xxx (2013) xxx-xxx



Fig. 16. Data set for Neuroscience. Pain intensity profile (psychophysical response) vs. time points, for some reference volunteers: subject #3 (solid black), subject #2 (dashed), subject #7 (dot-dashed) and subject #9 (dotted).

NPLS weights carry the information about the relevance of variables
 and sign of correlation with the modeled responses. However, they re quire to be discussed together with the Y loadings and per component.

553 In this respect, the VIPs are complementary pointing to the most influ-

ential variables for each property on taking into account all components

555 but of course require inspection of weights to assess the direction of the

effect (increasing or decreasing the response values). Finally, the results

obtained on the Neuroscience dataset were found to be in line with those published with a completely different method belonging to the

those published with a completely different method belonging to the machine learning field [26] as far as ranking of the most influential

ROIs is concerned, while the use of multi-way models added the possi-

bility to discuss both the common pattern to all volunteers in pain per-

ception as well as the peculiar behavior of specific ones.

### Q10 563 6. Uncited reference

564 [23]

### 565 Acknowledgments

The "Functional Neuroimaging" group, of the Department of Biomedical Sciences, Metabolic and Neuroscience, University of Modena and Reggio Emilia are kindly acknowledged for providing the fMRI data and for the support given in the discussion and interpretation of the results on this data set.

#### 571 References

575

576

- Q11 572 [1] R. Bro, Multi-way calibration. Multi-linear PLS, Journal of Chemometrics 10 (1996) 573 47–61. 574 [2] R. Bro, A.K. Smilde, S. De Jong, On the difference between low-rank and subspace
  - [2] R. Bro, A.K. Smilde, S. De Jong, On the difference between low-rank and subspace approximation: improved model for multi-linear PLS regression, Chemometrics and Intelligent Laboratory Systems 58 (2001) 3–13.
  - [3] A. Smilde, Comments on multilinear PLS, Journal of Chemometrics 11 (1997)
     367–377.
  - 579 [4] S. De Jong, Regression coefficients in multilinear PLS, Journal of Chemometrics 12 (1998) 77–81.
  - [5] A. Smilde, R. Bro, P. Geladi, Multi-way analysis with applications, Multi-way Analysis in the Chemical Sciences, Chpt 3.1, John Wiley & Sons, 2004, pp. 35–45.
  - [6] E. Salvatore, M. Bevilacqua, R. Bro, F. Marini, M. Cocchi, Classification methods for multi-way arrays as a basic tool for food PDO authentication, in: M. De La Guardia, A. Gonzalvez Illueca (Eds.), Food Protected Designation of Origin: Methodologies & Applications, Wilson & Wilson's Comprehensive Analytical Chemistry, vol. 60, Elsevier B.V., 2013
  - [7] P. Nomikos, J.F. MacGregor, Multivariate SPC charts for monitoring batch processes, Technometrics 37 (1995) 41–59.
  - [8] C. Durante, M. Cocchi, M. Grandi, A. Marchetti, R. Bro, Application of N-PLS to gas chromatographic and sensory of traditional balsamic vinegars of Modena, Chemometrics and Intelligent Laboratory Systems 83 (2006) 54–65.

- [9] C.M. Rubingh, M.J. van Erk, S. Wopereis, T. van Vliet, E.R. Verheij, N.H.P. Cnubben, 593 B. van Ommen, J. van der Greef, H.F.J. Hendriks, A.K. Smilde, Discovery of subtle effects in a human intervention trial through multilevel modeling, Chemometrics 595 and Intelligent Laboratory Systems 106 (2011) 108–114. 596
- and Intelligent Laboratory Systems 106 (2011) 108–114. 596 10] A. Conesa, J.M. Prats-Montalbán, S. Tarazona, M.J. Nueda, A. Ferrer, A multiway 597 approach to data integration in systems biology based on Tucker3 and N-PLS, 598 Chemometrics and Intelligent Laboratory Systems 104 (2010) 101–111. 599
- Chemometrics and Intelligent Laboratory Systems 104 (2010) 101–111.
   599

   [11] T. Mehmood, K.H. Liland, L. Snipen, S. Sæbø, A review of variable selection
   600

   methods in partial least squares regression, Chemometrics and Intelligent Labo 601

   ratory Systems 118 (2012) 62–69.
   602
- [12] L. Norgaard, A. Saudland, J. Wagner, J.P. Nielsen, L. Munck, S.B. Engelsen, Interval par-603 tial least squares regression (iPLS): a comparative chemometric study with an ex-604 mple from part informed constructory (Applied Searcharcory), 54 (2000) 414, 410.
- ample from near-infrared spectroscopy, Applied Spectroscopy 54 (2000) 413–419. 605
   [13] R. Leardi, A.L. Gonzales, Genetic algorithms applied to feature selection in PLS re- 606
   gression: how to use them, Chemometrics and Intelligent Laboratory Systems 41 607
   (1998) 195–207. 608
- [14] L. Stordrange, T. Rajalahtia, F.O. Libnau, Multiway methods to explore and model 609 NIR data from a batch process, Chemometrics and Intelligent Laboratory Systems 610 70 (2004) 137–145. 611
- P.J. Odman, C. Lindvald, L. Olsson, et al., Sensor combination and chemometric 612 variable selection for online monitoring of *Streptomyces coelicolor* fed-batch culti-613 vations, Applied Microbiology and Biotechnology 86 (2010) 1745–1759.
   G. Lorho, F. Westad, R. Bro, Generalized correlation loadings. Extending correla-615
- [16] G. Lorho, F. Westad, R. Bro, Generalized correlation loadings. Extending correlation loadings to congruence and to multi-way models, Chemometrics and Intelligent Laboratory Systems 84 (2006) 119–125. 617
- [17] S. Wold, E. Johansson, M. Cocchi, PLS: partial least squares projections to latent 618 structures, in: Hugo Kubinyi (Ed.), 3D QSAR in Drug Design: Theory, Methods 619 and Applications, ESCOM Science Publishers, Leiden, ISBN: 90-72199-14-6, 620 1993, pp. 523–550. 621
- S. Wold, M. Sjöström, L. Eriksson, PLS-regression: a basic tool of chemometrics, 622 Chemometrics and Intelligent Laboratory Systems 58 (2001) 109–130.
   I.G. Chong, C.H. Jun, Performance of some variable selection methods when 624
- [19] I.G. Chong, C.H. Jun, Performance of some variable selection methods when 624 multicollinearity is present, Chemometrics and Intelligent Laboratory Systems 78 625 (2005) 103–112. 626
- [20] R. Gosselin, D. Rodrigue, C. Duchesne, A bootstrap-VIP approach for selecting 627 wavelength intervals in spectral imaging applications, Chemometrics and Intelli 628 gent Laboratory Systems 100 (2010) 12–21.
   [21] T. Rajalahti, R. Arneberg, F.S. Berven, K.M. Myhr, R.J. Ulvik, O.M. Kvalheim, Bio 630
- [21] T. Rajalahti, R. Arneberg, F.S. Berven, K.M. Myhr, R.J. Ulvik, O.M. Kvalheim, Bio- 630 marker discovery in mass spectral profiles by means of selectivity ratio plot, 631 Chemometrics and Intelligent Laboratory Systems 95 (2009) 35–48. 632
- [22] E. Acar, C. Aykut-Bingol, H. Bingol, R. Bro, A.L. Ritaccio, B. Yener, Modeling and de- 633 tection of epileptic seizures using multi-modal data construction and analysis, 634 Technical Report, 2008, (Computer Science Department at RPI (http://www.cs. 635 rpi.edu/research/tr.html)). 636
- [23] J. Nilsson, S. De Jong, A.K. Smilde, Multiway calibration in 3D QSAR, Journal of 637 Chemometrics 11 (1997) 511. 638
- M. Li Vigni, M. Cocchi, Near infrared spectroscopy and multivariate analysis to 639 evaluate wheat flour doughs leavening and bread properties, Analytica Chimica 640 Acta 764 (2013) 17–23.
   C. Durante, R. Bro, M. Cocchi, A classification tool for N-way array based on SIMCA 642
- [25] C. Durante, R. Bro, M. Cocchi, A classification tool for N-way array based on SIMCA '642' methodology, Chemometrics and Intelligent Laboratory Systems 106 (2011) 73–85. 643
- [26] M. Prato, S. Favilla, L. Zanni, CA. Porro, P. Baraldi, A regularisation algorithm for 644 decoding perceptual temporal profiles from fMRI data, NeuroImage 56-1 (2011) 645 258–267. 646
- [27] A.J. Burnham, J.F. MacGregor, R. Viveros, Interpretation of regression coefficients under 647 a latent variable regression model, Journal of Chemometrics 15 (2001) 265–284.
   [28] J. Trygg, S. Wold, Journal of Chemometrics 17 (2003) 53–64.
  - 649 650