



**HAL**  
open science

# Geographical discrimination of red garlic (*Allium sativum* L.) using fast and non-invasive Attenuated Total Reflectance-Fourier Transformed Infrared (ATR-FTIR) spectroscopy combined with chemometrics

Alessandra Biancolillo, F. Marini, A. d'Archivio

► **To cite this version:**

Alessandra Biancolillo, F. Marini, A. d'Archivio. Geographical discrimination of red garlic (*Allium sativum* L.) using fast and non-invasive Attenuated Total Reflectance-Fourier Transformed Infrared (ATR-FTIR) spectroscopy combined with chemometrics. *Journal of Food Composition and Analysis*, 2020, 86, pp.103351. 10.1016/j.jfca.2019.103351 . hal-02609770

**HAL Id: hal-02609770**

**<https://hal.inrae.fr/hal-02609770>**

Submitted on 21 Jul 2022

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

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

Copyright

**Geographical discrimination of red garlic (*Allium sativum* L.) using fast and non-invasive Attenuated Total Reflectance-Fourier Transformed Infrared (ATR-FTIR) spectroscopy combined with chemometrics**

*Alessandra Biancolillo,<sup>a,b\*</sup> Federico Marini<sup>a</sup>, Angelo Antonio D'Archivio<sup>c</sup>*

<sup>a</sup>Dipartimento di Chimica, Università degli Studi di Roma "La Sapienza" Piazzale Aldo Moro,  
00185, Roma, Italy

<sup>b</sup> ITAP, Irstea, Montpellier SupAgro, University of Montpellier, Montpellier, France

<sup>c</sup>Dipartimento di Scienze Fisiche e Chimiche, Università degli Studi dell'Aquila,  
Via Vetoio, 67010 Coppito, L'Aquila, Italy

\*Corresponding author:

dr. Alessandra Biancolillo

Dept. of Chemistry

University of Rome La Sapienza

P.le Aldo Moro 5

I-00185 Rome

Italy

Tel +39 06 49913680

Fax +39 06 4969 3292

e-mail: [alessandra.biancolillo@uniroma1.it](mailto:alessandra.biancolillo@uniroma1.it)

Declarations of interest: none

1 **Geographical discrimination of red garlic (*Allium sativum* L.) using fast and non-invasive**  
2 **Attenuated Total Reflectance-Fourier Transformed Infrared (ATR-FTIR) spectroscopy**  
3 **combined with chemometrics**

4 Declarations of interest: None

5

6 **Abstract**

7 Four varieties of red garlic (*Allium sativum* L.) cultivated in different Italian territories, Sulmona  
8 (Abruzzo), Proceno and Castelliri (Lazio), and Nubia (Sicily), were analysed by Attenuated Total  
9 Reflectance-Fourier Transformed Infrared (ATR-FTIR) spectroscopy. ATR-FTIR spectra of bulbils  
10 and bulbil tunics were separately acquired and processed by Partial Least Squares Discriminant  
11 Analysis (PLS-DA) with the aim of classifying the garlic samples on the basis of their geographical  
12 origin. Finally, two multi-block strategies (based on Sequential and Orthogonalized Partial Least  
13 Squares and Sequential and Orthogonalized Covariance Selection, coupled with Fisher's Linear  
14 Discriminant Analysis) have been applied in order to test whether a joint analysis of data could lead  
15 to higher prediction rates. Eventually, the best results were achieved by the multi-block approach  
16 based on SO-PLS, which allows obtaining a total classification rate of 95% (corresponding to one  
17 misclassified sample over 20) in external validation.

18

19 **Keywords:** *Garlic; ATR-FTIR; geographical classification; PLS-DA; Multi-block; Multi-block; SO-*  
20 *PLS; SO-CovSel.*

21

22 **1. Introduction**

23 Garlic (*Allium sativum* L.) has been worldwide employed as food condiment and herbal medicine  
24 for millennia. Apart from the common culinary use of fresh leaves or cloves, commercial products  
25 obtained by various processing methods, including oil maceration, dehydration and lyophilisation,  
26 are today marketed for therapeutic purposes (Ramirez et al., 2017). Historical references dating  
27 back to 4000 years mention diffuse use of garlic in ancient civilizations, from religious and  
28 superstition rituals to prevention and cure of infections and diseases (Corzo-Martínez et al., 2007).  
29 Bio-activity of garlic, including anti-inflammatory, antimicrobial, cardioprotective, anticancer and  
30 antidiabetic action, has been demonstrated in the last decades by epidemiological and clinical  
31 studies (Corzo-Martínez et al., 2007; Martins et al., 2016; Yun et al., 2014; Shukla & Kalra, 2007).

32 Unique pungent aroma and most of the medical properties attributed to garlic by traditional and  
33 modern medicine are related to distinctive organosulfur compounds.

34 Great attention is paid to the relation between the cultivar and geographical origin of garlic and its  
35 metabolomic profile, with specific reference to the aroma precursors and other bio-active  
36 constituents (Lu et al., 2011; Beato et al., 2011; Khar et al., 2011; Montaña et al., 2011). Moreover,  
37 traditional garlic varieties cultivated in given territories are appreciated by an increasing number of  
38 consumers because of their peculiar taste, aroma and functional properties compared to commercial  
39 products. In Europe, quality and geographical identity of some traditional garlic varieties of Italy  
40 (*Aglio di Voghiera* and *Aglio Bianco Polesano*), France (*Ail violet de Cadours*, *Ail fumé d'Arleux*,  
41 *Ail blanche de Lomagne* and *Ail de la Drôme*) and Spain (*Ajo Morado de Las Pedroñeras*) have  
42 been officially recognised in recent years through the attribution of PDO (Protected Designation of  
43 Origin) or PGI (Protected Geographical Indication) mark (European Commission, Agriculture and  
44 Rural Development, 2019). In particular, in the last years, different consortia in Italy have been  
45 constituted to valorize and preserve the traditional varieties of specific territories. Beside the  
46 protection provided by institutions (through laws and regulations) a wide effort has been put in  
47 developing analytical methodologies aimed at authenticating and tracing food specialties awarded  
48 of quality marks, for example (Lastra-Mejías et al., 2020; D'Archivio, et al. 2019a; Biancolillo et al.,  
49 2018a; Rocha et al., 2019; Giannetti et al., 2019; Mora et al., 2020; Firmani et al. 2019). In this  
50 context, analytical/chemometric approaches for the classification of garlic according to the cultivar  
51 and/or the geographical origin are essential tools to unveil commercial frauds arising from the  
52 intentional substitution of varieties cultivated in specific territories by commercial products.

53 Various analytical techniques, such as  $^1\text{H}$  high resolution magic angle spinning-nuclear magnetic  
54 resonance spectroscopy (Ritota et al., 2012), infrared spectroscopy (Lu et al., 2011), high  
55 performance liquid chromatography (Montaña et al., 2011), high resolution mass spectrometry  
56 (Hrbek et al., 2018) and electronic nose (Trirongjitmoah et al., 2015), were applied to characterise  
57 garlic for traceability purposes. In these investigations, the organo-sulphur compounds and other

58 metabolome components, such as amino acids, fatty acids, organic acids and sugars were  
59 recognised as promising traceability indicators to assess the garlic provenance or variety. In  
60 addition, garlic cultivated in different countries (Smith, 2005; Vasi et al., 2016) or close regions  
61 (D'Archivio et al., 2019b) were well discriminated using the trace multi-element profile determined  
62 by atomic spectroscopy. However, most of the above analytical methods are relatively complex,  
63 expensive, time-consuming and require specialized skills. Moreover, a preliminary sample  
64 treatment is often necessary, which, apart from further increasing complexity and cost of the  
65 characterisation method, may also alter the metabolomic profile of garlic. Attenuated Total  
66 Reflectance-Fourier Transformed Infrared (ATR-FTIR) spectroscopy, by contrast, is a relatively  
67 simple, fast, cheap and non-invasive technique applicable to both liquids and solids without any  
68 complex sample pre-treatment. About garlic, ATR-FTIR spectroscopy was previously used to  
69 quantify the total phenol content and antioxidant activity with the aim of differentiating the samples  
70 grown in different US states (Lu et al., 2011), but the ATR-FTIR spectra, rather than directly on the  
71 garlic samples, were acquired from methanolic extracts.

72 **In the light of these considerations, the aim of the present work is to test the potentiality of ATR-**  
73 **FTIR for geographical traceability purposes. In particular, red garlic varieties cultivated in four**  
74 **distinct areas of Italy, namely Sulmona (Abruzzo), Castelliri and Proceno (Lazio), and Nubia**  
75 **(Sicily), were analyzed and the observed spectra were handled by chemometrics. These red garlic**  
76 **ecotypes have been chosen because of their valuable characteristics; additionally,** those cultivated in  
77 Sulmona and Nubia were also included by Slow Food Foundation for Biodiversity (Slow Food  
78 Foundation, 2019) in the list of local plant varieties to safeguard. At first, ATR-FTIR spectra of  
79 bulbils and tunics were separately analysed by Partial Least Squares-Discriminant Analysis (PLS-  
80 DA) (Sjöström et al., 1986); this approach is widely and satisfactorily applied for authentication of  
81 agro-food, in particular handling **Infrared (IR)** data (Biancolillo & Marini, 2018b). Finally, the two  
82 data blocks were jointly analyzed by multi-block classifiers, in order to test whether data fusion  
83 strategies would provide more accurate models, allowing a deeper comprehension of the system.

84 Consequently, IR spectra were analyzed by Sequential and Orthogonalized Partial Least Square  
85 (SO-PLS) (Næs et al., 2011) or Sequential and Orthogonalized Covariance Selection (SO-CovSel)  
86 (Biancolillo et al., 2019a) coupled with Fisher's Linear Discriminant Analysis (LDA).

87

## 88 **2. Materials and methods**

### 89 *2.1. Garlic samples*

90 Bulbs of red garlic varieties cultivated in 2017 in four Italian sites, Sulmona (Abruzzo), Castelliri  
91 and Proceno (Lazio), and Nubia (Sicily), were kindly donated by producers working in the  
92 respective territories that assured the geographical origin of the samples. **In particular, 81 garlic**  
93 **cloves were analyzed with the tunics; of these, 20 were from Castelliri, 20 from Proceno, 20 from**  
94 **Nubia and 21 from Sulmona. Eventually, spectra were collected on 82 skinned samples (21 were**  
95 **from Castelliri, 20 from Proceno, 19 from Nubia and 22 from Sulmona). Only 69 cloves were**  
96 **analyzed on both compartments (i.e., either with and without tunic). Of these, 19 were from**  
97 **Castelliri, 20 from Proceno, 17 from Nubia and 13 from Sulmona.** Samples were acquired in July-  
98 September 2017, stored under typical domestic conditions (**bulbils were located into a box which**  
99 **allowed their perspiration, not exposed under direct light and in a cool room**) and analyzed before  
100 December to avoid variations in the composition due to aging or sprouting.

101

### 102 *2.2 ATR FT-IR measurements*

103 The infrared spectra of garlic cloves and clove tunics were separately recorded on a PerkinElmer  
104 Spectrum Two™ (PerkinElmer, Waltham MA, USA) FT-IR spectrometer consisting in a deuterated  
105 triglycine sulfate (DTGS) detector and a PerkinElmer Universal Attenuated Total Reflectance  
106 (uATR) accessory equipped with a single bounce diamond crystal. Each spectrum was registered  
107 from 4000 cm<sup>-1</sup> to 400 cm<sup>-1</sup> with 1 cm<sup>-1</sup> instrumental resolution and ten scans were averaged per  
108 spectral replicate. The background was collected with the crystal exposed to the air. Before each  
109 measurement, the ATR crystal was cleaned with methanol and air-dried. ATR FT-IR spectra were

110 collected on intact freshly peeled cloves by contacting a flat part of the clove with the ATR crystal.  
111 A consistent force was applied using the pressure monitoring system integrated with the instrument  
112 to maximize the spectrum intensity but avoiding crushing the clove. The ATR FT-IR spectra of  
113 clove tunics were recorded separately following the same procedure.

114 -----Insert Figure 1 approx. Here-----

115 The observed ATR-FTIR spectra (shown in Figure 1) display the typical vibration patterns of the  
116 plant constituents (proteins, fats and sugars) (Schulz & Baranska, 2007; Movasaghi et al., 2008) and  
117 reflect the composition of garlic clove and tunics. The garlic clove is mainly composed by water  
118 (65%), followed by carbohydrates (28%, mainly fructans), sulphur compounds (1-4%), proteins  
119 (2%), fibres (1.5%) and free amino acids (1-1.5%) (Rahman, 2003), while polysaccharides  
120 (cellulose, hemicellulose and pectin) and lignin are the main constituents of the bulbil skin (Kallel  
121 et al., 2015; Reddy & Rhim, 2014). The broad band centred at about  $3290\text{ cm}^{-1}$  can be assigned to  
122 the N-H stretching of proteins and O-H stretching of carbohydrates and water, while the two sharp  
123 signals at  $2920$  and  $2850\text{ cm}^{-1}$  are associated to symmetric and antisymmetric C-H stretching  
124 vibrations, respectively. The spectral region between  $1200$  and  $900\text{ cm}^{-1}$ , although showing  
125 different intensity and fine structure in the spectra of cloves and skins, takes origin from coupled C-  
126 C, C-O stretching and C-O-H, C-O-C deformation modes of oligo- and polysaccharides. The  
127 distinctive band at  $1025\text{ cm}^{-1}$  in particular can be assigned to the vibrational frequency of  $\text{CH}_2\text{OH}$   
128 groups of carbohydrates. In the same spectral region, the S=O stretching of sulfoxides may  
129 contribute to the signal at about  $1090\text{ cm}^{-1}$  (Nikolić et al., 2011), well visible in cloves but not in  
130 skin samples. The band at about  $1160\text{ cm}^{-1}$  arises from the glycosidic linkage (C-O-C) vibrations.  
131 The weak bands in the region  $880-900\text{ cm}^{-1}$  can be attributed to the C-O-C skeletal modes of  
132 carbohydrates and polysaccharides, the signal at  $894\text{ cm}^{-1}$  in particular being diagnostic of  $\alpha$ -(1 →  
133 4)-glycosidic bonds. The band at  $1735-1740\text{ cm}^{-1}$  observed in both clove and skin samples are  
134 typical of the C=O stretching vibration of polysaccharides and cellulose. The signal at  $1640\text{ cm}^{-1}$   
135 observed in the spectrum of cloves can be assigned to the bending vibration of water and to the

136 stretching of carbonyl of proteins (amide I), while the signals at  $1552\text{ cm}^{-1}$  and  $1252\text{ cm}^{-1}$  can be  
137 attributed to the amide II and amide III bands (associated with coupled C–N stretching and N–H  
138 bending vibrations of the peptide group). Additionally, the stretching vibrations of aliphatic and  
139 aromatic double bonds fall in the region  $1640\text{--}1500\text{ cm}^{-1}$  as well. The signal at  $1225\text{--}1230\text{ cm}^{-1}$  of  
140 the skin spectra can be attributed to the stretching vibrations of C–O bonds in lignin (Stark et al.,  
141 2016). The band observed at about  $1600\text{ cm}^{-1}$  in skin samples, partially overlapped to the amide I  
142 band of proteins in the clove spectra, can be assigned to the asymmetric stretching of carboxylate  
143 groups of amino acids, proteins or polysaccharides. The signals due to the O–H stretching of  
144 adsorbed water and asymmetric stretching of lignin aryl rings also fall in this spectral region, while  
145 the band at about  $1510\text{ cm}^{-1}$  can be assigned to the symmetric stretch of the aromatic groups of  
146 lignin (Stark et al., 2016). The two sharp signals at  $1460$  and  $1470\text{ cm}^{-1}$  observed in clove samples  
147 and sometimes superimposed, and the bands at  $1427\text{ cm}^{-1}$  and  $1328\text{ cm}^{-1}$  are due to O–H  
148 deformation vibration, and to various vibrational modes of  $\text{CH}_2$  groups of lipids, polysaccharides  
149 and proteins. The two weak signals at about  $720$  and  $727\text{ cm}^{-1}$  observed in the clove spectra but not  
150 in the skin samples can be attributed to C–S stretching vibrations of di-alkyl sulphides and  
151 disulphides (Minzhen et al., 2015) and to C–H deformation vibration.

152

## 153 2.3 Chemometric Analysis

154 The ATR-FTIR signals collected as described in the previous section have been analyzed by means  
155 of chemometric tools, in order to classify samples according to their geographical origin. To  
156 achieve this goal, three classification methods have been employed: PLS-DA, in order to handle the  
157 two data blocks individually, and SO-PLS-LDA or SO-CovSel-LDA, to achieve a simultaneous  
158 analysis of both sets of signals.

159

### 160 2.3.1 Partial Least Square-Discriminant Analysis (PLS-DA)

161 Discriminant classification methods discern samples on the basis of their mutual differences.  
162 Applying these approaches, the multi-dimensional samples-space is entirely divided into class-  
163 regions, and each object will be assigned to one specific category. One of the first discriminant  
164 classifiers proposed is the Linear Discriminant Analysis by Fisher. This approach, despite it  
165 performs well and it is still widely used, presents a considerable limitation: it can be used only when  
166 the data matrix is invertible. This condition is rarely met, in particular working with instrumental  
167 data, where the number of variables is likely higher than the number of samples.

168 Among the different methods developed in order to enable the application of discriminant analysis  
169 on ill-conditioned data matrices, Partial Least Square-Discriminant Analysis (PLS-DA) (Sjöström et  
170 al., 1986; Ståhle & Wold, 1987) is probably one of the most widely applied. One of the main  
171 reasons for its diffusion is that it is suitable to handle highly correlated variables (e.g., spectroscopic  
172 data), which makes it applicable on non-invertible data blocks. PLS-DA exploits the PLS algorithm  
173 (Wold et al., 1983) to solve a classification problem as if it were a regression one (Barker, &  
174 Rayens, 2003); mathematically, this corresponds to estimate Eq.1:

$$175 \quad \mathbf{Y} = \mathbf{X}\mathbf{b} + \mathbf{e} \quad (1)$$

176 Where  $\mathbf{X}$  is the data matrix of measures collected on samples,  $\mathbf{b}$  and  $\mathbf{e}$  are the regression  
177 coefficients and the residuals, respectively, and  $\mathbf{Y}$  is the so-called *Dummy Matrix*, a binary matrix  
178 (of dimensions  $N \times G$ , where  $N$  is the number of the analyzed samples and  $G$  is the number of  
179 categories present into the system) encoding the class-belonging. Once Eq.1 is solved (i.e., once the  
180 calibration model is built), whenever new measures ( $\mathbf{X}_{new}$ ) are collected on unknown samples, it is  
181 possible to predict their class-belonging (solving  $\hat{\mathbf{Y}}_{new} = \mathbf{X}_{new}\mathbf{b}$ ). Nevertheless,  $\hat{\mathbf{Y}}_{new}$ , i.e., the matrix  
182 collecting the responses predicted on new samples, is not categorical, but made of continuous  
183 values, so that it is necessary to define a criterion for class-attribution. Accordingly, classification  
184 may be accomplished, e.g., by assigning the sample to the category corresponding to the highest  
185 value of the predicted response, by application of LDA either on scores or on predicted responses or  
186 by a bayesian approach based on Gaussian mixture modeling (Perez et al., 2009).

187

### 188 **2.3.2 Sequential and Orthogonalized Partial Least Square-Linear Discriminant Analysis (SO-** 189 **PLS-LDA)**

190 Sequential and Orthogonalized Partial Least Square-Linear Discriminant Analysis (SO-PLS-LDA)  
191 (Biancolillo et al., 2015) is a multi-block discriminant classifier whose algorithm has been  
192 developed combining a multi-block regression method, SO-PLS (Næs et al., 2011), and Fisher's  
193 Linear Discriminant Analysis, where SO-PLS is used to reduce the dimensionality of the data  
194 blocks (solving problems related to ill-conditioned matrices) prior the application of the  
195 discriminant approach.

196 In order to create a SO-PLS-LDA model, first of all it is necessary to build the SO-PLS one; given  
197 two predictor blocks,  $\mathbf{X}$  and  $\mathbf{Z}$ , and a dummy  $\mathbf{Y}$ , this can be done applying the procedure described  
198 in (Næs et al., 2011). Briefly,  $\mathbf{Y}$  is fitted to  $\mathbf{X}$  by PLS (obtaining the  $\mathbf{X}$ -scores  $\mathbf{T}_X$ ),  $\mathbf{Z}$  is  
199 orthogonalized with respect to  $\mathbf{T}_X$  and then the resulting matrix ( $\mathbf{Z}_{Orth}$ ) is used to predict (by PLS)  
200 the residuals from the previous regression. Finally, the predictive model is calculated summing up  
201 the outcomes of the two PLS-models (the reader is addressed to (Næs et al. 2011) for more details).  
202 Once the SO-PLS model is built, classification is achieved by applying LDA on the predicted  $\mathbf{Y}$  (or  
203 on the row-augmented scores) (Biancolillo & Næs, 2019b).

204

### 205 **2.3.3 Sequential and Orthogonalized-Covariance Selection-Linear Discriminant Analysis (SO-** 206 **CovSel-LDA)**

207 As the name suggests, Sequential and Orthogonalized-Covariance Selection-Linear Discriminant  
208 Analysis (SO-CovSel-LDA) (Biancolillo et al., 2019a) is a multi-block classification method  
209 strictly linked to SO-PLS-LDA. In fact, the two approaches have similar algorithm, but in SO-  
210 CovSel-LDA the feature reduction is operated by a variable selection method called Covariance  
211 Selection (CovSel) (Roger et al., 2011). CovSel is a feature reduction approach developed to select  
212 variables in a regression context; in fact, it points out the predictors which contribute the most to the

213 estimation of a response. Briefly, considering a predictor block  $\mathbf{X}$ , used to estimate a response  $\mathbf{Y}$ ,  
214 Covariance Selection iteratively selects the  $\mathbf{X}$ -variables presenting the highest covariance with the  
215 response. As a consequence, the main divergence between SO-PLS and SO-CovSel, is that, in the  
216 latter, calculations are based on the original variables instead on the scores (Biancolillo et al.,  
217 2019a). Considering the two predictor blocks case above-mentioned ( $\mathbf{X}$  and  $\mathbf{Z}$ ), in order to create a  
218 SO-CovSel model, the first step consists in selecting variables from the  $\mathbf{X}$ -block (by CovSel),  
219 obtaining the reduced matrix  $\mathbf{X}_{red}$ , which is used to estimate the  $\mathbf{Y}$  by ordinary least squares (OLS).  
220 Then,  $\mathbf{Z}$  is orthogonalized with respect to  $\mathbf{X}_{red}$ , obtaining  $\mathbf{Z}_{Orth}$ . Covariance Selection is then used  
221 to select the  $\mathbf{Z}_{Orth}$ -variables which contribute the most to the prediction of the residuals ( $\mathbf{E}_Y$ ) from  
222 the previous regression, obtaining  $\mathbf{Z}_{Orth,red}$ , which is used to estimate  $\mathbf{E}_Y$  by OLS. Finally, the  
223 predicted  $\mathbf{Y}$  is calculated by summing up the individual predictions made by the two regressions.  
224 Once the SO-CovSel model is created, LDA can be applied on the predicted  $\mathbf{Y}$ . Also in this case, in  
225 order to solve the classification problem, the  $\mathbf{Y}$  is a binary matrix encoding class-belongings.

226

### 227 **3. Results and discussion**

228 After the collection of spectra, IR signals were exported in MatLab (The Mathworks, Natick, MA;  
229 version 2015b) for the analysis. In order to pursue external validation of the models, spectra are  
230 reorganized into a training set, for the optimization of the calibration model, and a test set, for  
231 validation, by the Duplex algorithm (Snee, 1977) (more details about the division are reported  
232 below).

233 In Figure 1 the average spectra collected on cloves (in red) and on tunics (in blue) are shown. From  
234 the figure it is straightforward the IR signals collected on the different compartments of garlic are  
235 slightly different, in particular in specific absorption ranges; a wider discussion over the  
236 interpretation of the spectra is reported in the related sections.

237

238 Independently on the classifier used, different pretreatments have been tested on training spectra:  
239 bare mean centering, 1<sup>st</sup> or 2<sup>nd</sup> derivative (following the Savitzky-Golay approach, using 19 points  
240 window and a second or third order for the interpolating polynomial, respectively) (Savitzky &  
241 Golay, 1964), Standard Normal Variate (SNV) (Barnes et al., 1989) and combinations of SNV and  
242 derivatives; the most suitable preprocessing approach (together with the optimal complexity, i.e.,  
243 the number of latent variables) has been defined as the one leading to the lowest classification error  
244 in a 7-fold cross-validation procedure. Even if not always explicitly mentioned, data is assumed to  
245 be mean-centered prior the creation of any model.

246

### 247 **3.1 PLS-DA analysis on tunics**

248 As above-mentioned, the IR signals were divided into a calibration and a validation set. The 81  
249 signals collected on tunics were divided into a training set of 61 samples (15 belonging to Class  
250 Castelliri, 15 pertaining to Class Proceno, 15 for Class Nubia and 16 from Sulmona) and a test set  
251 of 20 samples (5 per each category). Then, six different pretreatments (listed in Table 1) have been  
252 tested on data, and an equal number of PLS-DA models were calculated (on training samples); the  
253 preprocessing approaches, the number of latent variables (LVs) extracted and the average cross-  
254 validated classification errors (%) are reported in Table 1.

255 -----Insert Table 1 approx. here-----

256

257 After inspection of Table 1, the optimal calibration model has been built on data pretreated by SNV,  
258 1<sup>st</sup> derivative and mean centering; this PLS-DA model, applied on the test set (preprocessed  
259 accordingly) provided a classification rate of 75% (corresponding to 5 misclassified test samples  
260 over 20; of these, 2 belong to class Castelliri, 2 to class Sulmona and 1 appertains to class Nubia).

261 After the creation of any PLS-based model, in order to give a depth insight into the data set under  
262 study, it is possible to calculate the **Variable Importance in Projection (VIP)** indices (Wold et al.,

263 1993) to understand which variables contribute the most to the model; generally, each spectral  
264 variable presenting a VIP index higher than 1 is considered relevant. A graphical representation of  
265 VIP analysis is shown in Figure 2; one plot per class is displayed in order to show which variables  
266 characterize each category.

267 -----Insert Figure 2 approx. here-----

268

269

270 In Figure 2, the black solid lines represent the average training signals (offset to avoid overlapping  
271 and make them visible) whereas the selected variables are highlighted as bold colored dashed lines:  
272 red for Class Castelliri, blue for Class Proceno, green for Class Nubia and cyan for Class Sulmona.

273 From the figure it is possible to see that, as expected, the most relevant instrumental features are  
274 more or less the same among the four categories; in fact, independently on the class, variables  
275 presenting VIP indices higher than 1 are those around  $2849\text{ cm}^{-1}$  and  $2920\text{ cm}^{-1}$ , ascribable to the  
276 symmetric and antisymmetric stretching of the C-H bond, those in the spectral range between  $1640$   
277  $\text{cm}^{-1}$  and  $1629\text{ cm}^{-1}$ , attributable to the stretching of carboxylate groups, variables around  $1417\text{cm}^{-1}$   
278 ascribable to CH bending and some from  $1181\text{ cm}^{-1}$  to  $914\text{ cm}^{-1}$ , linked to the absorptions caused by  
279 CH deformations, skeletal stretching of C-O and C-C.

280 Despite the similarities clearly visible in the plot, it is also possible to spot some variables that are  
281 selected for the characterization of some categories, but not for others. For instance, spectral  
282 variables from  $1192\text{ cm}^{-1}$  to  $1292\text{ cm}^{-1}$  present VIP indices higher than 1 for all the categories  
283 except for Class Castelliri, probably indicating a different composition in oligo- and  
284 polysaccharides among the diverse categories. More details about spectral absorptions can be found  
285 in Section 2.2 and in the related literature.

286

287 **3.2 PLS-DA analysis on cloves**

288 The PLS-DA analysis of spectra collected on cloves has been carried out in the same way as  
289 described above for tunics. The 82 signals were divided into a training set of 62 objects (containing  
290 16 samples per Class Castelliri, 15 Proceno, 14 Nubia, and 17 per Class Sulmona) and a test set of  
291 20 samples (5 per category). The same pretreatments have been tested on data, and, also in this  
292 case, the most suitable has been defined (together with the optimal complexity) in cross-validation;  
293 results are reported in Table 2.

294 -----Insert Table 2 approx. here-----

295

296 In this case, the model leading to the lowest classification error in cross-validation is the one built  
297 on data after 1<sup>st</sup> derivative (and mean centering). When the optimal model was applied to the test  
298 set, only 3 over 20 samples were misclassified (correct classification rate: 85%); among these, 2  
299 belong to Class Castelliri and 1 to Class Sulmona.

300 Also in this case the VIP analysis was pursued in order to inspect the variables contributing the  
301 most to the observed differentiation among the geographical origin; the agreement regarding the  
302 selected variables among the different classes was strong (plot is not shown); the most relevant  
303 variables from the classification point of view were those from 2944 cm<sup>-1</sup> to 2837 cm<sup>-1</sup>, and some  
304 around 1413 cm<sup>-1</sup>, and 1757 cm<sup>-1</sup>.

305 Despite the results obtained by the individual analysis of data blocks were quite satisfactory, the  
306 above-mentioned multi-block approaches have been used, testing whether it would be possible to  
307 improve predictions.

308

309

### 310 **3.3. Multi-block analysis**

311 Unfortunately, the IR analysis of both tunics and cloves was not available for all the samples  
312 discussed for the PLS-DA analysis; consequently, the multi-block data set has been reduced to the  
313 69 samples which have been analyzed on both compartments.

314 In order to divide samples into a training and a test set taking into account both blocks of measures,  
315 two PCA models have been calculated, one per each set of data. Then, the first 5 principal  
316 components extracted by each PCA model were row-wise concatenated; finally, samples were  
317 divided into training and test set by the Duplex algorithm (Snee, 1977) calculating sample distances  
318 in the scores-space defined by the PCs.

319 Due to the reduction of the available samples, the training set included 49 samples (14 belonging to  
320 Class Proceno, 15 from Class Castelliri, 12 appertaining to Class Nubia and 8 of Class Sulmona),  
321 and the test set was made of 20 objects (5 objects per class).

322 As anticipated, two multi-blocks classifiers have been applied at this stage of the work;  
323 independently of the approach used, spectra collected on cloves have been used as first input block,  
324 while signals collected on tunics have been used as the second one.

325 Classification models have been built also using an inverted input order (tunics-data modelled as  
326 first block and clove-data as second) but results were slightly worse; consequently, these will not be  
327 discussed in the following sections.

328

#### 329 **3.3.1 SO-PLS-LDA analysis**

330 Building the SO-PLS models, all the possible combinations of the above-mentioned preprocessing  
331 approaches have been tested on the two data blocks in a cross-validation procedure (7 cancellation  
332 groups). The optimal calibration model is the one calculated using the mean-centered spectra  
333 collected on cloves as first input block and signals on tunics as the second one (pretreated by 2<sup>nd</sup>  
334 derivative and mean centering). The number of components extracted from the two blocks are 4 and  
335 11 for spectra on cloves and tunics, respectively. This model has been applied on the test set

336 (pretreated accordingly) and it provided 100% of correct classification for all categories except for  
337 class Sulmona, whose correct classification rate was 80% (corresponding to 1 misclassified  
338 sample). The results are graphically shown in Figure 3, where samples are projected onto the space  
339 of the first two canonical variates.

340 -----Insert Figure 3 approx. here-----

341

342 Looking at the plot, it is possible to recognize a quite clear distinction among the four different  
343 classes. In particular, the first canonical variate allows discriminating samples belonging to class  
344 Castelliri (red circles, at negative values) from the other three categories; while the second  
345 canonical variate allows distinguishing the objects belonging to class Proceno (blue squares) and  
346 class Sulmona (green triangle), at negative scores, from those belonging to class Castelliri (red  
347 circles) and Nubia (black diamonds) at positive values. From the figure, it is easy to spot the  
348 misclassified sample from class Sulmona (green triangle), in fact, this falls closer to the centroid of  
349 class Nubia rather than to the one of its own category.

350 VIP analysis was pursued on the SO-PLS model, following the embedded strategy described in  
351 (Biancolillo et al., 2016); nevertheless, the results were not relevantly different from those  
352 previously described for the individual PLS-DA analysis, and therefore they are not reported.

353

### 354 **3.3.2 SO-CovSel-LDA analysis**

355 Similarly to SO-PLS-LDA analysis, also for SO-CovSel-LDA, several multi-block models have  
356 been built (in a cross-validation procedure) in order to test different combinations of pretreatments;  
357 simultaneously, also the number of variables to be selected per each block is chosen. The optimal  
358 model has been calculated on the clove-block pretreated by 1<sup>st</sup> derivative whereas the signals  
359 collected on tunics were preprocessed by SNV and 1<sup>st</sup> derivative. The number of selected variables  
360 is 1 and 7 on cloves- and tunics-block, respectively. The application of the calibration model to the

361 test set led to a correct classification rate of 85%, corresponding to 3 misclassified samples in total  
362 (2 object from Class Castelliri and 1 belonging to Class Sulmona assigned to Class Proceno).

363 As described in the Section 2.3.3, SO-CovSel-LDA naturally provides information about the  
364 variables contributing the most to the classification. A visual representation of the selected variables  
365 is reported in Figure 4, in particular, the mean spectra collected on cloves and tunics are reported in  
366 in Figure 4a and Figure 4b, respectively; selected variables are highlighted by red circles.

367 -----Insert Figure 4 approx. Here-----

368

369 As expected, the selection provided by CovSel is sharper **than** the one achieved by VIP analysis;  
370 nevertheless, the two are in agreement. In fact, CovSel selects the variable at 2843  $\text{cm}^{-1}$  in the  
371 spectra collected on cloves, and those at 2917  $\text{cm}^{-1}$  and 3285  $\text{cm}^{-1}$  on tunics, probably associated to  
372  $\text{CH}_3$ ,  $\text{CH}_2$ , O-H and N-H stretching; variables at 1013  $\text{cm}^{-1}$ , 1034  $\text{cm}^{-1}$ , 1588  $\text{cm}^{-1}$ , 400  $\text{cm}^{-1}$  and 945  
373  $\text{cm}^{-1}$  associable to polysaccharides (for more details the reader is addressed to Section 2.2 and to the  
374 related literature).

375

#### 376 **4. General overview of the results**

377 In general, all the classification models provided acceptable results, indicating ATR FT-IR coupled  
378 with discriminant classifiers could represent a suitable approach for assessing the geographical  
379 origin of the investigated cultivars of red garlic. The best results, from the prediction point of view,  
380 are provided by a multi-block approach; this outcome is somehow expected, because data fusion  
381 strategies are supposed to provide comparable or better results than models built on the individual  
382 data blocks (Biancolillo, et al. 2019c). In order to ease a comprehensive overview of the  
383 classification rates provided by the different models, they are reported all together in Table 3.

384 -----Inset Table 3 approx. here-----

385 From the table, it is straightforward the most suitable methodology to solve the classification  
386 problem under study is SO-PLS-LDA. Concerning the single-block analysis, the best results are

387 provided by the model built on data collected on the cloves. This latter achievement suggests this  
388 compartments contains more information suitable for distinguishing the different red garlic  
389 ecotypes.

390

## 391 **5. Conclusions**

392 The aim of the present work was to develop a non-destructive approach suitable for distinguishing  
393 different cultivars of red garlic according to their geographical origin. In order to achieve this goal,  
394 samples of red garlic harvested in four different Italian towns (Castelliri, Proceno, Nubia and  
395 Sulmona) were analyzed by ATR-FTIR spectroscopy and classified. Spectra collection was pursued  
396 on both tunics of bulbils and cloves, avoiding any other physical-chemical pretreatment of samples.  
397 The data-block obtained were individually analyzed by PLS-DA and involved in multi-block  
398 models by the application of SO-PLS-LDA and SO-CovSel-LDA. In general, all the approaches  
399 provided good results from the prediction point of view. Concerning the classification pursued on the  
400 individual data blocks, the lowest classification error was provided by the PLS-DA model  
401 calculated on spectra collected on cloves, which led to the misclassification of three test objects  
402 over twenty. Nevertheless, the best results have been provided by a data-fusion strategy, the SO-  
403 PLS-LDA approach, which allowed achieving extremely satisfactory results, misclassifying only  
404 one sample over the 20 constituting the validation set.

405

406

## 407 **References**

408 Barker, M., & Rayens, W. (2003). Partial least squares for discrimination. *Journal of*  
409 *Chemometrics*, 17, 166–173. <https://doi.org/10.1002/cem.785>

410

411 Barnes, R. J., Dhanoa, M. S., & Lister, S. J. (1989). Standard normal variate transformation and de-  
412 trending of near-infrared diffuse reflectance spectra. *Applied Spectroscopy*, 43, 772-777.  
413 <https://doi.org/10.1366/0003702894202201>

414

415 Beato, V.M., Orgaz, F., Mansilla, F., & Montaña, A. (2011). Changes in Phenolic Compounds in  
416 Garlic (*Allium sativum* L.) Owing to the Cultivar and Location of Growth, *Plant Foods for Human*  
417 *Nutrition*, 66, 218-223 <https://doi.org/10.1007/s11130-011-0236-2>.

418

419 Biancolillo, A., Måge, I., & Næs, T. (2015). Combining SO-PLS and linear discriminant analysis  
420 for multi-block classification. *Chemometrics and Intelligent Laboratory Systems*, 141, 58–67.  
421 <https://doi.org/10.1016/j.chemolab.2014.12.001>

422

423 Biancolillo, A., Hovde Liland, K., Måge, I., Næs, T., & Bro, R. (2016). Variable selection in multi-  
424 block regression. *Chemometrics and Intelligent Laboratory Systems* 156, 89–101.  
425 <https://doi.org/10.1016/j.chemolab.2016.05.016>

426

427 Biancolillo, A., De Luca, S., Bassi, S., Roudier, L., Bucci, R., Magrì, A.D., & Marini, F.  
428 (2018a). Authentication of an Italian PDO hazelnut (“Nocciola Romana”) by NIR spectroscopy,  
429 *Environmental Science and Pollution Research* 25 28780-28786.

430

431 Biancolillo, A., & Marini, F. (2018b). Chemometrics Applied to Plant Spectral Analysis. In J.  
432 Lopes, & C. Sousa (Eds.), *Vibrational Spectroscopy for Plant Varieties and Cultivars*  
433 *Characterization, Comprehensive Analytical Chemistry*, 80 (pp. 69-104). Amsterdam: Elsevier.

434

435 Biancolillo, A., Marini, F., & Roger, J-M. (2019a). SO-COVSEL: a novel method for variable  
436 selection in a multi-block framework, *Journal of Chemometrics*, e3120,  
437 <https://doi.org/10.1002/cem.3120>

438

439 Biancolillo, A., & Naes, T. (2019b). The sequential and orthogonalised PLS regression (SO-PLS)  
440 for multi-block regression: theory, examples and extensions. In M. Cocchi (Ed.), *Data Handling in*  
441 *Science and Technology*, 31 (pp. 157-177). Amsterdam: Elsevier.

442

443 Biancolillo, A., Boqué, R., Cocchi, M., & Marini, F. (2019c). Data Fusion strategies in food  
444 analysis. In M. Cocchi (Ed.), *Data Handling in Science and Technology*, 31 (pp. 271-310).

445 Amsterdam: Elsevier.

446

447 D'Archivio, A.A., Di Vacri, M.L., Ferrante, M., Maggi, M.A., Nisi, S., & Ruggieri, F. (2019a).  
448 Geographical discrimination of saffron (*Crocus sativus* L.) using ICP-MS elemental data and class  
449 modeling of PDO Zafferano dell'Aquila produced in Abruzzo (Italy). *Food Analytical Methods*, 1-  
450 10.

451

452 D'Archivio, A.A., Foschi, M., Aloia, R., Maggi, M.A., Rossi, L. & Ruggieri, F. (2019b).  
453 Geographical discrimination of red garlic (*Allium sativum* L.) produced in Italy by means of  
454 multivariate statistical analysis of ICP-OES data. *Food Chemistry*, 275, 333–338.  
455 <https://10.1016/j.foodchem.2018.09.088>.

456

457 European Commission, Agriculture and Rural Development. DOOR database,  
458 [ec.europa.eu/agriculture/quality/door/list.html?locale=en](http://ec.europa.eu/agriculture/quality/door/list.html?locale=en) Accessed 4 April 2019.

459

460

461 Firmani, P., Bucci, R., Marini, F., Biancolillo, A. (2019). Authentication of “Avola almonds” by  
462 near infrared (NIR) spectroscopy and chemometrics. *Journal of Food Composition and Analysis*, 82  
463 2019, 103235.

464

465 Giannetti, V., Boccacci Mariani, M., Marini, F., Torrelli, P., Biancolillo, A. (2019). Flavour  
466 fingerprint for the differentiation of Grappa from other Italian distillates by GC-MS and  
467 chemometrics. *Food Control*, 105 123–130.

468

469 Hrbek, V., Rektorisova, M., Chmelarova, H., Ovesna, J., & Hajslova, J. (2018) Authenticity  
470 assessment of garlic using a metabolomic approach based on high resolution mass spectrometry,  
471 *Journal of Food Composition and Analysis*, 67, 19–28. <https://doi.org/10.1016/j.jfca.2017.12.020>.

472

473 Kallel, F., Driss, D., Bouaziz, F., Belghith, L., Zouari-Ellouzi, S., Chaari, F., Haddar, A.,  
474 Chaabouni, S.E., & Ghorbel, R. (2015). Polysaccharide from garlic straw: Extraction, structural  
475 data, biological properties and application to beef meat preservation, *RSC Advances*, 5, 6728–6741.  
476 <https://10.1039/c4ra11045e>.

477

478 Khar, A., Banerjee, K. Jadhav, M.R., & Lawande, K.E. (2011). Evaluation of garlic ecotypes for  
479 alliin and other allyl thiosulphinates, *Food Chemistry*, *128*, 988–996.  
480 <https://doi.org/10.1016/j.foodchem.2011.04.004>.  
481

482 Lastra-Mejías, M., González-Flores, E., Izquierdo, M., Cancilla, J.C., & Torrecilla, J.S. (2020).  
483 Cognitive chaos on spectrofluorometric data to quantitatively unmask adulterations of a PDO  
484 vinegar. *Food Control*, *108*, 106860.  
485

486 Lu, X., Ross, C.F., Powers, J.R. Aston, D.E. & Rasco, B.A. (2011). Determination of total phenolic  
487 content and antioxidant activity of garlic (*Allium sativum*) and elephant garlic (*Allium*  
488 *ampeloprasum*) by attenuated total reflectance-fourier transformed infrared spectroscopy, *Journal of*  
489 *Agricultural and Food Chemistry*, *59*, 5215–5221. <https://doi.org/10.1021/jf201254f>.  
490

491 Martins, N., Petropoulos, S., & Ferreira, I.C.F.R. (2016). Chemical composition and bioactive  
492 compounds of garlic (*Allium sativum* L.) as affected by pre- and post-harvest conditions: A review.  
493 *Food Chemistry*, *211*, 41–50. <https://doi.org/10.1016/j.foodchem.2016.05.029>.  
494

495 Mora, M., Elzo-Aizarna, J., Rozas-Fuertes, S., Velilla-Echeita, L., Vázquez-Araújo, L. (2020).  
496 Implicit reaction vs explicit emotional response: Protected designation of origin in apple cider.  
497 *Food Quality and Preference*, *79*, 103773.  
498

499 Minzhen, S., Lun, L., Chuanyun, Z., & Deqing, Z. (2015) Identification and Characterization of  
500 volatile Organic Compounds of Fresh Plant Using Headspace Combined with Surface-Enhanced  
501 Raman Scattering. *Journal of Food Processing & Technology*, *6*, 1-6. [https://doi.org/10.4172/2157-](https://doi.org/10.4172/2157-7110.1000520)  
502 [7110.1000520](https://doi.org/10.4172/2157-7110.1000520)

503 Montaña, A., Beato, V.M., Mansilla, F., & Orgaz, F. (2011). Effect of Genetic Characteristics and  
504 Environmental Factors on Organosulfur Compounds in Garlic (*Allium sativum* L.) Grown in  
505 Andalusia, Spain. *Journal of Agricultural and Food Chemistry*, *59*, 1301-1307.  
506 <https://doi.org/10.1021/jf104494j>.  
507

508 Movasaghi, Z., Rehman, S., & Rehman, I. (2008). Transform Infrared (FTIR) Spectroscopy of  
509 Biological Tissues. *Applied Spectroscopy Reviews*, *43*, 134-179.

510 <https://doi.org/10.1080/05704920701829043>.

511

512 Næs, T., Tomic, O., Mevik, B.-H., & Martens, H. (2011). Path modelling by sequential PLS  
513 regression. *Journal of Chemometrics*, 25, 28–40. <https://doi.org/10.1002/cem.1357>

514

515 Nikolić, V. D., Ilić, D.P., Nikolić, L.B. Stanković, M.Z., Stanojević, L., P., Savić, I., M., & Savić,  
516 I., M. (2012). The synthesis and structure characterization of deoxyalliin and alliin. *Advanced*  
517 *technologies*, 1, 38-46.

518

519 Pérez, N.F., Ferré, J., & Boqué, R. (2009). Calculation of the reliability of classification in  
520 discriminant partial least-squares binary classification. *Chemometrics and Intelligent Laboratory*  
521 *System*, 95, 122-128. <https://doi.org/10.1016/j.chemolab.2008.09.005>

522

523 Rahman, K. (2003). Garlic and aging: New insights into an old remedy. *Ageing Research Reviews*,  
524 2, 39-56. [https://10.1016/S1568-1637\(02\)00049-1](https://10.1016/S1568-1637(02)00049-1).

525

526 Ramirez, D.A., Locatelli, D.A., González, R.E., Cavagnaro, P.F., & Camargo, A.B. (2017).  
527 Analytical methods for bioactive sulfur compounds in Allium: An integrated review and future  
528 directions. *Journal of Food Composition and Analysis*, 61, 4–19.  
529 <https://doi.org/10.1016/j.jfca.2016.09.012>

530

531 Reddy, J.P., & Rhim, J.W. (2014). Isolation and characterization of cellulose nanocrystals from  
532 garlic skin, *Materials Letters*, 29, 20-23. <https://10.1016/j.matlet.2014.05.019>.

533

534 Ritota, M., Casciani, L., Han, B.-Z., Cozzolino, S., Leita, L., Sequi, P., & Valentini, M. (2012).  
535 Traceability of Italian garlic (*Allium sativum* L.) by means of HRMAS-NMR spectroscopy and  
536 multivariate data analysis. *Food Chemistry*, 135, 684–693.  
537 <https://doi.org/10.1016/j.foodchem.2012.05.032>.

538

539 Rocha, S., Pinto, E., Almeida, A., & Fernandes, E. (2019). Multi-elemental analysis as a tool for  
540 characterization and differentiation of Portuguese wines according to their Protected Geographical  
541 Indication. *Food Control*, 103, 27-35.

542

543 Roger, J-M., Palagos, B., Bertrand, D., & Fernandez-Ahumada, E. (2011). CovSel: variable  
544 selection for highly multivariate and multi-response calibration, application to IR spectroscopy.  
545 *Chemometrics and Intelligent Laboratory Systems*, 106, 216-223.  
546 <https://doi.org/10.1016/j.chemolab.2010.10.003>  
547

548 Savitzky, A., & Golay, M.J.E. (1964). Smoothing and differentiation of data by simplified least  
549 squares procedures. *Analytical Chemistry*, 36, 1627-1639. <https://doi.org/10.1021/ac60214a047>  
550

551 Schulz, H., & Baranska, M. (2007). Identification and quantification of valuable plant substances by  
552 IR and Raman spectroscopy. *Vibrational Spectroscopy*, 43, 13-25  
553 <https://doi.org/10.1016/j.vibspec.2006.06.001>  
554

555 Shukla, Y., Kalra, & N., (2007). Cancer chemoprevention with garlic and its constituents, *Cancer*  
556 *Letters*, 247, 167-181. <https://doi.org/10.1016/j.canlet.2006.05.009>.  
557

558 Sjöström, M., Wold, S., & Söderström, B. (1986). PLS discriminant plots. In E.S. Gelsema, &  
559 L.N.Kanal (Eds.), *Pattern recognition in practice* (pp. 461-470). Amsterdam: Elsevier.  
560

561 Slow Food Foundation. <https://www.fondazione Slow Food.com/en/> Accessed 4 April 2019.  
562

563 Smith, R.G. (2005). Determination of the country of origin, of garlic (*Allium sativum*) using trace  
564 metal profiling, *Journal of Agricultural and Food Chemistry*, 53, 4041–4045.  
565 <https://10.1021/jf040166+>.  
566

567 Snee, R. D. (1977). Validation of regression models: methods and examples. *Technometrics*, 19,  
568 415-428.  
569

570 Stark, N.M., Yelle, D.J., Agarwal, U.P. (2016). Techniques for Characterizing Lignin. In O. Faruk  
571 & M. Sain (Eds.), *Lignin in Polymer Composites* (pp. 49-66) Amsterdam: Elsevier  
572 DOI: <https://doi.org/10.1016/B978-0-323-35565-0.00004-7>  
573

574 Ståhle, L., & Wold, S. (1987). Partial least squares analysis with cross-validation for the two-class  
575 problem: a Monte Carlo study. *Journal of Chemometrics*, *1*, 185-196.  
576 <https://doi.org/10.1002/cem.1180010306>

577

578 Trirongjitmoah, S., Juengmunkong, Z., Srikulnath, K., & Somboon, P. (2015). Classification of  
579 garlic cultivars using an electronic nose, *Computers and Electronics in Agriculture*, *113*, 148–153.  
580 <https://10.1016/j.compag.2015.02.007>.

581

582 Vasi, S., Alfa, M., Salvo, A., Bua, G., Giofrè, S., Corsaro, C., Mallamace, D., Cicero, N., Mottese,  
583 A., Vadalà, R., & Dugo, G. (2016). Statistical Analysis of Mineral Concentration for the  
584 Geographic Identification of Garlic Samples from Sicily (Italy), Tunisia and Spain. *Foods*, *5*, 20.  
585 <https://10.3390/foods5010020>.

586

587 Yun, H.M., Ban, J.O., Park, K.R., Lee, C.K., Jeong, H.S., Han, S.B., & Hong, J.T. (2014). Potential  
588 therapeutic effects of functionally active compounds isolated from garlic. *Pharmacology &*  
589 *Therapeutics*, *142*, 183–195. <https://doi.org/10.1016/j.pharmthera.2013.12.005>.

590

591 Wold, S., Johansson, E., & Cocchi, M. (1993). PLS: partial least squares projections to latent  
592 structures. In H. Kubinyi (Ed.), *3D QSAR in Drug Design: Theory, Methods and Applications* (pp.  
593 523-550), Leiden, The Netherlands: KLUWER ESCOM Science Publisher.

594

595 Wold, S., Martens, H., & Wold, H. (1983). The multivariate calibration problem in chemistry  
596 solved by the PLS method. In A. Ruhe, B. Kågström (Eds.), *Proceedings of the Conference on*  
597 *Matrix Pencils. Lecture Notes in Mathematics* (pp 286-293), Heidelberg, Germany: Springer  
598 Verlag.

599

600

## 601 **Figure captions**

602 Figure 1 Mean spectra collected on cloves (in red) and on tunics (in blue).

603 Figure 2 VIP analysis. Solid black lines represent mean spectra for the four different categories  
604 (offset to make them visible) while bold variables highlight selected features. The upmost plot  
605 refers to spectra collected on samples from Sulmona, lines in the middle are mean spectrum for

606 class Nubia and Proceno, respectively; the lowest lines represent the mean signal for samples  
607 belonging to class Castelliri.

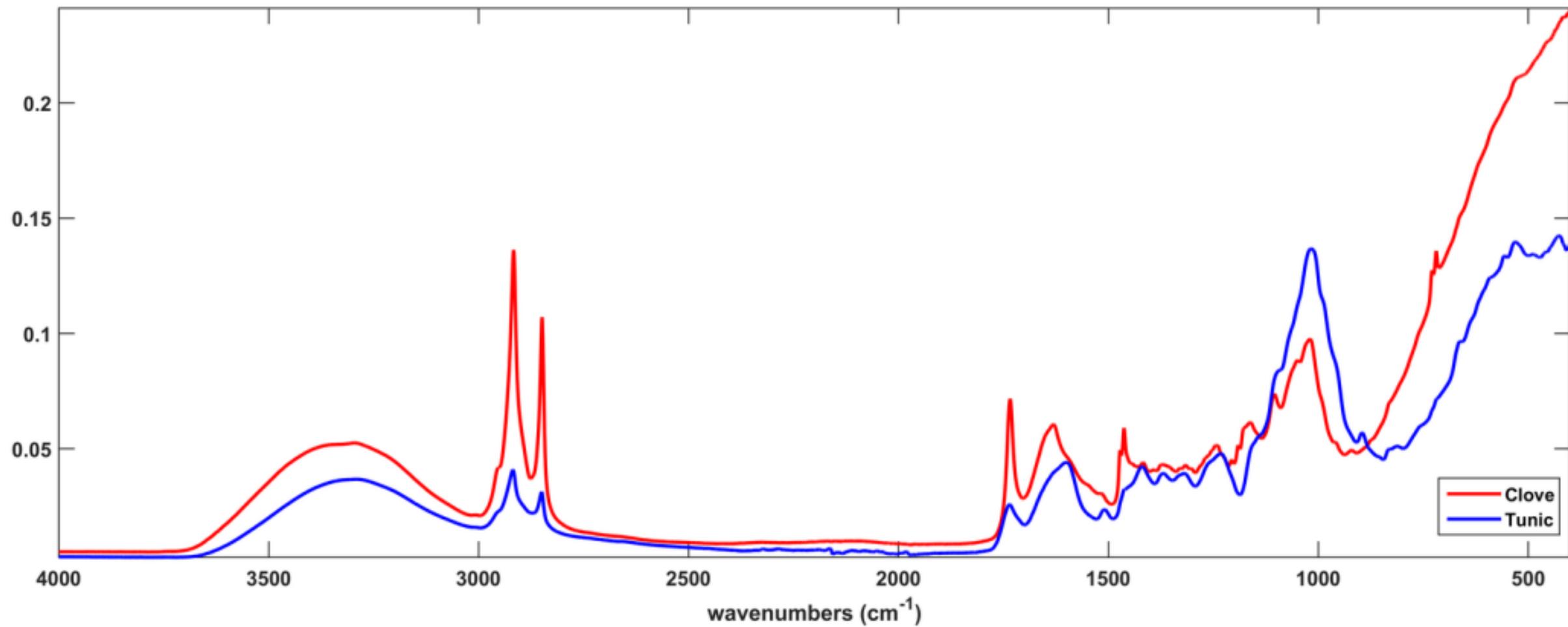
608

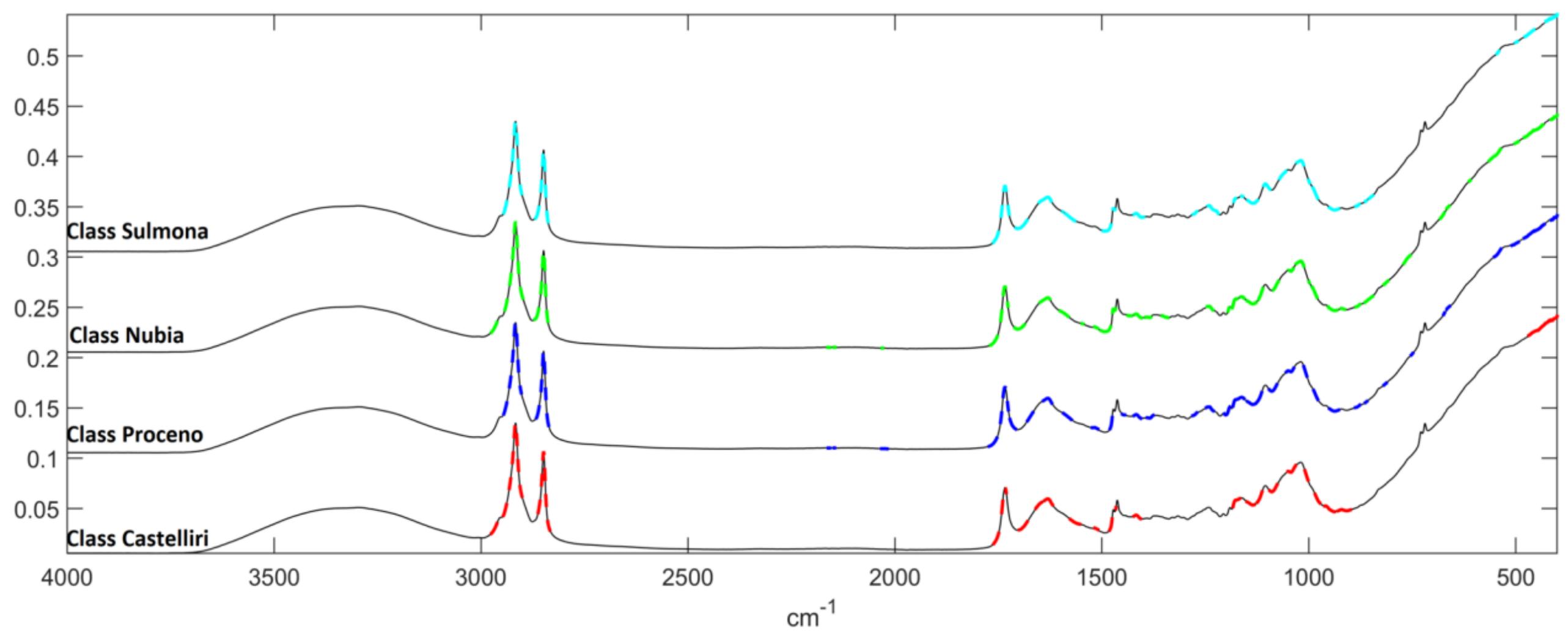
609 Figure 3 SO-PLS-LDA analysis: Samples project onto the first to canonical variates. Legend: Red  
610 circles: Class Castelliri; Blue squares: Class Proceno; Black diamonds: Class Nubia; Green  
611 triangles: Class Sulmona. Empty and filled symbols represent training and test samples,  
612 respectively.

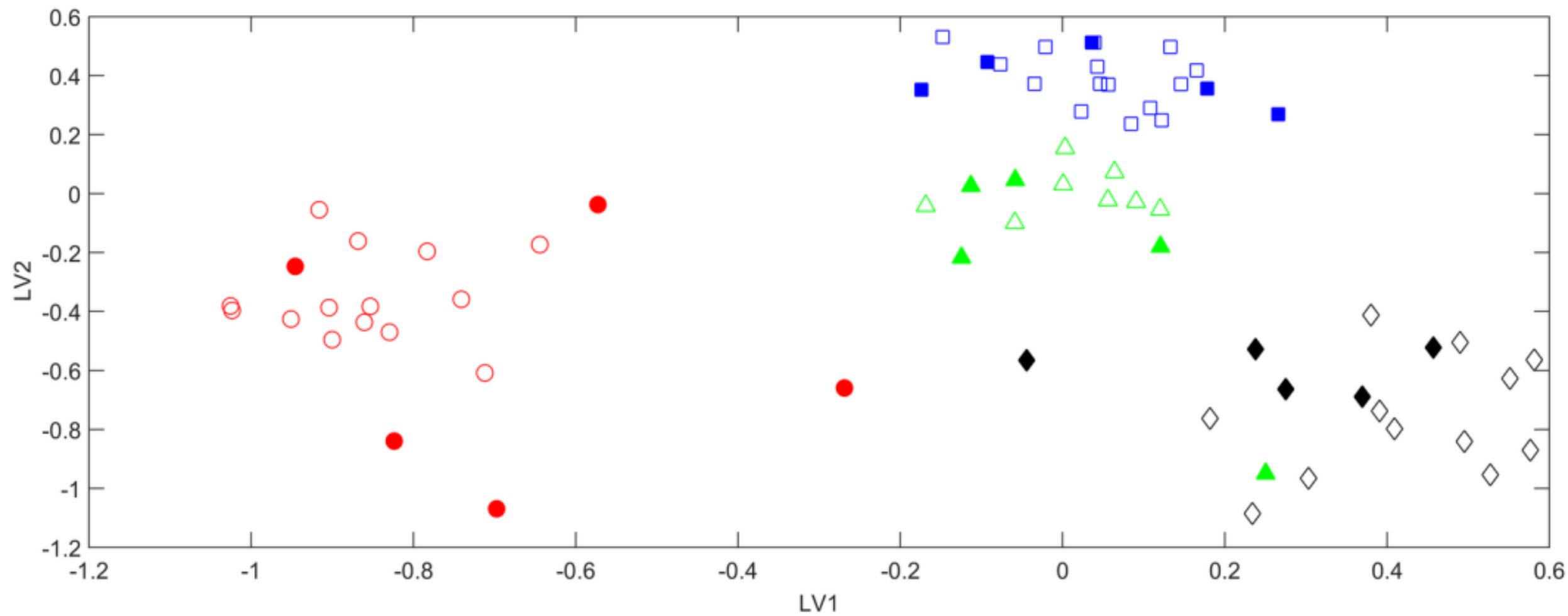
613

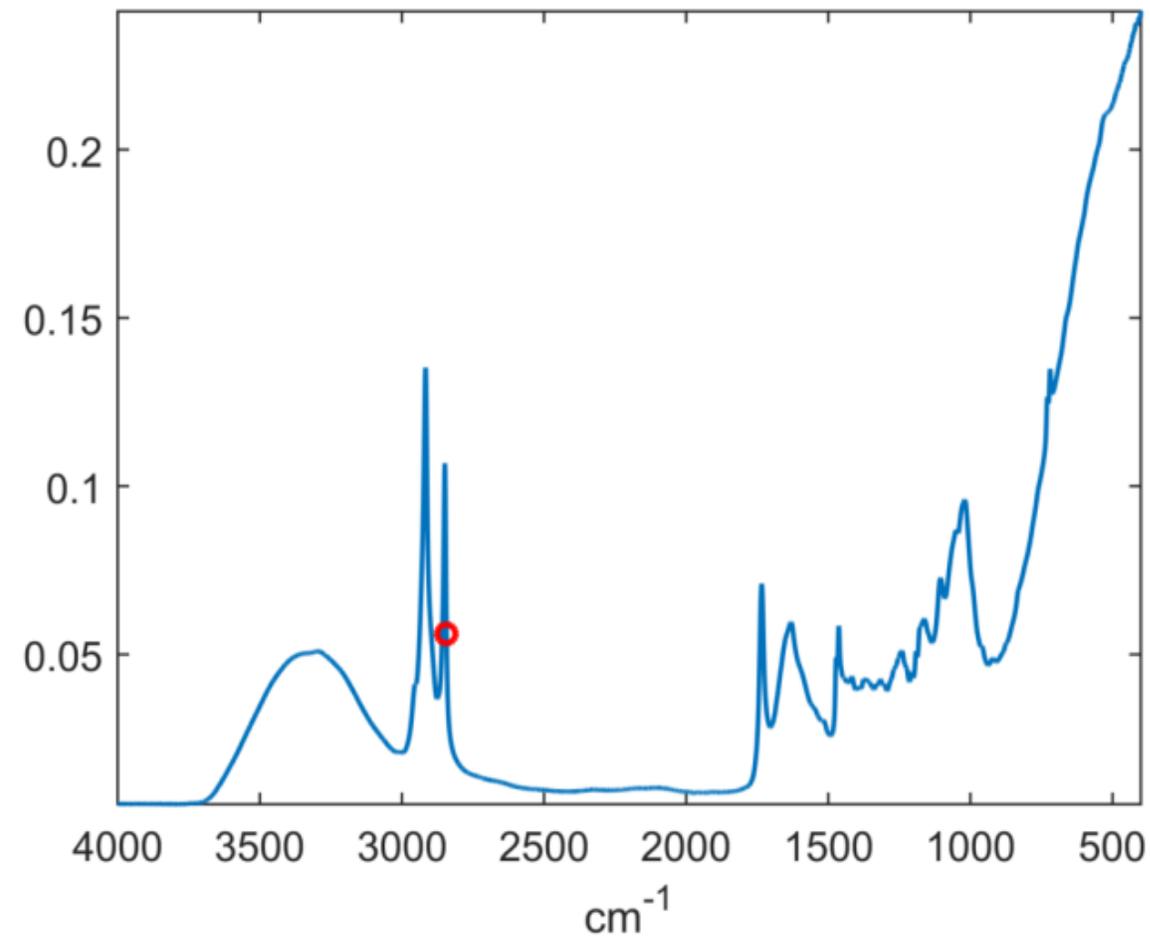
614 Figure 4 SO-CovSel Analysis: blue line represents mean spectrum collected on a) cloves b) tunics.

615 Variables selected by SO-CovSel are circled in red.

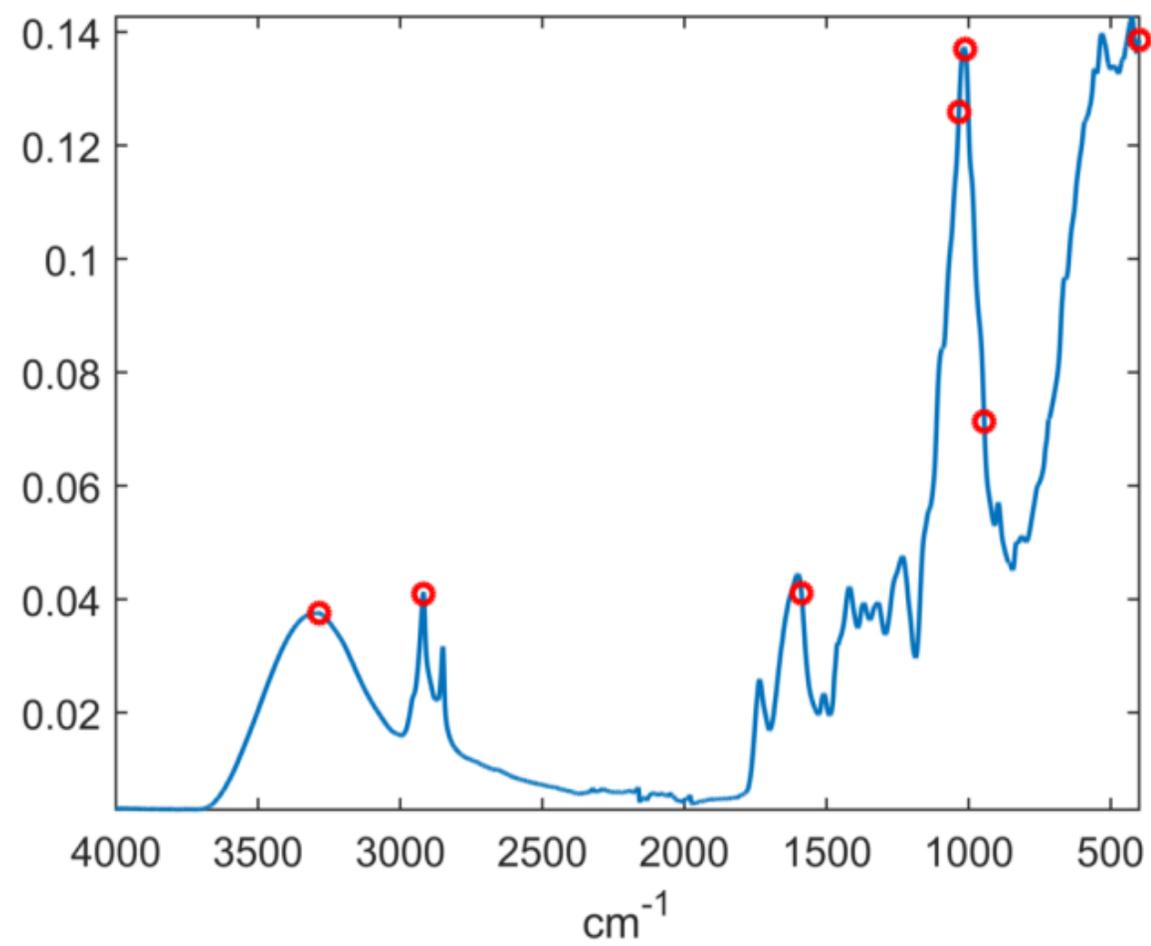








**a)**



**b)**

Table 1 PLS-DA analysis on tunics: Cross-validated mean classification errors (%) as function of preprocessing and complexity (LVs).

Preprocessing	LVs	Average classification errors (%-CV)
Mean Centering (MC)	10	12.3
1 <sup>st</sup> Derivative + MC	12	10.2
2 <sup>nd</sup> Derivative + MC	7	12.5
SNV + MC	10	11.8
SNV + 1 <sup>st</sup> Derivative + MC	8	8.7
SNV + 2 <sup>nd</sup> Derivative + MC	10	12.0

Table 2 PLS-DA analysis on cloves: Cross-validated mean classification errors (%) as function of preprocessing and complexity (LVs).

Preprocessing	LVs	Average classification errors (%)
Mean Centering (MC)	13	13.4
1 <sup>st</sup> Derivative + MC	15	10.5
2 <sup>nd</sup> Derivative + MC	13	10.9
SNV + MC	16	12.8
SNV + 1 <sup>st</sup> Derivative + MC	14	12.3
SNV + 2 <sup>nd</sup> Derivative + MC	9	16.4

Table 3 Classification rates (on the test set) for the four proposed strategies.

Method	Classification rates on the test set (%)
PLS-DA on tunics	75.0
PLS-DA on cloves	85.0
SO-PLS-LDA	95.0
SO-CovSel-LDA	85.0