

1 **Detection of adulterants in grape nectars by attenuated total reflectance Fourier-transform mid-infrared**
2 **spectroscopy and multivariate classification strategies**

3

4 **Carolina Sheng Whei Miaw^{a,b,c}, Marcelo Martins Sena^d, Scheilla Vitorino Carvalho de Souza^a, Maria Pilar Callao^{c,*} and Itziar**
5 **Ruisanchez^c**

6 ^a Department of Food Science, Faculty of Pharmacy (FAFAR), Federal University of Minas Gerais (UFMG), Av. Antônio Carlos, 6627, Campus da UFMG, Pampulha, 31270-010, Belo Horizonte, MG, Brazil.

7 ^b CAPES Foundation, Ministry of Education of Brazil, 70040-020, Brasília, DF, Brazil.

8 ^c Chemometrics, Qualimetric and Nanosensors Grup, Department of Analytical and Organic Chemistry, Rovira i Virgili University, Marcel·lí Domingo s/n, 43007 Tarragona, Spain

9 ^d Department of Chemistry, Institute of Exact Sciences (ICEX), Federal University of Minas Gerais (UFMG), Av. Antônio Carlos, 6627, Campus da UFMG, Pampulha, 31270-010, Belo Horizonte, MG, Brazil.

10 * Corresponding author: mariapilar.callao@urv.cat

11

12 **Abstract**

13 There is no any doubt about the importance of food fraud control, as it has implications in food safety and in consumer health. Focusing
14 on fruit beverages, some types of adulterations have been detected more frequently, such as substitution with less expensive fruits. A
15 methodology based on attenuated total reflectance Fourier-transform mid-infrared spectroscopy (ATR-FTIR) and multivariate
16 classification was applied to detect whether grape nectars were adulterated by substitution with apple juice or cashew juice. A total of
17 126 samples were obtained and analyzed. Two strategies were proposed: one-class and multiclass approaches. Soft independent
18 modeling of class analogy (SIMCA), partial least squares discriminant analysis (PLS-DA) and partial least squares density modeling
19 (PLS-DM) were used to build the models. Among them, PLS-DA presented the best performance with a sensitivity and specificity of
20 nearly 100%. The multiclass strategy was preferred if the adulterants to be studied are known because it provides additional information.

21

22 **Highlights**

- 23 • The detection of grape nectar adulteration with cashew or apple was studied.
- 24 • One-class and multiclass approaches were implemented.
- 25 • The multivariate classification methods SIMCA, PLS-DA and PLS-DM were compared.
- 26 • PLS-DA provided better performance to detect grape nectar adulterations.

27

28 **Keywords**

29 Food adulteration, Fruit nectar, PLS-DA, SIMCA, One-class classification, Multiclass classification.

30

31 **1. Introduction**

32 Because of the highly competitive market, drink industries are always searching for product diversification, and in recent years,
33 the largest increase in production was of fruit nectar (Neves, Trombin, Lopes, Kalaki, & Milan, 2012). Nectar is defined as an
34 unfermented beverage produced by the dilution in water of the edible part of fruits or vegetables or their extracts with the addition of
35 sugars, intended for direct consumption (Brazil, 2009). In Brazil, Standards of Identity and Quality (SIQ) are established for fruit
36 nectars and cover the minimum percentages of pulp that must be used in each type of nectar. For some fruits, the minimum parameters
37 include soluble solids (SS), total titratable acidity (TTA), total sugars (TS) and ascorbic acid (AA) (MAPA, 2003, 2013). According to
38 the Brazilian Association of Soft Drinks and Non-Alcoholic Beverages (ABIR), the most consumed nectar in Brazil is grape flavor
39 (ABIR, 2015).

40 Considering the issue of adulteration of fruit-based beverages, the most frequent practices include substitution with cheaper
41 ingredients, such as simple dilution with water or sugar syrup, and undeclared addition of different species, which can be botanically
42 related, or not, to the main fruit in question (Asadpoor, Ansarin, & Nemati, 2014).

43 Methods based on spectroscopic techniques are generally rapid, non-destructive, simple and require little or no sample
44 preparation. However, they have the disadvantage of low specificity. Therefore, powerful tools for adulteration testing can be created
45 by combining these techniques with multivariate chemometric methods, while some authors applied just basic statistical techniques (El
46 Darra et al., 2017). Classification methods are particularly suitable for food fraud detection. They can be differentiated in discriminant
47 and class-modeling methods. The most common discriminant method is partial least squares discriminant analysis (PLS-DA), while
48 the most used class-modeling method is soft independent modeling of class analogy (SIMCA) (Bevilacqua et al., 2013).

49 The necessity of food quality control was reflected in a specific review concerning the development of an effective food
50 traceability system to reduce the numerous cases of food safety incidents and fraudulence. (Dandage, Badia-Melis, & Ruiz-García,
51 2017). In that sense, reviews have been recently published addressing the use of multivariate classification methods to authenticate or
52 detect adulteration in food (Callao & Ruisánchez, 2018; Esteki, Shahsavari & Simal-Gandara, 2018; Szymańska et al., 2015).
53 Multivariate classification methods have been successfully applied to elucidate specific problems of authenticity or adulteration in
54 different types of food. Examples are wines (Sen & Tokatli, 2016), oils (Georgouli, Del Rincon, & Koidis, 2017), , milk (Gondim,

55 Junqueira, Souza, Ruisánchez, & Callao, 2017), hazelnut pastes (López, Trullols, Callao, & Ruisánchez, 2014), coffee (Bona et al.,
56 2017), mushrooms (Xu et al., 2016), vinegar (Ríos-Reina, Callejón, Oliver-Pozo, Amigo, & García-González, 2017) and whiskies
57 (Martins, Talhavini, Vieira, Zacca, & Braga, 2017).

58 Comparatively, the application of these techniques to studies involving authentication or detection of frauds in fruits and
59 derivatives is more limited. For this aim, articles have developed for multivariate classification or calibration models employing
60 different analytical techniques, such as UV-VIS spectroscopy (Boggia, Casolino, Hysenaj, Oliveri, & Zunin, 2013), spectrofluorometry
61 (Ammari, Redjda, & Rutledge, 2015), nuclear magnetic resonance (NMR) (Cuny et al., 2008) and mid-infrared spectroscopy (He,
62 Rodriguez-Saona, & Giusti, 2007; Miaw et al., 2018; Shah, Cynkar, Smith, & Cozzolino, 2010; Shen et al., 2016;).

63 In the present study, the detection of grape nectar adulteration with apple and cashew juices was studied by means of attenuated
64 total reflectance Fourier-transform mid-infrared spectroscopy (ATR-FTIR) and classification methods. Apple juice has commonly been
65 used as filler for economic gain by beverage industries (Singhal, Kulkarni, & Rege, 1997), but it is now also being used to replace some
66 of the added sugar. Furthermore, cashew and apple are fruits suspected of being utilized for adulterations by fraudulent industries,
67 justifying the importance of the development of analytical methods to detect these potential adulterants in the most popular beverage
68 products, such as the grape nectar matrix.

69 In this paper, two approaches were proposed considering their different purposes: one-class and multiclass approaches, utilizing
70 discriminant or class-modeling methods. One-class classification is adequate when the goal is to test whether a sample is adulterated,
71 regardless of which adulterant might be present (López et al., 2014). If the adulterant is known, the multiclass strategy can be chosen,
72 since it gives additional information, such as multiple assignments and samples not assigned to any class (Gondim et al., 2017).

73 In recent years, some authors have criticized the predominance in the chemometric literature of the use of discriminant methods,
74 such as PLS-DA, to food authentication problems (Rodionova, Oliveri, & Pomerantsev, 2016; Oliveri, 2017). This criticism has noted
75 that classification results will be unreliable when the model is used to predict a new sample from an untrained class. In response, other
76 authors have combined PLS-DA with outlier detection, identifying samples from untrained classes based on large *Hotelling T²* and *Q*
77 residues (Martins et al., 2017). However, as class-modeling models are developed using only the information concerning one-class
78 samples at a time, they are unable to ensure the model specificity for the detection of various food frauds (Xu et al., 2016). Considering
79 all these relevant discussions, it is important to compare the alternatives for developing supervised classification models for detecting
80 food fraud. Thus, SIMCA and PLS-DA, as the most used class-modeling and discriminant methods, respectively, were applied to the
81 authentication of grape nectars. In addition, a recently proposed one-class modeling method, partial least squares density modeling

82 (PLS-DM) (Oliveri et al., 2014), was also applied. The three classification methods were compared through the evaluation of sensitivity
83 and specificity.

84

85 2. Materials and methods

86

87 2.1 Formulation of nectars

88 Grape nectars samples, were prepared starting from reliable raw materials and rigorously meeting the established regulations
89 (MAPA, 2003, 2013), at the Food Science Laboratory and at the Technology Laboratory, both located in the Food Department of the
90 Faculty of Pharmacy of the Federal University of Minas Gerais (UFMG).

91 Isabel grape samples were obtained from EMBRAPA (the Brazilian Agricultural Research Corporation) Grape & Wine, located
92 in Petrolina, PE, Brazil. Red cashews and Fuji apples were acquired from the Minas Gerais Supply Center (CEASA) in Contagem,
93 MG, Brazil. The selection of fruits took into account the absence of mechanical and phytopathological damage, the degree of maturation
94 and other typical physical characteristics of each fruit, such as size, color and texture (Paltrinieri & Figuerola, 1998). The fruits were
95 stored in the refrigerator at 4-7°C until the preparation of nectars (EMBRAPA, 2016).

96 The fruits were sanitized with 100 mg/L of sodium hypochlorite solution (Vetec Química Fina, Ltda, Rio de Janeiro, RJ, Brazil)
97 for 2 min and washed. The juices/pulps of grape, apple and cashew were obtained as described below:

- 98 • grapes were heated under constant steam for 1 to 2 h in an autoclave (Fanem, São Paulo, SP, Brazil) at 100 °C, pressed and
99 sieved to obtain the juice;
- 100 • apples were peeled and cut into eight pieces, and the seeds were removed. The fruits were scalded in boiling water for 3 min,
101 followed by immersion in water with ice until cooling;
- 102 • cashews had their chestnut removed and the fruits were cut into four pieces.

103 Apples and cashews were individually pulped in an industrial blender (Fisatom 752, São Paulo, SP, Brazil) and sieved (1 mm
104 sieve).

105 For the formulation of grape nectars, the only SIQ parameter recommended in the Brazilian legislation (MAPA, 2003) is a
106 minimum of 50 % of pulp, which was considered in the formulations of the unadulterated nectars. The amounts of pulp/juice and syrup
107 were estimated as described in Equation (1).

$$108 \frac{a \times A}{100} + \frac{b \times B}{100} + \frac{c \times C}{100} = \frac{m \times (A+B+C)}{100} \quad (1)$$

109 where "a" represents pulp Brix, "A" represents percentage of pulp that must be present in the nectar, "b" represents syrup Brix, "B"
110 represents percentage of syrup, "c" represents adulterant pulp Brix, "C" represents percentage of adulterant, "m" represents final nectar
111 Brix, and "A + B + C" is equal to 100 (pulp + syrup + adulterant) (Tressler & Joslyn, 1961).

112 The quantity of additives added was 0.25 g/100 g, 15 mg/100 g and 0.075 g/100 mL for citric acid, ascorbic acid and guar gum
113 (Pryme Foods, Sorocaba, SP, Brazil), respectively. Syrup at 20 °Brix was prepared and added to the additives in adequate proportions
114 to produce nectars with 11 to 13 °Brix. These values were within the ranges permitted by Brazilian legislation and based on preliminary
115 experiments involving commercial nectars (Miaw et al., 2018).

116 Juices were added to the syrup, homogenized (Fisatom 752, São Paulo, SP, Brazil) and filled in labelled amber glass bottles
117 (250-mL) with plastic screw caps (both previously sterilized by autoclaving at 100 °C for 10 min). Nectars were pasteurized in the
118 autoclave at 100 °C for 10 min. Bottles were hermetically sealed and left at room temperature (Paltrinieri & Figuerola, 1998). After
119 being opened for analysis, the nectar bottles were refrigerated (4 - 7 °C).

120 As illustrated in **Fig. 1**, a set of 42 samples of grape nectar were prepared for each of the three studied classes: unadulterated,
121 adulterated with cashew, and adulterated with apple.

122 First, seven representative batches of each class were prepared according to the following formulations:

- 123 a) unadulterated batches were formulated with 50 % of grape, sugar syrup and additives (corresponding to the other
124 50 %).
- 125 b) batches adulterated with cashew were formulated with 40 % grape, 10 % of cashew juice, sugar syrup and additives
126 (corresponding to the other 50 %).
- 127 c) batches adulterated with apple were formulated with 40 % grape, 10 % of apple juice, sugar syrup and additives
128 (corresponding to the other 50 %).

129 Then, to obtain the 42 representative samples of each class, the 7 above described batches were mixed taking 3 of them in
130 almost the same proportion (35/35/30) to give the additional 35 samples. The final number of samples was 126.

131

132 *2.2 Instrumentation and software*

133 The Brix degrees of each juice/pulp produced was measured using a refractometer (Hanna Instruments Brasil, Barueri, SP,
134 Brazil).

135 Samples were analyzed by ATR-FTIR in an IRAffinity-1 FTIR (Shimadzu, Kyoto, Japan) spectrophotometer with a DLATGS
136 detector (Deuterated Triglycine Sulfate Doped with L-Alanine) equipped with a horizontal ATR accessory with a ZnSe prism (PIKE

137 Technologies, Madison, WI, USA) of 20 internal reflections. For each sample, 1.5mL were pipetted onto the ATR cell surface and
138 three readings were recorded with 16 scans, 4 cm⁻¹ resolution, generating spectra between 4000 to 937 cm⁻¹. A background correction
139 was performed after each measurement to avoid atmospheric interference and reduce instrumental noise.

140 Multivariate analysis was conducted using MATLAB software version 8.0.0.783 - R2012b (Natick, MA, USA) and PLS Toolbox
141 7.0.2 (Eigenvector Research Inc., Wenatchee, WA, USA).

142

143 *2.3 Data analysis*

144

145 *2.3.1 Pre-processing and exploratory analysis*

146 Multiplicative scatter correction (MSC) (Rinnan, Berg, & Engelsen, 2009) was applied to correct the spectra baseline
147 deviations. Principal component analysis (PCA) was used as an unsupervised exploratory analysis tool to visualize the sample
148 distribution in the multivariate space, to identify any natural clustering in the samples that could influence the subsequent multivariate
149 analysis and to identify possible outliers.

150

151 *2.3.2 Classification methods*

152 Multivariate classification methods are supervised techniques. They can be divided, among other criteria, into class-modeling
153 and discriminant methods. Discriminant methods define delimiters in the hyperspace of the variables, separating the samples into a
154 number of regions corresponding to the number of predefined classes, and focusing on the differences between the samples from each
155 class. Class-modeling methods build an individual model for each predefined class regardless of the information for the other classes
156 or categories and focusing on the similarities between samples from the same class (Bevilacqua et al., 2013)

157 SIMCA is a modelling technique based on Principal Component Analysis (PCA) in which each class is modelled independently
158 from all others (Bevilacqua et al., 2013). Each sample is characterized by two scalar statistics, Hotelling T² and Q, which measures the
159 information from each sample included or not included in the model, respectively. Class frontiers (*Hotelling T_{lim}²* and *Q_{lim}*) are
160 calculated for each pre-defined class (class model), at a specific significance level (α), usually set at 0.05 (Rius, Callao & Rius, 1997)

161 Historically, various criteria have been used for the classification of samples in SIMCA models. A common criterion assigns
162 samples to classes based on their reduced values *Hotelling T_r²* and *Q_r*. These values are the ratios between the statistics of sample *i* (*T_i²*
163 and *Q_i*) and the corresponding statistical limits for each class. A sample must have values lower than 1.0 for both the reduced parameters
164 to be considered within the class model. The most used criterion is a slight variation of the former. A sample *i* is assigned based on its

165 distance from class j (d_{ij}), which is defined as a combination of its reduced parameters (Equation (2)) (Márquez, López, Ruisánchez,
166 & Callao, 2016). In this last case, the class boundary for a sample to be assigned as within the model is a semi-circle with a radius 1.0
167 (d equal to or lower than 1.0), so this criterion is more restrictive than considering *Hotelling* T_r^2 and Q_r statistics independently.

$$168 \quad d_{ij} = \sqrt{(Q_{r,i})^2 + (T_{r,i}^2)^2} \quad (2)$$

169 PLS-DA is a discriminant method that adapts PLS regression to a classification task. It establishes a linear regression between
170 a matrix of independent variables (\mathbf{X}) and an array of dependent variables (\mathbf{Y}). \mathbf{Y} contains binary dummy variables that indicate the
171 class to which each sample belongs, where 1 indicates membership and 0 does not (Barker & Rayens, 2003). Since this paper aimed to
172 differentiate and classify between three classes, class 1 samples were encoded as (1,0,0), class 2 as (0,1,0) and class 3 as (0,0,1).

173 The PLS-DA model predicts the class for each sample, assigning values approximately 0 or 1. Bayesian statistics are used to
174 calculate the threshold value above which the sample is considered to belong to the class (Bylesjö et al., 2006). The Bayesian threshold
175 considers that y predicted values of the PLS-DA model are normally distributed, selecting the y value in which the number of false
176 results are minimal (false-negatives and false-positives) (Pulido, Ruisánchez, Boqué, & Rius, 2003). Thus, predicted values above or
177 below this threshold mean that a sample does or does not belong to the class, respectively.

178 PLS-DM is a one-class method that adapts PLS regression to a classification task. Its particularity is that PLS-DM computes
179 the response vector \mathbf{y} as an estimation of sample density, based on inter-sample distances in the multivariate space. With the algorithm
180 used in this work (Oliveri, 2017; Oliveri et al., 2014), for each sample in the training set, the response vector \mathbf{y} is calculated as the sum
181 of Euclidean distances between k samples with the lowest distance in the multivariate space. The algorithm applies all possible
182 combinations using the parameter distance of k nearest neighbors, smoothing coefficient α (for the definition of the class space in the
183 PLS score domain), the number of latent variables (LV) and the pre-processing suitable for the \mathbf{X} matrix. Then, the best combination
184 is chosen with the adjustment of the number of LV using efficiency criteria (geometric mean of sensitivity and specificity) and with
185 the evaluation of the other parameters.

186 For this model, the specificity is calculated in the presence of the non-target class, which can be composed of more than one
187 extraneous class. In this case, the specificity obtained is calculated from the overall alternative class. If the specificity of each specific
188 alternative class is required, it must be calculated for each non-target class separately (Rodionova, Oliveri, & Pomerantsev, 2016).

189

190 *2.3.3 Performance Parameters*

191 The performance parameters are measurable attributes that indicate the quality of the analytical method (López, Callao, &
192 Ruisánchez, 2015). For qualitative methods, the most common parameters are sensitivity, specificity and the more recently proposed

193 inconclusive ratio. The first two are based on probabilities regarding four possible binary responses: true positive (TP) (positive
194 response for a sample that is positive), false positive (FP) (positive response for a sample that is negative), true negative (TN) (negative
195 response for sample that is negative) and false negative (FN) (negative response for a sample that is positive). The expressions to
196 calculate these values are presented below.

197 Sensitivity (SEN) indicates the likelihood of recognizing samples that truly belong to the modeled class (samples from class j ,
198 $n^{\circ}S_j$, that have been properly predicted by the model as belonging to class j).

$$199 \quad \mathbf{SEN}_j = \mathbf{TP}_j / \mathbf{n}^{\circ}\mathbf{S}_j \quad (3)$$

200 Specificity (SPE) indicates the likelihood of recognizing samples that are truly different from the modeled class (samples that
201 are not from class j , $n^{\circ}S_{not\ j}$, that have been properly predicted as not belonging to class j).

$$202 \quad \mathbf{SPE}_j = \mathbf{TN}_j / \mathbf{n}^{\circ}\mathbf{S}_{not\ j} \quad (4)$$

203 Inconclusive ratio (IR) indicates the percentage of samples that cannot be undoubtedly assigned to class j , and thus considers
204 no assignation to any class and the multiple assignation (López et al., 2014).

$$205 \quad \mathbf{IR}_j = (\mathbf{NA}_j + \mathbf{MA}) / \mathbf{n}^{\circ}\mathbf{S}_j \quad (5)$$

206 where NA_j means unassigned samples (samples that are from class j that are not assigned to class j or any other class); MA means
207 multiple assignation samples (samples from class j assigned to more than one class) and $n^{\circ}S_j$ means the total number of samples that
208 really belong to class j .

209

210 3. Results and discussion

211 **Fig. 2** shows the mean pre-processed spectra of each predefined class under study. As previously observed (Miaw et al., 2018),
212 the intense band near 3300 cm^{-1} and the sharp peak at 1640 cm^{-1} present in all samples are related to the O-H absorption of water (He
213 et al., 2007; Shen et al., 2016). The region between 1700 and 1000 cm^{-1} incorporates the typical bands for phenolic compounds, such
214 as the C=C-C aromatic ring stretching, the phenol OH bending, the aromatic C-H in-plane bending, and the C-O stretching of phenol
215 (Bureau, Ścibisz, Le Bourvellec, & Renard, 2012). Additionally, in this region, sugars and organic acids are present showing the
216 characteristic bands (between 1500 and 950 cm^{-1}) (Shah et al., 2010; Shen et al., 2016). The low-intensity bands between 1500 and
217 1200 cm^{-1} were related to the deformations of CH_2 , C-C-H and H-C-O (Shah et al., 2010; Vardin, Tay, Ozen, & Mauer, 2008). For the
218 fingerprint region (1200 to 900 cm^{-1}), the stretching vibrations of C-C and C-O bonds correspond to the presence of sugars and organic

219 acids (He et al., 2007; Shah et al., 2010; Vardin et al., 2008). These described components are present in all the nectars, justifying the
220 similarities among the spectra of the three classes showed in Fig. 2.

221 First, an exploratory analysis by Principal Component Analysis (PCA) was performed on all the samples from the three classes
222 studied. The scores plot of the first two principal components (PC1 \times PC2), accounting for 90.32 % of the total variance, are illustrated
223 in **Fig. 3**. It can be seen that PC1 could not distinguish between the 3 classes. Along the PC2, samples adulterated with apple (squares)
224 presented negative scores values and were clearly separated from unadulterated samples, which presented positive score values
225 (triangles). Samples adulterated with cashew (circles) appeared to clearly overlap with the unadulterated samples, and just a few of
226 them appeared to overlap with the apple adulterated samples.

227 For the supervised classification modeling, each class was separated into training and test sets using the Kennard and Stone
228 algorithm (28 samples for training and 14 for test set) which selects representative and uniformly distributed samples into the
229 multivariate space (Kennard & Stone, 1969).

230 Initially, a multiclass strategy was implemented by applying SIMCA and PLS-DA classification techniques to establish the
231 three classes: unadulterated (UN), adulterated with cashew (CAS) and adulterated with apple (APP). SIMCA models were
232 independently established for each class using the training set and the optimal numbers of PCs were selected based on the lowest value
233 of RMSECV (root mean square error of cross validation). The models were validated using leave-one-out cross validation as well as
234 predictions of the test set. Three PCs for each class were necessary to build the SIMCA model, accounting for 95.20, 93.79 and 90.58 %
235 of total variance, for UN, CAS and APP classes, respectively.

236 PLS-DA models were also built with the three classes. The model was validated using leave-one-out cross validation and the
237 number of LV, chosen based on the smallest cross validation classification errors, was 6, accounting for 95.13 % of variance in the **X**
238 block and 82.94 % in the **Y** block. The threshold values were 0.25 for the UN class, 0.14 for the CAS class and 0.09 for the APP class,
239 as can be observed from **Fig. 4**.

240 The summarized class assignments obtained by applying SIMCA and PLS-DA models are presented in **Table 1**. Regarding the
241 results obtained with SIMCA, as expected considering the PCA model shown in **Fig. 3**, samples from UN and CAS classes were
242 multiply assigned to each other. Almost all unadulterated samples, in both training and test sets, were doubly assigned to their class
243 and to the CAS class. To a lesser extent, samples adulterated with cashew, five from the training set and seven from the test set, were
244 also doubly assigned to their class and as unadulterated (UN class), and seven of the 28 training samples were not assigned to any class.
245 Finally, as expected, samples adulterated with apple were properly recognized by their class model, with no wrong or multiple
246 assignment to other classes. Only six of the 28 samples from the training set were not assigned to any class, while all samples from the

247 test set were correctly assigned. As a result of the assignments, high inconclusive ratios were obtained for all three classes, and the
248 unadulterated class was the one with the highest ratio (**Table 1**).

249 For results obtained with PLS-DA (**Table 1** and **Fig. 4**), no incorrect assignments were obtained. In addition, few inconclusive
250 assignments, all corresponding to samples adulterated with cashew, were observed: one sample from the training set that was doubly
251 assigned, and one sample from the training and three from test set that were not assigned to any class. Notably, in no cases were
252 adulterated samples assigned as unadulterated; this outcome means that no false-negative errors were obtained. From the perspective
253 of food fraud, false-negative errors are the most important to control, as they correspond to errors related to not detecting the
254 contaminant when it is present.

255 The next step was the implementation of a one-class strategy, in which only the target class was established by all three
256 classification methods. The UN class was considered the target class and CAS and APP samples were jointly the non-target class. This
257 SIMCA model was similar to the previous one established for the multiclass approach. The only difference is reflected in the calculation
258 of specificity, since CAS and APP samples were modeled together in a single class. PLS-DA was established for two contrasting
259 classes, encoded as (1,0), with 1 as the UN class and 0 as the CAS+APP class. This model was built as in the multiclass approach,
260 namely, the number of LV 5, which accounted for 94.29% of variance in the **X** block and 37.49% in the **Y** block.

261 As has been explained in the theory section (2.3.2), PLS-DM implies the optimization of several parameters: the number of
262 nearest neighbors' k , from 1 to 6; pre-processing type; smoothing coefficient α of the potential function, from 0.3 to 0.8; and the number
263 of LV, from 1 to 10. The optimization step was applied in the training set and, as a result, a matrix of sensitivity, specificity and
264 efficiency values (data not shown) was obtained for all studied values of these parameters. The optimal combination of these results
265 was evaluated considering the highest efficiency and an odd number of k nearest neighbors. Even k values can lead to ambiguous
266 classifications, which is the reason why odd numbers are preferred. The optimal parameter values were set as $k = 3$, mean-center
267 preprocessing, $\alpha = 0.6$ and $LV = 4$.

268 The classification results for these three methods in terms of sensitivity and specificity, according to the one-class strategy, are
269 summarized in **Table 2**. PLS-DA presented the best predictions, since both the sensitivity and the specificity of the training and the
270 test set was 100%. Regarding the results of both SIMCA and PLS-DM, they cannot be considered satisfactory, especially in relation to
271 specificity, since a significant percentage of adulterated samples were predicted as not adulterated (25% for SIMCA and 32% for PLS-
272 DM).

273 When the two strategies are compared, it can be stated that the multiclass classification would be preferable, because it provides
274 more specific information about the adulterations. Many samples in the one-class strategy were erroneously assigned, and in the multi-
275 class were considered inconclusive; therefore, a confirmatory analysis is required.

276 Regarding the comparison among the three classification methods, PLS-DA, SIMCA and PLS-DM, the best performance was
277 clearly provided by the discriminant PLS-DA model. This superior performance of discriminant over class-modeling methods is
278 consistent with observations in the chemometric literature (Bylesjö et al, 2006). Class-modeling methods, such as SIMCA, search for
279 data directions of the highest variance, which might be distinct from the variance direction responsible for the separation of classes. A
280 specific explanation for the worse results provided by class-modeling methods (SIMCA and PLS-DM) in our case is the similarity
281 between UN and CAS samples, which was verified by observing their highly overlapped clusters in the PCA model shown in **Fig. 3**.

282

283 **4. Conclusions**

284 The combination of ATR-FTIR and classification techniques allowed the detection of adulterations of grape nectars with apple
285 and cashew juices. The entire analytical procedure was very simple and rapid, and it did not require sample pretreatment or the
286 consumption of reagents or solvents. All 126 samples used in this study were obtained from reliable raw ingredients and prepared in
287 strict compliance with Brazilian regulations, except for the intended adulterations.

288 Three different classification models (SIMCA, PLS-DA and PLS-DM) were developed, and two approaches were considered:
289 the one-class approach with all three methods, and the multiclass approach with SIMCA and PLS-DA. The one class approach is
290 adequate if the main interest is only to detect whether a sample is adulterated, regardless of the type of the adulterant. For the problem
291 under study, PLS-DA provided excellent results, classifying all samples correctly. SIMCA and PLS-DM produced less satisfactory
292 results, with specificity for the test set of 75% and 68%, respectively.

293 The multiclass approach is the proper choice when the main interest is to investigate the possible presence of known adulterants.
294 It provides more specific information, since in addition to the percentage of samples correctly or incorrectly assigned, information
295 related to the inconclusive assignments is also available. Samples inconclusively classified could be submitted in the sequence to
296 undergo confirmatory analyses. Among the multiclass models, PLS-DA also presented the best performance, with no false-negative
297 predictions, i.e., no adulterated samples were classified as unadulterated. In food fraud analysis, it is essential to avoid false-negative
298 results, since the analyst could declare a sample as unadulterated when it is actually adulterated. For the multiclass approach, the

299 SIMCA model was not able to differentiate unadulterated samples from samples adulterated with cashew. Nonetheless, the apple class
300 was well characterized by SIMCA.

301 Finally, we can suggest this type of application as a potential tool to assist the beverage industry and regulatory organisms in
302 the field of food quality control, allowing detection in fruit nectars through direct, fast and reliable screening analyzes. Further research
303 could be the implementation of the developed classification techniques to detect grape nectar samples adulterated with blends of more
304 than one adulterant.
305

306 **Acknowledgements**

307 The authors acknowledge CAPES for providing the sandwich PhD scholarship (Proc., nº 88881. 132172/2016-01); Prof. Adriana
308 Silva França from the Biofuels Laboratory of UFMG for enabling the use of an ATR-FTIR Spectrophotometer and MSc Andréia H.
309 Suzuki for assistance with this equipment; Giuliano Elias Pereira from EMBRAPA Grape & Wine/Semiarid (Petrolina, Brazil) for
310 providing the Isabel grapes; and Prof. Paolo Oliveri (University of Genova, Italy) for providing the PLS-DM routine and assisting with
311 its use.

312

313 **5. References**

- 314 ABIR (2015). Associação Brasileira das Indústrias de Refrigerantes e Bebidas não Alcoólicas. Available in <https://abir.org.br/>, accessed
315 in March 2018.
- 316 Ammari, F., Redjdal, L., & Rutledge, D. N. (2015). Detection of orange juice frauds using front-face fluorescence spectroscopy and
317 independent components analysis. *Food Chemistry*, *168*, 211-217.
- 318 Asadpoor, M., Ansarin, M., & Nemati, M. (2014). Amino acid profile as a feasible tool for determination of the authenticity of fruit
319 juices. *Advanced Pharmaceutical Bulletin*, *4*, 359.
- 320 Barker, M., & Rayens, W. (2003). Partial least squares for discrimination. *Journal of Chemometrics*, *17*, 166-173.
- 321 Bevilacqua, M., Bucci, R., Magri, A. D., Magri, A. L., Nescatelli, R., & Marini, F. (2013). Classification and class-modelling. In F.
322 Marini (Ed.), *Data handling in science and technology* (Vol. 28, pp. 171-233). Amsterdam: Elsevier.
- 323 Boggia, R., Casolino, M. C., Hysenaj, V., Oliveri, P., & Zunin, P. (2013). A screening method based on UV-Visible spectroscopy and
324 multivariate analysis to assess addition of filler juices and water to pomegranate juices. *Food Chemistry*, *140*, 735-741.
- 325 Bona, E., Marquetti, I., Link, J. V., Makimori, G. Y. F., Arca, V. C., Lemes, A. L. G., Ferreira, J. M. G., Scholz, M. B. S., Valderrama,
326 P., & Poppi, R. J. (2017). Support vector machines in tandem with infrared spectroscopy for geographical classification of
327 green arabica coffee. *LWT-Food Science and Technology*, *76*, 330-336.
- 328 BRASIL. Ministério da Agricultura, Pecuária e Abastecimento. Decreto nº 6.871, de 04 de junho de 2009. Regulamenta a Lei n. 8.918,
329 de 14 de julho de 1994. Dispõe sobre a padronização, a classificação, o registro, a inspeção, a produção e a fiscalização de
330 bebidas.
- 331 Bureau, S., Ścibisz, I., Le Bourvellec, C., & Renard, C. M. (2012). Effect of sample preparation on the measurement of sugars, organic
332 acids, and polyphenols in apple fruit by mid-infrared spectroscopy. *Journal of Agricultural and Food Chemistry*, *60*, 3551-
333 3563.
- 334 Bylesjö, M., Rantalainen, M., Cloarec, O., Nicholson, J. K., Holmes, E., & Trygg, J. (2006). OPLS discriminant analysis: combining
335 the strengths of PLS-DA and SIMCA classification. *Journal of Chemometrics*, *20*, 341-351.
- 336 Callao, M. P., & Ruisánchez, I. (2018). An overview of multivariate qualitative methods for food fraud detection. *Food Control*, *86*,
337 283-293.
- 338 Cuny, M., Vigneau, E., Le Gall, G., Colquhoun, I., Lees, M., & Rutledge, D. (2008). Fruit juice authentication by ¹H NMR spectroscopy
339 in combination with different chemometrics tools. *Analytical and Bioanalytical Chemistry*, *390*, 419-427.

340 Dandage, K., Badia-Melis, R., & Ruiz-García, L. (2017). Indian perspective in food traceability: A review. *Food Control*, 71, 217-227.

341 El Darra, N., Rajha, H. N., Saleh, F., Al-Oweini, R., Maroun, R. G., & Louka, N. (2017). Food fraud detection in commercial
342 pomegranate molasses syrups by UV–VIS spectroscopy, ATR-FTIR spectroscopy and HPLC methods. *Food Control*, 78,
343 132-137.

344 EMBRAPA. (2016). Sistema de Produção do Caju. Retrieved from
345 [https://www.spo.cnptia.embrapa.br/conteudo?p_p_id=conteudoportlet_WAR_sistemasdeproducaof6_lga1ceportlet&p_p_lif](https://www.spo.cnptia.embrapa.br/conteudo?p_p_id=conteudoportlet_WAR_sistemasdeproducaof6_lga1ceportlet&p_p_lifecycle=0&p_p_state=normal&p_p_mode=view&p_p_col_id=column-1&p_p_col_count=1&p_r_p_76293187_sistemaProducaoId=7705&p_r_p_-996514994_topicoId=10320)
346 [ecycle=0&p_p_state=normal&p_p_mode=view&p_p_col_id=column-1&p_p_col_count=1&p_r_p_](https://www.spo.cnptia.embrapa.br/conteudo?p_p_id=conteudoportlet_WAR_sistemasdeproducaof6_lga1ceportlet&p_p_lifecycle=0&p_p_state=normal&p_p_mode=view&p_p_col_id=column-1&p_p_col_count=1&p_r_p_76293187_sistemaProducaoId=7705&p_r_p_-996514994_topicoId=10320)
347 [76293187_sistemaProducaoId=7705&p_r_p_-996514994_topicoId=10320](https://www.spo.cnptia.embrapa.br/conteudo?p_p_id=conteudoportlet_WAR_sistemasdeproducaof6_lga1ceportlet&p_p_lifecycle=0&p_p_state=normal&p_p_mode=view&p_p_col_id=column-1&p_p_col_count=1&p_r_p_76293187_sistemaProducaoId=7705&p_r_p_-996514994_topicoId=10320)

348 Esteki, M., Shahsavari, Z., & Simal-Gandara, J. (2018). Use of spectroscopic methods in combination with linear discriminant
349 analysis for authentication of food products. *Food Control*, 91, 100-112.

350 Georgouli, K., Del Rincon, J. M., & Koidis, A. (2017). Continuous statistical modelling for rapid detection of adulteration of extra
351 virgin olive oil using mid infrared and Raman spectroscopic data. *Food Chemistry*, 217, 735-742.

352 Gondim, C., Junqueira, R. G., Souza, S. V. C., Ruisánchez, I., & Callao, M. P. (2017). Detection of several common adulterants in raw
353 milk by MID-infrared spectroscopy and one-class and multi-class multivariate strategies. *Food Chemistry*, 230, 68-75.

354 He, J., Rodriguez-Saona, L. E., & Giusti, M. M. (2007). Midinfrared spectroscopy for juice authentication rapid differentiation of
355 commercial juices. *Journal of Agricultural and Food Chemistry*, 55, 4443-4452.

356 Kennard, R. W., & Stone, L. A. (1969). Computer aided design of experiments. *Technometrics*, 11, 137-148.

357 López, M. I., Callao, M. P., & Ruisánchez, I. (2015). A tutorial on the validation of qualitative methods: From the univariate to the
358 multivariate approach. *Analytica Chimica Acta*, 891, 62-72.

359 López, M. I., Trullols, E., Callao, M. P., & Ruisánchez, I. (2014). Multivariate screening in food adulteration: Untargeted versus
360 targeted modelling. *Food Chemistry*, 147, 177-181.

361 MAPA (2003). Secretaria de Defesa Agropecuária. Ministério da Agricultura, Pecuária e Abastecimento. Instrução Normativa No. 12,
362 Brazil.

363 MAPA (2013). Secretaria de Defesa Agropecuária. Ministério da Agricultura, Pecuária e Abastecimento. Instrução Normativa No. 42,
364 Brazil.

365 Márquez, C., López, M. I., Ruisánchez, I. & Callao, M. P. (2016). FT-Raman and NIR spectroscopy data fusion strategy for multivariate
366 qualitative analysis of food fraud. *Talanta*, 161, 80-86.

367 Martins, A. R., Talhavini, M., Vieira, M. L., Zacca, J. J., & Braga, J. W. B. (2017). Discrimination of whisky brands and counterfeit
368 identification by UV–Vis spectroscopy and multivariate data analysis. *Food Chemistry*, 229, 142-151.

369 Miaw, C. S. W., Assis, C., Silva, A. R. C. S., Cunha, M. L., Sena, M. M., & de Souza, S. V. C. (2018). Determination of main fruits in
370 adulterated nectars by ATR-FTIR spectroscopy combined with multivariate calibration and variable selection methods. *Food*
371 *Chemistry*, 254, 272-280.

372 Neves, M. F., Trombin, V. G., Lopes, F. F., Kalaki, R., & Milan, P. (2012). *The orange juice business: A Brazilian perspective:*
373 Wageningen Academic Publishers.

374 Oliveri, P. (2017). Class-modelling in food analytical chemistry: Development, sampling, optimisation and validation issues—A tutorial.
375 *Analytica Chimica Acta*, 982, 9-19.

376 Oliveri, P., López, M. I., Casolino, M. C., Ruisánchez, I., Callao, M. P., Medini, L., & Lanteri, S. (2014). Partial least squares density
377 modeling (PLS-DM)—A new class-modeling strategy applied to the authentication of olives in brine by near-infrared
378 spectroscopy. *Analytica Chimica Acta*, 851, 30-36.

379 Paltrinieri, G., & Figuerola, F. (1998). Small-scale processing of native and introduced Amazonian fruits and vegetables. Technical
380 manual. Santiago (Chile): Food and Agriculture Organization (FAO).

381 Pulido, A., Ruisanchez, I., Boqué, R., & Rius, X. F. (2003). Uncertainty of results in routine qualitative analysis. *TrAC Trends in*
382 *Analytical Chemistry*, 22, 647-654.

383 Rinnan, Å., Berg, F. v. d., & Engelsen, S. B. (2009). Review of the most common pre-processing techniques for near-infrared spectra.
384 *TrAC Trends in Analytical Chemistry*, 28, 1201-1222.

385 Ríos-Reina, R., Callejón, R. M., Oliver-Pozo, C., Amigo, J. M., & García-González, D. L. (2017). ATR-FTIR as a potential tool for
386 controlling high quality vinegar categories. *Food Control*, 78, 230-237

387 Rius, A., Callao, M. P., & Rius, F. X. (1997). Multivariate statistical process control applied to sulfate determination by sequential
388 injection analysis. *Analyst*, 122, 737-741.

389 Rodionova, O. Y., Oliveri, P., & Pomerantsev, A. L. (2016). Rigorous and compliant approaches to one-class classification.
390 *Chemometrics and Intelligent Laboratory Systems*, 159, 89-96.

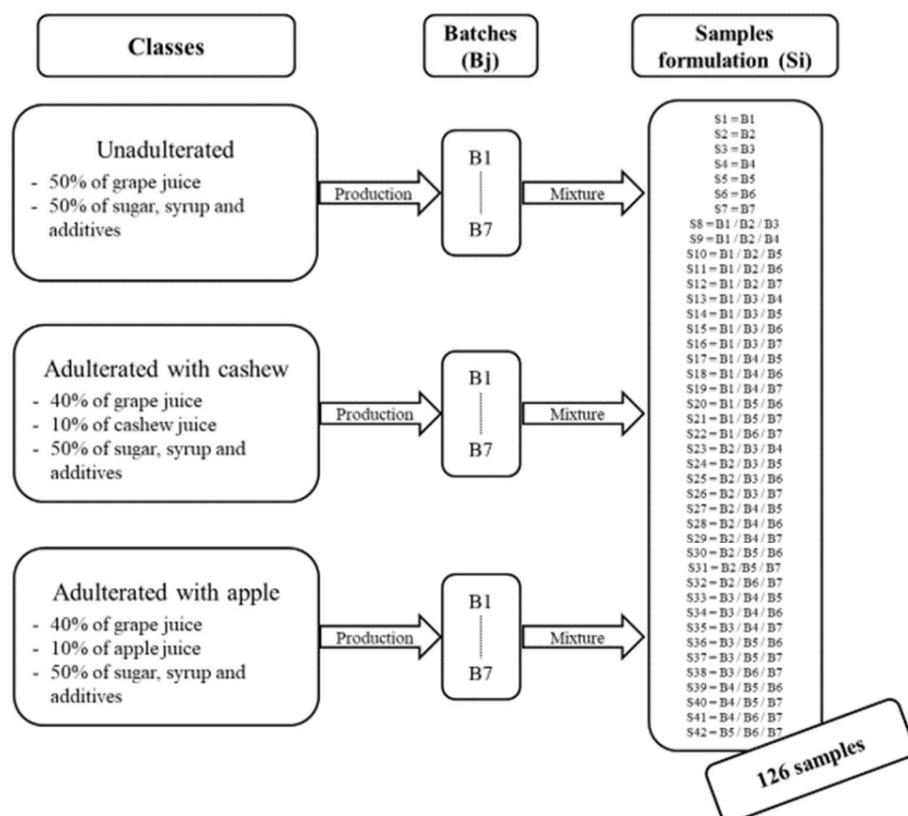
391 Sen, I., & Tokatli, F. (2016). Differentiation of wines with the use of combined data of UV–visible spectra and color characteristics.
392 *Journal of Food Composition and Analysis*, 45, 101-107.

393 Shah, N., Cynkar, W., Smith, P., & Cozzolino, D. (2010). Use of attenuated total reflectance midinfrared for rapid and real-time analysis
394 of compositional parameters in commercial white grape juice. *Journal of Agricultural and Food Chemistry*, 58, 3279-3283.

- 395 Shen, F., Wu, Q., Su, A., TAng, P., ShAo, X., & Liu, B. (2016). Detection of Adulteration in Freshly Squeezed Orange Juice by
396 Electronic Nose and Infrared Spectroscopy. *Czech Journal of Food Science*, 34, 224-232.
- 397 Singhal, R. S., Kulkarni, P. R., & Rege, D. V. (1997). Chapter 3 - Fruit and Vegetable Products. In *Handbook of Indices of Food*
398 *Quality and Authenticity* (pp. 77-130). Oxford: Woodhead Publishing.
- 399 Szymańska, E., Gerretzen, J., Engel, J., Geurts, B., Blanchet, L., & Buydens, L. M. (2015). Chemometrics and qualitative analysis have
400 a vibrant relationship. *TrAC Trends in Analytical Chemistry*, 69, 34-51.
- 401 Tressler, D. K., & Joslyn, M. A. (1961). *Fruit and vegetable juice processing technology*. Westport: AVI Publishing Company.
- 402 Vardin, H., Tay, A., Ozen, B., & Mauer, L. (2008). Authentication of pomegranate juice concentrate using FTIR spectroscopy and
403 chemometrics. *Food Chemistry*, 108, 742-748.
- 404 Xu, L., Fu, H. Y., Yang, T. M., Li, H. D., Cai, C. B., Chen, L. J., & She, Y. B. (2016). Enhanced specificity for detection of frauds by
405 fusion of multi-class and one-class partial least squares discriminant analysis: geographical origins of Chinese shiitake mushroom.
406 *Food Analytical Methods*, 9, 451-458.

407

408 **Figure 1.** Scheme of grape nectar samples formulation.

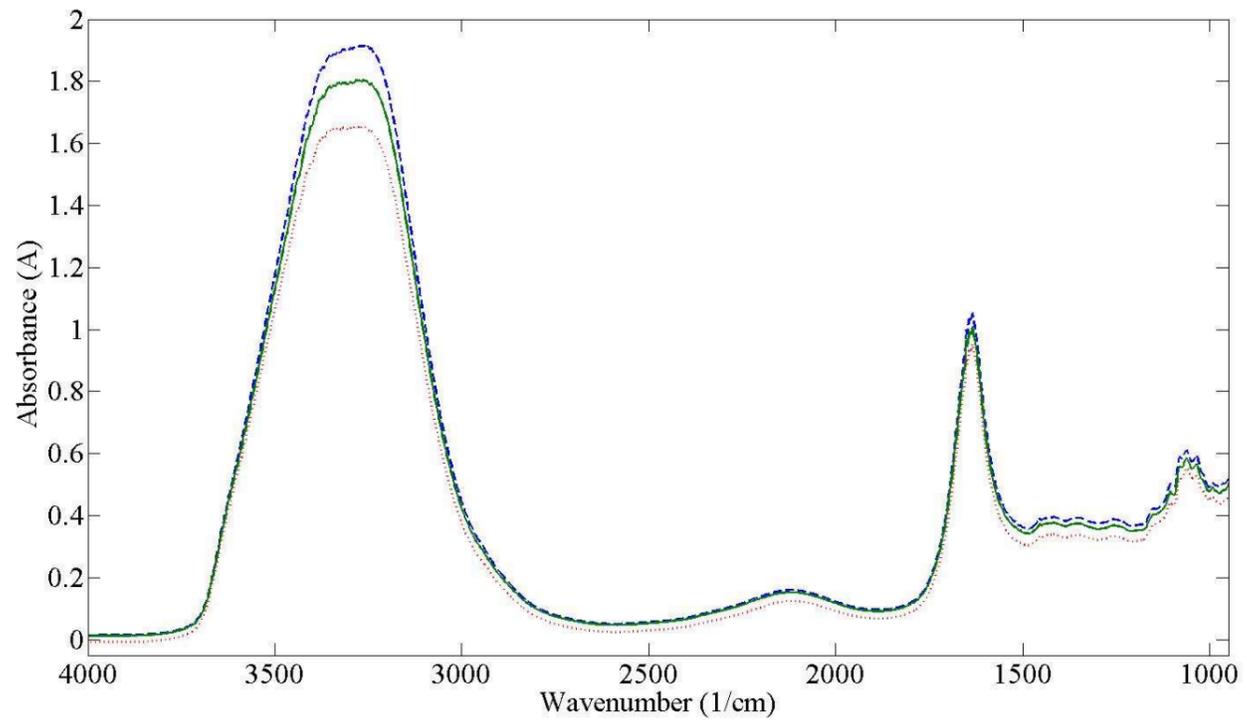


409

410

411

Figure 2. Mean preprocessed spectra of unadulterated class (dashed line), adulterated with cashew class (solid line) and adulterated with apple class (dashed-dot line).

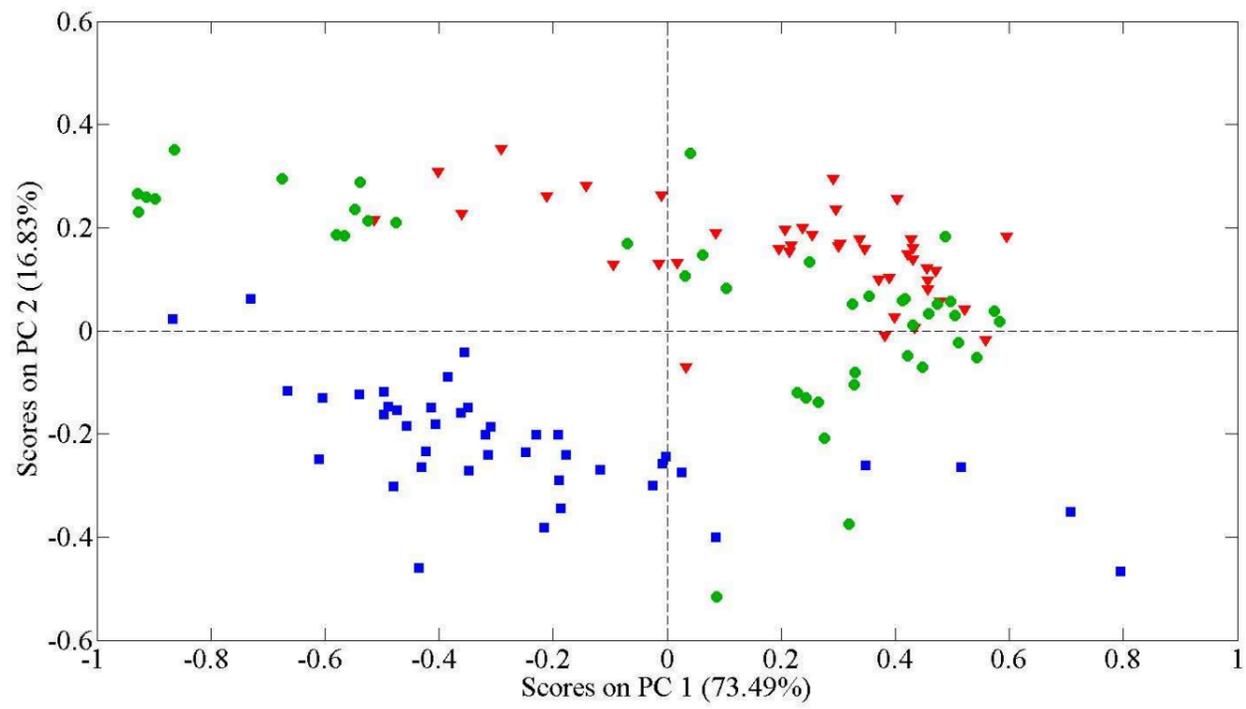


412

413

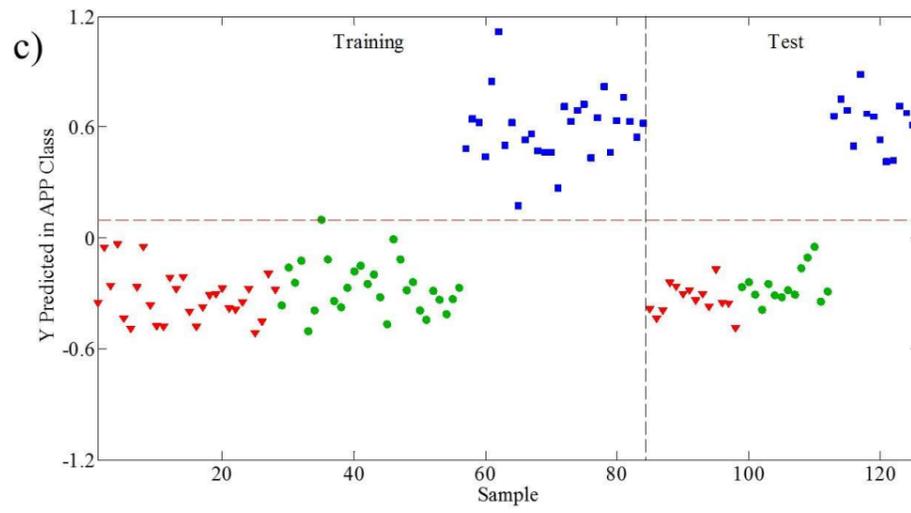
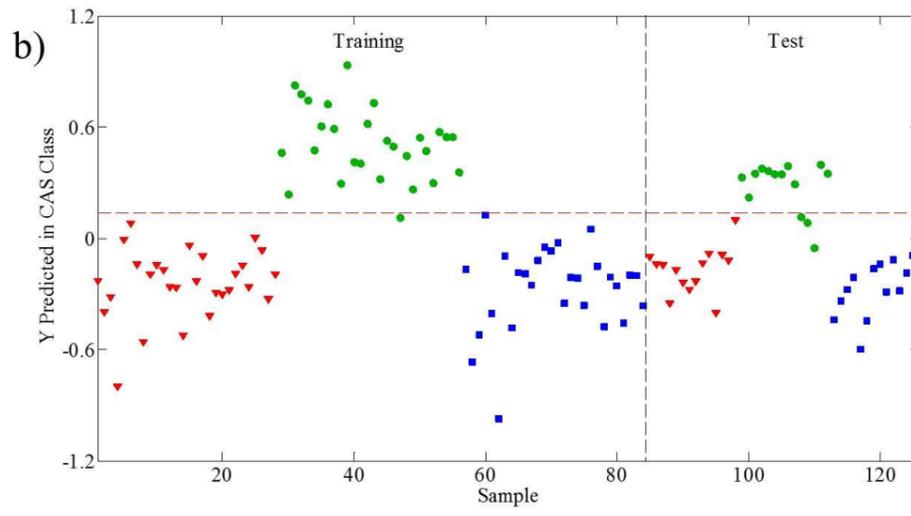
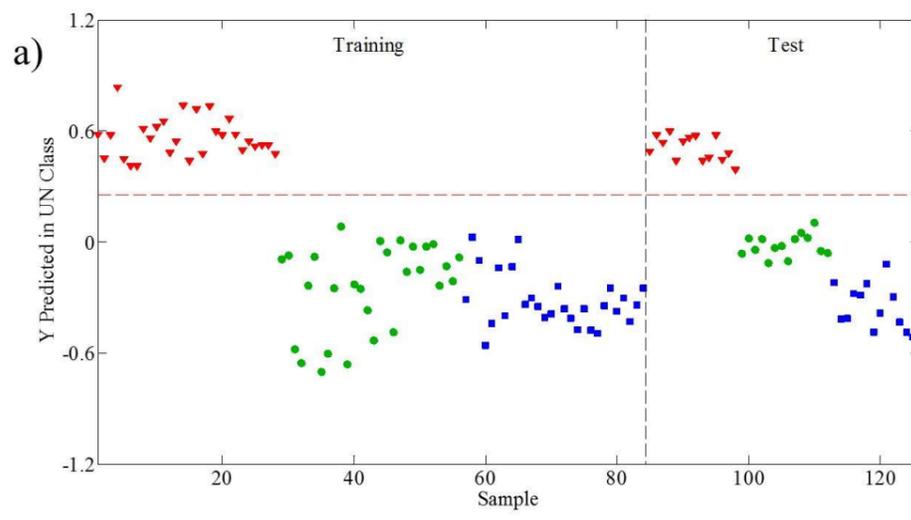
414

Figure 3. Scores of PC1 *versus* PC2 of unadulterated (down triangles), adulterated with cashew (circles) and adulterated with apple (squares) grape nectar samples.



415

416 **Figure 4.** PLS-DA predictions for each class: a) unadulterated (UN), b) adulterated with cashew (CAS) and c) adulterated with
417 apple (APP). Horizontal dashed lines indicate the threshold class and the vertical dashed lines separate training and test samples.
418 Samples symbols: down triangles for unadulterated, circles for adulterated with cashew and squares for adulterated with apple.



420 **Table 1.** SIMCA and PLS-DA multi-class predictions of samples from the unadulterated class (UN), the adulterated with cashew class
 421 (CAS) and the adulterated with apple class (APP) for training and test set. n°S: number of samples; NA: not assigned; MA: multiple
 422 assignments; IR: inconclusive ratio

423

424

425

426

427

428

429

430

431

432

433

434

435

436

437

438

Table 2. Sensitivities and specificities for the one-class strategy

| Method | Set | Class | n°S | Classified as | | | | | |
|--------|----------|-------|-----|---------------|-----|-----|----|----|--------|
| | | | | UN | CAS | APP | NA | MA | IR (%) |
| SIMCA | Training | UN | 28 | 26 | 28 | 0 | 0 | 26 | 92.86 |
| | | CAS | 28 | 5 | 21 | 0 | 7 | 5 | 42.86 |
| | | APP | 28 | 0 | 0 | 22 | 6 | 0 | 21.43 |
| | Test | UN | 14 | 14 | 14 | 0 | 0 | 14 | 50.00 |
| | | CAS | 14 | 7 | 14 | 0 | 0 | 7 | 25.00 |
| | | APP | 14 | 0 | 0 | 14 | 0 | 0 | 0.00 |
| PLS-DA | Training | UN | 28 | 28 | 0 | 0 | 0 | 0 | 0.00 |
| | | CAS | 28 | 0 | 27 | 1 | 1 | 1 | 7.14 |
| | | APP | 28 | 0 | 0 | 28 | 0 | 0 | 0.00 |
| | Test | UN | 14 | 14 | 0 | 0 | 0 | 0 | 0.00 |
| | | CAS | 14 | 0 | 11 | 0 | 3 | 0 | 21.43 |
| | | APP | 14 | 0 | 0 | 14 | 0 | 0 | 0.00 |

439

440

| Method | Set | Sensitivity (%) | Specificity (%) |
|--------|----------|-----------------|-----------------|
| SIMCA | Training | 93 | 91 |
| | Test | 100 | 75 |
| PLS-DA | Training | 100 | 100 |
| | Test | 100 | 100 |
| PLS-DM | Training | 82 | 91 |
| | Test | 100 | 68 |