

# Analytical Methods

Accepted Manuscript



This is an *Accepted Manuscript*, which has been through the Royal Society of Chemistry peer review process and has been accepted for publication.

*Accepted Manuscripts* are published online shortly after acceptance, before technical editing, formatting and proof reading. Using this free service, authors can make their results available to the community, in citable form, before we publish the edited article. We will replace this *Accepted Manuscript* with the edited and formatted *Advance Article* as soon as it is available.

You can find more information about *Accepted Manuscripts* in the [Information for Authors](#).

Please note that technical editing may introduce minor changes to the text and/or graphics, which may alter content. The journal's standard [Terms & Conditions](#) and the [Ethical guidelines](#) still apply. In no event shall the Royal Society of Chemistry be held responsible for any errors or omissions in this *Accepted Manuscript* or any consequences arising from the use of any information it contains.

1  
2  
3 1 **COMPARISON OF DIFFERENT ANALYTICAL CLASSIFICATION SCENARIOS:**  
4  
5 2 **APPLICATION FOR GEOGRAPHICAL ORIGIN OF EDIBLE PALM OIL BY**  
6  
7 3 **STEROLIC (NP)HPLC FINGERPRINTING**  
8  
9 4

10 5 Estefanía PÉREZ-CASTAÑO <sup>a,\*</sup>, Cristina RUIZ-SAMBLÁS <sup>a</sup>, Santiago MEDINA-  
11 6 RODRÍGUEZ <sup>a,b</sup>, Verónica QUIRÓS-RODRÍGUEZ <sup>a</sup>, Ana M. JIMÉNEZ-CARVELO <sup>a</sup>, Lucía  
12 7 VALVERDE-SOM <sup>a</sup>, Antonio GONZÁLEZ-CASADO <sup>a</sup>, Luis CUADROS-RODRÍGUEZ <sup>a</sup>  
13  
14 8

15  
16  
17 9 <sup>a</sup> Department of Analytical Chemistry, University of Granada, c/ Fuentenueva, s.n. E-18071  
18 10 Granada, Spain.

19  
20 11 <sup>b</sup> Department of Signal Theory, Networking and Communications, CITIC-UGR, University of  
21 12 Granada, c/ Periodista Rafael Gómez, E-18071 Granada, Spain.  
22  
23 13  
24  
25 14

26  
27 15 **Abstract**

28  
29 16 This work shows how the best scenario, resulting to apply two chemometric classifiers on  
30 17 different analytical data set from the same sample set, could be chosen according to the  
31 18 classification results. On this way, several classification quality features such as sensitivity  
32 19 (or recall), specificity, positive (or precision) and negative predictive values, Youden index,  
33 20 positive and negative likelihood ratios, F- measure (or F- score), discriminant power,  
34 21 efficiency (or accuracy), AUC (area under the receiver operating curve), Matthews correlation  
35 22 coefficient, Kappa coefficient, overall agreement probability, overall agreement probability  
36 23 from chance and overall Kappa coefficient are described and discussed. As application  
37 24 example, two sterolic chromatographic fingerprints obtained from two different normal-phase  
38 25 HPLC systems are used to discern the geographical origin (South-East Asia, West Africa and  
39 26 South America) of edible palm oil. In each case, two conventional and well-known  
40 27 chemometric classification methods are applied: soft independent modelling by class analogy  
41 28 (SIMCA) and partial least squares-discriminant analysis (PLS-DA).  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52

53 30 **Keywords**

54 31 Classification scenario comparison; liquid chromatography sterolic fingerprints; food  
55 32 authenticity; palm oil.  
56  
57 33  
58 34

59 35 \* Corresponding author: phone: +34 958240797; fax: +34 958243328; email: [stefani@ugr.es](mailto:stefani@ugr.es).  
60

## 1. INTRODUCTION

Palm oil is one of the most consumed edible oils in the world for its price and its properties. It is obtained from the fruit of the palm (*Elaeis guineensis*). The plant is native from Western Africa, later extended to South America in the XVI century, and more recently, in the XIX century, it was introduced in Eastern Asia from America. Crude palm oil is a semi-solid fat at room temperature which has a high stability compared with other vegetable oils<sup>1,2,3</sup>. Knowing the geographical origin of the palm oil and the traceability of palm oil supply chain<sup>4</sup> is interesting from an environmental point of view, because the palm oil exploitation is cause of deforestation and loss of biodiversity in some tropical countries<sup>5</sup>. Furthermore, an indication of the geographical origin of food products is an increasingly important consumer demand on food labelling because consumers consider it to be an added-value to the product. Other possible reasons of this growing interest are: specific culinary, organoleptic qualities, or purported health benefits associated with regional products<sup>6</sup>.

For the characterization and verification of the geographical origin of palm oils, as well as other food products, it is necessary to find specific qualities derived from its place of production (markers) and which are subject to specific local factors such as climate and terrain<sup>7</sup>. An example of this approach is the use of molecular markers as the DNA fingerprinting. Alternatively chemical markers could be also used. In literature sterol profiles and total sterol contents have been used as tools to value the oil authenticity<sup>8,9,10,11</sup> because each vegetal species has a characteristic compositional profile of sterols<sup>12,13</sup>. This suggests that sterol profiles might be suitable candidates to develop an analytical tool to verify the geographical origin of palm oil.

Phytosterols are a group of bioactive compounds, with a derived cyclopentane-perhydrophenanthrene four-ring molecular structure. These compounds are present in plants and they are differentiated by the number of carbon atoms of the side chain, and by the nature of the same. There are three classes of sterols: 4-des, 4-mono and 4,4'-dimethylsterols that could be found in free form or esterified with fatty acids and other conjugates<sup>14,15</sup>. Crude palm oil contains about 0.7-0.8 g/kg of total sterols, however it must be taken into mind that the refining processes affect the concentration and the compositional profile of sterols because occur hydrolysis and oxidation processes that destroy sterols.

There are different methods and techniques (chromatographic and non-chromatographic) for the analysis of sterols in vegetable oil. Chromatographic methods are the most commonly used against the non-chromatographic ones because much more information is obtained about the sterols composition present in the sample<sup>15,16,17,18,19</sup>.

Depending on the information that is wanted, different sample preparation steps could be applied. Generally, phytosterol analysis includes: an extraction of the lipid fraction, acid

1  
2  
3 73 hydrolysis or basic hydrolysis (saponification) to release phytosterols, an extraction of the  
4 74 unsaponifiable fraction, the separation or partial purification of sterols and finally a  
5 75 chromatographic analysis. However, in some specific cases, in addition to the above steps,  
6 76 the formation of derivatives of phytosterols is previously required to chromatographic  
7 77 analysis. The routine methods employ a saponification reaction although it could be replaced  
8 78 by a transesterification reaction with similar results<sup>20</sup>.

9  
10 79 Analytical liquid chromatography has not been much applied for determining the sterol  
11 80 composition of vegetable oils but it is usually performed in reverse phase<sup>16</sup>. Normal-phase  
12 81 HPLC methods can also be used for the absolute quantification of the total amount of  
13 82 phytosterols but these methods show a poor chromatographic resolution and do not provide  
14 83 precise information on the sterol composition<sup>16</sup>. In addition, several conventional detection  
15 84 methods, such as ultraviolet absorption, refractive index, and evaporative light-scattering  
16 85 have been applied. In 2004 the Corona Charged Aerosol detector (CAD) was developed as  
17 86 an alternative, which has ability to accurately measure a wide range of analytes<sup>21</sup>.  
18 87 Nonetheless, only a method has been just reported using CAD for the determination of  
19 88 sterols in vegetable oils<sup>22</sup>.

20 89 If the chromatographic conditions are properly optimized and a suitable detector is coupled,  
21 90 the yielded chromatogram contains specific and relevant information about the considered  
22 91 product, which could be used for authentication purposes. When the chromatogram is well  
23 92 resolved, the information for each chemical component could be extracted from each peak.  
24 93 Instead, if the chromatogram is shaped on a broad and comprehensive band, the intrinsic  
25 94 information is not evident and the chromatographic fingerprinting methodology should be  
26 95 then applied<sup>23,24</sup>. By reaching that point, the application of multivariate chemometric tools is  
27 96 required to extract the useful information from the chromatographic raw data<sup>25,26,27</sup>. Some  
28 97 examples of the use of chromatographic fingerprints merging with classical chemometric  
29 98 methods have been recently reported by our research group with satisfactory results in the  
30 99 authentication of vegetable edible oils<sup>28,29,30,31</sup>. In addition, data mining classification methods  
31 100 have also been used for edible vegetable oils authenticity applications<sup>32</sup>.

32 101 In the last few years, some papers have been published using classic supervised pattern  
33 102 recognition methods for vegetable oil classification and authentication based on their sterol  
34 103 composition. In most of them, the data matrices are made up from the sterol contents  
35 104 (concentrations or compositional data)<sup>8,10,33</sup> and only a work set the data matrix from the  
36 105 sterolic chromatographic fingerprint<sup>11</sup>. In the same way, there is not enough background  
37 106 about the authentication of palm oil by using this methodology and, as far as we know, only  
38 107 two studies have been published, but they use fingerprintings from volatile compounds<sup>34</sup> and  
39 108 triacylglycerols<sup>35</sup>.

1  
2  
3 109 In this study, two sterolic chromatographic fingerprints obtained from two different normal-  
4  
5 110 phase HPLC systems are used to discern the geographical origin (South-East Asia, West  
6  
7 111 Africa and South America) of edible palm oil. The aim of this paper is to show how the  
8  
9 112 obtained results from different classification scenarios would be compared in order to select,  
10  
11 113 if possible, the best combination of classification chemometric methods and/or measured  
12  
13 114 analytical data set. The comparison is based on different quality classification metrics which  
14  
15 115 are defined in this work, such as sensitivity (or recall), specificity, positive (or precision) and  
16  
17 116 negative predictive values, Youden index, positive and negative likelihood ratios, F-measure  
18  
19 117 (or F-score), discriminant power, efficiency (or accuracy), AUC (area under the receiver  
20  
21 118 operating curve), Matthews correlation coefficient, Kappa coefficient, overall agreement  
22  
23 119 probability, overall agreement probability from chance and overall Kappa coefficient.  
24  
25 120  
26  
27 121

## 24 122 **2. MATERIALS AND METHODS**

### 27 124 **2.1. Instrumentation**

28  
29 125 The analyses were performed using two different HPLC systems. The first one was a Konik  
30  
31 126 Model 560 (Konik-Tech, Sant Cugat del Valle, Barcelona, Spain) with a quaternary pump, a  
32  
33 127 column oven, an autosampler with a 20  $\mu$ L loop, and an UV-Vis detector.

34 128 The second one was an Agilent 1100 Series (Agilent Technologies, Santa Clara, CA, USA)  
35  
36 129 with a quaternary pump, degasser, autosampler and thermostatted HPLC column  
37  
38 130 compartment Eppendorf CH-30 (Eppendorf, Hamburg, Germany). Detection was carried out  
39  
40 131 with a Corona CAD (ESA Biosciences Inc., Chelmsford, MA, USA).  
41  
42 132

### 42 133 **2.2. Chemicals**

43  
44 134 Solid standards of stigmasterol, campesterol and cholestanol (internal standard, IS) were  
45  
46 135 provided from Sigma-Aldrich (Steinheim, Germany) and  $\beta$ -amyirin was from  
47  
48 136 EXTRASYNTHESE (Genay, France).

49 137 Sodium methoxide, citric acid monohydrate and anhydrous sodium sulphate were provided  
50  
51 138 from Alfa-Aesar (Karlsruhe, Germany), Sigma-Aldrich (Steinheim, Germany), and Panreac  
52  
53 139 Quimica (Barcelona, Spain) respectively. The solvents employed (n-hexane, 2-propanol,  
54  
55 140 methanol and methyl tert-butyl ether (MTBE), BDH Prolabo, HPLC grade), were purchased  
56  
57 141 from VWR International (Madrid, Spain). All the aqueous solutions were prepared with Milli-  
58  
59 142 Q deionized water (Millipore, Bedford, MA).

60 143 The nitrogen (99.9999%) used for CAD detector was provided from Air Liquid (Madrid,  
144  
145 Spain).

### 2.3. Samples

A total of 102 crude palm oil samples were supplied from RIKILT-Institute of Food Safety Wageningen University, (Wageningen, The Netherland). The samples coming from the main continents of palm oil production: South-East Asia (56 samples from Malaysia, Indonesia, Papua New Guinea and Salomon Islands), West Africa (30 samples, from Ghana, Guinea, Cote d'Ivoire, Nigeria and Cameroon) and South America (16 samples from Brazil). Table 1 shows in detail the origin of the different samples tested.

TABLE 1
---------

### 2.4. Sample preparation

Prior to chromatographic analysis, a methylation reaction was applied on the palm oil samples. This reaction replaces the usual saponification and isolation processes, and it has the advantage of being less time-consuming and requiring less sample amount. The applied procedure is similar to the one previously described by Biederman<sup>20</sup> and Kamm<sup>36</sup>. The transesterified sample solutions were frozen (-20 °C) and kept in the dark until analysis. Just before the chromatographic analysis, 500 µL of this solution was added in a 2 mL HPLC vial, and then 120 µL of 0.05% (w/w) cholestanol solution in n-hexane was added as control internal standard. Finally the mixture was diluted with 1000 µL of n-hexane. The vial was sealed and vortexed for 20 s. This solution was prepared just for analysis.

### 2.5. LC Conditions

HPLC analysis is carried out on a (250 x 4 mm i.d., 5 µm) column Lichrospher® 100 CN maintained at 25 °C. The composition of the eluent was n-hexane/2-propanol (99:1, v/v) at a flow rate of 1.2 mL/min and a run time of 20 min. No gradient was applied. UV detection (Konik equipment) was performed at 202 nm. For CAD monitoring (Agilent equipment), a 100 pA output range was used and nitrogen gas pressure was adjusted to 35 psi. Chromatographic data handling were performed by a Konikrom Plus software (version 3.0.5) for HPLC Konik, and ChemStation software (version A.10.02) for HPLC Agilent.

### 2.6. Chemometrics

The raw data files for each chromatogram were exported in a CSV file (*comma-separated values*) from the instrument software to the MATLAB environment (version 7.8, R2009a, The Mathworks Inc. MA, USA). Initially, each chromatographic fingerprinting is coded in a two-data vectors (time/intensity) with 4500 (UV-Vis Konik) and 2400 (CAD ESA) elements (variables), depending of the data acquisition rate of each HPLC detector.

1  
2  
3 181 All the intensity data vectors, one for each oil sample, from the same chromatographic  
4 182 detector are merged in a single data matrix (X- block matrix) composed of 102 rows (palm oil  
5 183 samples) and a certain number of columns (variables) that varies depending on the  
6 184 measuring time and the rate of acquisition of data of each HPLC detector, as it has been  
7 185 described in the previous paragraph. The elements of X-matrix are the intensities values of  
8 186 the chromatographic signals. In addition, a new column is added to each data matrix  
9 187 specifying the class (geographical continent) of each sample, titled by an alphanumeric code,  
10 188 for example "AF", AM" and "AS" for Africa, America and Africa. This column set up the Y-  
11 189 block of the data matrix.

12  
13 190 Next, a preprocessing of each X-matrix was carried out using a home-made MATLAB  
14 191 function, named "MEDINA" (version 07). This function makes use of some of the functions  
15 192 contained in MATLAB Bioinformatics Toolbox™ software to improve the quality of raw  
16 193 chromatographic data, and also the "icoshift" algorithm (version 1.2) for solving signal  
17 194 alignment problems in chromatographic data<sup>37</sup>. Basically, the chromatographic data  
18 195 processing consists of the following stages (see Supplementary Information for details): (1)  
19 196 selection of the interval of interest from chromatograms; (2) decimation of the raw  
20 197 chromatographic data. It makes possible to resample the signal into a more manageable  
21 198 chromatographic data vector, preserving the information contained in the chromatogram (in  
22 199 this case, a decimation factor of 2 was used); (3) de-noising and smoothing of the  
23 200 chromatographic signal using a least-squares digital polynomial filter (i.e., a Savitzky-Golay  
24 201 filter); (4) baseline correction using the "msbackadj" function (available in the above  
25 202 mentioned MATLAB toolbox); (5) alignment of the chromatographic profiles with the "icoshift"  
26 203 algorithm. Finally, a mean centring of the chromatographic data matrix was applied (i.e., the  
27 204 subtraction of the mean from each data vector) prior to the statistical analysis. Once the  
28 205 chromatographic data preprocessing was carried out, it was then possible to use  
29 206 classification and statistical learning tools to create classifiers.

30 207 For multivariate chemometric pattern recognition PLS\_Toolbox (version 7.5.2, Eigenvector  
31 208 Research, Wenatchee, WA) was used. The performance features of each classifier,  
32 209 described in the next section, were calculated on the validation test from a home-designed  
33 210 MS Excel™ spreadsheet (version 14.0, 2010).

34 211  
35 212 *Exploratory analysis and classification methods*

36 213 Principal components analysis (PCA) is a type of exploratory data non-supervised analysis  
37 214 which can be applied to any X-matrix<sup>38,39,40</sup>. The main aim of PCA is the dimension reduction  
38 215 when the variables are correlated. A few new variables are defined, named principal  
39 216 components (PCs), as linear combination of the original variables in order to explain as much  
40 217 variability as possible with the smallest number of PCs.

Two habitually classification methods were then applied. A venetian blinds object out cross-validation procedure was adopted to optimize all the built models.

Soft independent modelling by class analogy (SIMCA) is a well-known class-modeling classification method based on principal component analysis<sup>27, 41</sup>. Each class is independently modelled by a PCA so that each model defines the boundary regions for each class. The number of principal components of each category was determined using the rule of thumb based on the cross-validation, which gives the model optimal prediction properties. The unknown samples are applied to the model of prediction; they are compared to the defined classes and assigned to a class according to their similarity (analogy). In this study, the recognition is made based on the distance  $d_{i,C}$  of each  $i$ -sample from each  $C$ -class. This is calculated by applying the following equation:

$$d_{i,C} = \sqrt{\left(\frac{Q_{i,C}}{Q_{C(0.95)}}\right)^2 + \left(\frac{T_{i,C}^2}{T_{C(0.95)}^2}\right)^2}$$

where  $Q_{i,C}$  and  $T_{i,C}^2$  are the computed statistics  $Q$ -residuals and the  $T^2$ -Hotelling respectively, calculated from the corresponding  $C$ -class PCA model, and  $Q_{C(0.95)}$  and  $T_{C(0.95)}^2$  are the values for a 95% confidence level. The chosen classification threshold<sup>42</sup> was  $d_{i,C} = \sqrt{2}$ . For the class assignment of a sample, the calculated distance to such class have to be lesser than or equal to  $\sqrt{2}$ . On the contrary, if the distance is always larger, the sample is unclassified (class no assigned). If a sample is simultaneously assigned to more of two classes (because the distance to both ones was lesser than or equal to  $\sqrt{2}$ ), the sample will be definitively assigned to the class whose distance value is lesser.

Partial least squares-discriminant analysis (PLS-DA) is a linear discrimination method based upon the classical PLS regression method<sup>43</sup> for building predictive models. The goal of PLS regression is to provide dimensionality reduction in an application where the response variable ( $Y$ -block) is related to the predictor variables ( $X$ -block). The used software is only able to perform binary classifications. So, for  $n$ -class classification is necessary build  $n$  two-class (binary) models<sup>44</sup>.

PLS-DA is applied to develop a model that predicts the representative class value (between 0 and 1) for each sample in each classification. To make a class assignment, the discrimination thresholds and the probability of a sample belonging to a specific class were calculated based on a Bayesian approach. The unknown samples will be correctly classified always than the assigned class value is equal or greater than the threshold value; otherwise the sample will be unclassified. A sample could be classified into two classes if it has a predicted value greater than the threshold value in both classifications. In this case, the sample will be assigned to the class whose predicted value is closer to 1.

1  
2  
3 252 *External validation of classifiers*

4  
5 253 In order to apply a proper external validation of the classification/prediction models, the  
6  
7 254 original data set was divided into two data sets: (1) a training set, used to establish the  
8  
9 255 chemometric models; and (2) a validation (or testing) set, in order to test the validity of the  
10  
11 256 models. Approximately 30% of the samples from each class were randomly chosen to  
12  
13 257 constitute the validation set. Table 2 shows, more specifically, the composition of the  
14  
15 258 samples that were used in in both sets of calibration and test.  
16  
17 259

18  
19 260  
20 261  
21 262  
22 263  
23 264  
24 265  
25 266  
26 267  
27 268  
28 269  
29 270  
30 271  
31 272  
32 273  
33 274  
34 275  
35 276  
36 277  
37 278  
38 279  
39 280  
40 281  
41 282  
42 283  
43 284

TABLE 2
---------

44 285  
45 286  
46 287  
47 288  
48 289  
49 290  
50 291  
51 292  
52 293  
53 294  
54 295  
55 296  
56 297  
57 298  
58 299  
59 300  
60 301

### 3. BACKGROUND: QUALITY PERFORMANCE OF CLASSIFICATION SCENARIOS

263 The empirical evaluation of classification scenarios is a matter of on-going debate between  
264 researchers where classification scenario is referred to the combination of classifier and  
265 analytical data in a particular case. The assessment of the quality classification performance  
266 without focusing on a class is the most general way of comparing the quality of the  
267 classification results. In order to quantify this quality, several performance features have  
268 been proposed as metrics<sup>45,46,47</sup>. The estimation of such metrics is based on measuring the  
269 classifier's ability to distinguish classes and, consequently, to avoid failure in classification.  
270 Although most performance features in use today focus on a classifier's ability to identify  
271 classes correctly, in certain cases, other properties such as failure avoidance or class  
272 discrimination may also be useful<sup>48</sup>.  
273 Quality features for classification are built from a contingency table which records correctly  
274 and incorrectly assigned examples for each class. The corresponding ternary contingency  
275 table is shown in Table 3.

46  
47 276  
48 277  
49 278  
50 279  
51 280  
52 281  
53 282  
54 283  
55 284  
56 285  
57 286  
58 287  
59 288  
60 289

TABLE 3
---------

278 The quality performance features of the different classifiers are calculated by reducing the  
279 ternary contingency table to three binary contingency tables<sup>49</sup> because the quality  
280 parameters are described to binary classification. A binary contingency table is a square of  
281 2x2 where the rows represent the number of classifier predictions and the columns are the  
282 actual value of class. Table 4 presents a standard contingency table for binary classification.

58  
59 284  
60 285

TABLE 4
---------

285 The final value is obtained from the average of the corresponding features obtained for  
 286 binary classifiers, weighting with respect to the number of samples in each class.

287 The different quality metrics used in this paper for evaluating the classification results are  
 288 shown below:

289  
 290 **Sensibility (SENS) (or recall).** It indicates the probability of classifying a sample as positive  
 291 really, *i.e.*, the confidence in a positive result for a sample of the label class is obtained. The  
 292 range of values for this feature is 0 to 1.

$$\text{SENS} = \frac{\text{TP}}{\text{TP} + \text{FN}} = \frac{\text{TP}}{\text{TL}}$$

293  
 294 **Specificity (SPEC).** It indicates the probability of classifying a sample as negative really, *i.e.*,  
 295 the confidence that a negative result for a sample of non-label class is obtained. It is also  
 296 ranged between 0 and 1.

$$\text{SPEC} = \frac{\text{TN}}{\text{TN} + \text{FP}} = \frac{\text{TN}}{\text{TnL}}$$

297 Sensitivity and specificity assess the effectiveness of the classifier on a single class, positive  
 298 and negative respectively.

299  
 300 **Positive predictive value (PPV) (or precision).** It estimates the predictive power of the  
 301 classifier; this metric quantifies the precision of the classifier to identify examples of a given  
 302 class. PPV measures the proportion of correctly assigned positive examples and its value  
 303 varies between 0 and 1.

$$\text{PPV} = \frac{\text{TP}}{\text{TP} + \text{FP}} = \frac{\text{TP}}{\text{AP}}$$

304  
 305 **Negative predictive value (NPV).** The complement of PPV in this context appears in the  
 306 form of the negative predictive value (NPV), which measures the proportion of correctly  
 307 assigned negative examples. The range of values is also between 0 and 1.

$$\text{NPV} = \frac{\text{TN}}{\text{TN} + \text{FN}} = \frac{\text{TN}}{\text{AN}}$$

308  
 309 **Youden's index (YOU).** It evaluates the classifier's ability to avoid failure; it is derived from  
 310 sensitivity and specificity. This parameter varies between 0 and 1.

$$\text{YOU} = \text{SENS} - (1 - \text{SPEC})$$

311

1  
2  
3 312 **The likelihood ratios (LR).** It is possible to distinguish between positive likelihood ratio,  
4 313 LR(+), and negative likelihood ratio, LR(-). The positive likelihood ratio represents the ratio  
5 314 between the probability to predict an example as positive when it is truly positive, and the  
6 315 probability to predict an example as positive when actually it is not positive:

$$LR(+) = \frac{SENS}{1 - SPEC}$$

10 316  
11  
12  
13 317 while the negative likelihood ratio is the ratio between the probabilities to predict an example  
14 318 as negative when it is actually positive, and the probability to predict an example as negative  
15 319 when it is truly negative:

$$LR(-) = \frac{1 - SENS}{SPEC}$$

19 320 Higher positive likelihood ratio and a lower negative likelihood ratio mean better performance  
20 321 on positive and negative classes respectively.

22 322  
23 323 **F-measure (F).** It is defined as the harmonic mean of precision and sensibility. It is a  
24 324 composite feature which benefits classifiers with higher sensitivity and challenges classifiers  
25 325 with higher specificity. This metric ranges between 0 and 1.

$$F = 2 \times \frac{SENS \times PPV}{SENS + PPV}$$

27 326  
28 327 **Discriminant power (DP).** It does exactly what its name implies: *i.e.*, it assesses how well a  
29 328 classifier distinguishes between positive and negative examples.

$$DP = \frac{\sqrt{3}}{\pi} \left( \log \frac{SENS}{1 - SENS} + \log \frac{SPEC}{1 - SPEC} \right)$$

31 329  
32 330 **Efficiency (EFFIC) (or accuracy).** The most common metric for classifier evaluation, it  
33 331 assesses the overall effectiveness of the classifier by estimating the probability of the true  
34 332 value of the class label. The EFFIC values are included between 0 and 1.

$$EFFIC = \frac{TP + TN}{TP + FP + FN + TN} = \frac{TP + TN}{T}$$

35 333  
36 334 **Area under the ROC curve (AUC) (or correct classification rate).** The area under the  
37 335 ROC (Receiver Operating Characteristic) curve is a summary indicator of ROC curve quality  
38 336 that can summarize the performance of a classifier into a single metric. Graphically, ROC is  
39 337 plotted as a curve that gives the true positive rate as a function of false positive rate for the  
40 338 same group. AUC is a measure of the ability of the classifier to avoid errors during

339 classification. The AUC varies between 0 and 1 although, in practice, its values should be  
340 larger than 0.5.

$$\text{AUC} = \frac{\text{SENS} + \text{SPEC}}{2}$$

341  
342 **Matthews correlation coefficient (MCC).** It measures the overall quality of a method  
343 classification since it considers mutually accuracies and error rates, and involve all values of  
344 the contingency table.

$$\text{MCC} = \frac{\text{TP} \times \text{TN} - \text{FP} \times \text{FN}}{\sqrt{(\text{TP} + \text{FP}) \times (\text{FN} + \text{TN}) \times (\text{TP} + \text{FN}) \times (\text{FP} + \text{TN})}}$$

345 MCC ranges from 1 for a perfect prediction to -1 for the worst possible prediction. MCC close  
346 to 0 indicates a model that performs randomly.

347  
348 **Kappa coefficient (K).** It indicates the proportion of agreement after the chance agreement  
349 is removed from consideration<sup>50</sup>. It calculated from a probability rate where the numerator is  
350 the percent of units in which beyond-chance agreement occurred, and the denominator is the  
351 percent of subjects for which one would not expect any agreement by chance.

$$K = \frac{P_a - P_c}{1 - P_c}$$

352 In this equation,  $P_a$  is the probability term from agreement and  $P_c$  is the probability term from  
353 chance.

$$P_a = \frac{\text{TP} + \text{TN}}{T} = \text{EFFIC} \quad ; \quad P_c = \frac{\text{AP} \times \text{TL} + \text{AN} \times \text{TnL}}{T^2}$$

354 The values of K are ranged between 1 (the classifiers are in complete agreement) and 0  
355 (there is no agreement among the classifiers other than what would be expected by chance,  
356 as defined by  $P_c$ ).

357  
358 An overall value for any of the parameters, which have been previously defined, could also  
359 be directly calculated as it is explained below.

360  
361 **Overall agreement probability (overall  $P_a$ ).**

$$\text{overall } P_a = \frac{\sum a_i}{T} = \frac{\sum(\text{TP} + \text{TN})}{T}$$

362  
363 **Overall chance agreement probability (overall  $P_c$ ).**

$$\text{overall } P_c = \frac{\sum(TC_i \times AC_i)}{T^2}$$

364

365 **Overall kappa coefficient (overall K).**

$$\text{overall } K = \frac{\text{overall } P_a - \text{overall } P_c}{1 - \text{overall } P_c}$$

366

367

#### 368 4. RESULTS AND DISCUSSION

369 An example of the chromatograms obtained from the same palm oil sample from Africa by  
370 both HPLC systems, is shown in Fig. 1. As it can be observed, the chromatograms are split  
371 in four regions. In order to identify the regions corresponding to the sterolic fraction, three  
372 aliquots of this palm oil sample were fortified each one with a representative sterol standard:  
373  $\beta$ -amyirin (a dimethylsterol), stigmasterol and campesterol (two desmethylsterols). Next they  
374 were analysed by applying the two analytical chromatographic methods. By inspecting where  
375 the height is increased, and in accordance with the assignment carried out by Biedermann<sup>20</sup>,  
376 it could be concluded that the sterols are divided in regions II and III: the region II is  
377 associated to the dimethylsterols while the region III contains the methylsterols and  
378 desmethylsterols and, possibly, other compounds as the fatty alcohols. The region I was not  
379 assigned although the large peak should be probably due to the fatty acids methyl esters,  
380 whereas region IV would be due to the terpenic alcohols.

381

FIGURE 1

382

383 As previous exploratory analysis, PCA was performed on the X-matrix in order to perceive  
384 similar, dissimilar, typical, or outlier samples. Two PCs were enough to explain 88.8% and  
385 90.6% of the cumulative variance from the HPLC-UV and HPLC-CAD fingerprint data,  
386 respectively. Both PC1-PC2 scores and PC1 loading plots for each X-matrix are shown in  
387 Figure 2.

388

FIGURE 2

389

390 Both scores plots allow to distinguish two groups separated on the first principal component  
391 which are correlated with the AMERICA (left) and ASIA samples (right). On the other hand,  
392 the AFRICA samples are not grouped and they do not show any trend but they are dispersed  
393 on the plotted space. This fact shows that the sterolic chromatographic profiles from AFRICA

1  
2  
3 394 samples have not a specific pattern and some of them are similar to the ones from samples  
4  
5 395 from AMERICA or ASIA. This fact could be explained by the African common origin of all  
6  
7 396 palm oils. The PC2 does not provide information about the geographical origin.

8 397 The PC1 loadings plot shows a profile that coincides with chromatographic region  
9  
10 398 corresponding to region III of the chromatogram, associated to the  $\Delta^5$ - and  
11 399  $\Delta^7$ -desmethysterols (see Figure 1). Therefore, this region contains the significant information  
12  
13 400 about the geographical origin of the palm oils and it will be selected for building the  
14  
15 401 classification models.

16 402 In order to apply a SIMCA classification, three PC models were built, from the corresponding  
17  
18 403 training set samples, for each class (AFRICA, AMERICA and ASIA). The number of chosen  
19  
20 404 PCs for each model was respectively 3, 5 and 5 from HPLC-UV fingerprint data, and 3, 4 and  
21  
22 405 4 from HPLC-CAD fingerprint data. In all cases, the percentage of explained variance was  
23  
24 406 higher than 96%.

25 407 In a similar way, the three-class PLS-DA model was trained. The number of latent variables  
26  
27 408 (LVs) chosen for each model was respectively 5 from HPLC-UV fingerprint data, and 3 from  
28  
29 409 HPLC-CAD fingerprint data, with percentages of explained variance for the X-block and  
30  
31 410 Y-block of the data matrix of 94% and 49% for the first model, and 94% and 43% for the  
32  
33 411 second one, respectively.

34 412 Once the classification models are defined, the more probable class is assigned to each  
35  
36 413 sample of the validation set. The contingency tables showing the results of the assignment  
37  
38 414 from each classifier are shown in Table 5.

39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
415  
416  
417  
418  
419  
420  
421  
422  
423  
424  
425  
426  
427

TABLE 5
---------

417 In general, the samples from AFRICA and ASIA are better classified than the samples from  
418 AMERICA. By comparing the results from the two HPLC fingerprinting, it seems that SIMCA  
419 classifies better than PLS-DA. However, only by having a look, it would be difficult to decide  
420 the best classifier. In order to have a set of appropriate metrics for making this decision,  
421 Table 6 collects the pooled performance features the four classification scenarios; each  
422 value has been calculated from the three reduced binary contingency tables by weighing the  
423 number of actual samples from each class. The empty cells are the consequence of the zero  
424 value obtained for some classification rates in the contingency tables. In addition Table 6  
425 collects the overall values of quality performance features (above mentioned).

TABLE 6
---------

1  
2  
3 428 The classifiers can now be arranged in decreasing order of performance as:

4 429 SIMCA(HPLC-UV) > SIMCA(HPLC-CAD) = PLS-DA(HPLC-CAD) > PLS-DA(HPLC-UV)

5 430 This ranking is easily established dealing the main features related to the overall  
6 431 classification reliability: efficiency, AUC and Kappa coefficient. These performance features  
7 432 are the best indicators of the classification ability, although practically the same ranking could  
8 433 be obtained starting from anyone of the tabulated features.

9 434 However, in strictly technical terms, none of the tested classifiers show an enough assurance  
10 435 as to be applied in order to discern the geographical origin of any sample of palm oil since, in  
11 436 the best case, one of each six-seven samples (15%) would be erroneously classified.

12 437

13 438

## 14 439 **5. CONCLUSIONS**

15 440 The chemometric classification methods are widely used for food authentication purposes.  
16 441 As input experimental data set, any unspecific chromatographic signal (formally named  
17 442 chromatographic fingerprint) could be used. To make the most of classification performance,  
18 443 different chemical fractions characteristic of the studied material, different chromatographic  
19 444 conditions and different classification methods could be tried. Later, the best classification  
20 445 scenario has to be select in order to be applied in a real framework. In this work, several  
21 446 classification performance quality features have been presented and discussed.

22 447 As application example, two classification methods are applied on two sterolic  
23 448 chromatographic fingerprints obtained from two different normal-phase HPLC systems in  
24 449 order to discern the geographical origin of edible palm oil. For each one of the four  
25 450 classification scenarios, the corresponding quality features have been calculated and used to  
26 451 select the best one. For HPLC-UV fingerprint data, the best classifier is SIMCA classification,  
27 452 whereas for the HPLC-CAD one both classification methods behave on a similar way.

28 453 Finally, it is remarkable that all above mentioned parameters are applied jointly to different  
29 454 chromatographic fingerprints for the first time.

30 455

## 31 456 **Acknowledgment**

32 457 The authors are very grateful to Rikilt Wageningen University and Research Centre to supply  
33 458 the palm oil samples.

34 459

35 460

## 36 461 **SUPPLEMENTARY INFORMATION**

37 462 Descriptions of the "MEDINA" function for preprocessing of chromatographic data. Details of  
38 463 various processing options, as well as illustrated examples of the effect of different  
39 464 processing steps on a set of chromatographic data are provided.

## 465 REFERENCES

1. Y. Basiron, in *Bailey's Industrial Oil and Fat Products. Edible Oil and Fat Products: Edible Oils*, ed. F. Shaidi, Wiley-Interscience, Hoboken, 6th edn., 2005, vol. 2, ch. 8, pp. 333–429.
2. L. Siew Wai, in *Vegetable Oils in Food Technology. Composition, Properties and Uses*, ed. Frank, D. G., Wiley-Blackwell, Chichester, 2nd edn., 2011, ch. 2, pp. 25–58.
3. O.M. Lai, in *Healthful Lipids*, eds. C. Akoh, O.M. Lai, AOCS Press, Urbana, 2005, pp 731–749.
4. G. van Duijn, *Lipid Technol.*, 2013, **25**, 15–18.
5. R.J. Orsato, S.R. Clegg, H. Falcão, *J. Change Manag.*, 2013, **13**, 444–459.
6. S. Kelly, K. Heaton, J. Hoogewerff, *Trends Food Sci. Tech.*, 2005, **16**, 555–567.
7. A. Tres, G. van der Veer, M. Alewijn, E. Kok, S. van Ruth, in *Oil Palm: Cultivation, Production and Dietary Components*, ed. S.A. Penna, Nova Science, 2011, ch. 1, pp.1–44.
8. M. R. Alves, S. C. Cunha, J. S. Amaral, J.A. Pereira, M. B. Oliveira, *Anal. Chim. Acta.*, 2005, **549**, 166–178.
9. M.J. Lerma García, G. Ramis Ramos, J.M. Herrero Martínez, E. Simo Alfonso, *Rapid Commun. Mass Sp.* 2008, **22**, 973–978.
10. M.J. Lerma García, E.F. Simó Alfonso, A. Méndez, J.L. Lliberia, J.M. Herrero Martínez, *Food Res. Int.* 2011, **44**, 103–108.
11. D. Gázquez Evangelista, E. Pérez Castaño, M. Sánchez Viñas, M.G. Bagur González, *Food Anal. Method.* 2014, **7**, 912–925.
12. W. Kamm, F. Dionisi, C. Hischenhuber, K.H. Engel, *Food Rev. Int.*, 2001, **17**, 249–290.
13. S. Azadmard-Damirchi, *Food Addit. Contam. A.*, 2010, **27**, 1–10.
14. E. Wasowicz, in *Chemical, Biological, and Functional Aspects of Food Lipids*, eds. Z.E. Sikorski, A. Kolakowska, CRC Press, Boca Raton, 2nd edn, 2011, ch.7, pp. 113–134.
15. H. Saussem, in *Handbook of Analysis of Active Compounds in Functional Foods*, eds. L. Nollet, F. Toldrá, CRC Press, Boca Raton, 2012, ch. 35, pp. 787–804.
16. L. Nyström, in *Analysis of Antioxidant-Rich Phytochemicals*, eds., Z.R. Xu, L. Howard, John Wiley & Sons, Oxford, 2012, pp. 313–351.
17. S.L. Abidi, *J. Chromatogr. A*, 2001, **935**, 173–201.
18. A. Maija Lampi, V. Piironen, J. Toivo, in *Phytosterols as Functional Food Components and Nutraceuticals*, ed. P.C. Dutta, Marcel Dekker, New York, 2004, chs. 2,3, pp. 33–73.
19. S. Azadmard-Damirchi, P.C. Dutta, in *Olives and Olive Oil in Health and Disease Prevention*, eds. V. Preedy, R. Watson, Elsevier, London, 2010, ch. 27, pp. 249–257.
20. M. Biedermann, K. Grob, C. Mariani, *Fat Sci Tech.* 1993, **95**, 127–133.
21. G.Tadeusz, L. Frederic, S. Roman, S. Pat, *Anal. Chem.* 2006, **78**, 3186–3192.
22. I. Acworth, B. Bailey, M. Plante, P. Gamache, Simple and direct analysis of phytosterols in red palm oil by reversed-phase HPLC and charged aerosol detection. Thermo Fisher Scientific, 2011.
23. J.M. Bosque Sendra, L. Cuadros Rodriguez, C. Ruiz Samblas, A.P. de la Mata, *Anal. Chim. Acta.*, 2012, **724**, 1–11.
24. D.I. Ellis, V.L. Brewster, W.B. Dunn, J.W. Allwood, A.P. Golovanov, R. Goodacre, *Chem. Soc. Rev.*, 2012, **41**, 5706–5727.
25. L.A. Berrueta, R.M. Alonso Salces, K. Héberger, *Trends Anal. Chem.*, 2007, **35**, 74–86.
26. R. Leardi in *Modern Techniques for Food Authentication*, ed., D-W. Sun, Academic Press / Elsevier, Burlington, 2008, ch. 16, pp. 585–616.
27. P. Oliveri, G. Downey, *J. Chromatogr. A*, 2011, **1158**, 196–214.

- 1  
2  
3  
4  
5 28. P. de la Mata Espinosa, J.M. Bosque Sendra, R. Bro, L. Cuadros Rodríguez, *Anal. Bioanal. Chem.* 2011, **399**, 2083–2092.
- 6  
7 29. P. de la Mata Espinosa, J.M. Bosque Sendra, R. Bro, L. Cuadros Rodríguez, *Talanta*, 2011, **85**, 177–182.
- 8  
9 30. C. Ruiz Samblás, L. Cuadros Rodríguez, A. González Casado, F. de Paula Rodríguez García, de P. la Mata Espinosa, J.M. Bosque Sendra, *Anal. Bioanal. Chem.*, 2011, **399**, 2093–2103.
- 10  
11 31. C. Ruiz Samblás, F. Marini, L. Cuadros Rodríguez, A. González Casado, *J. Chromatogr. B.*, 2012, **910**, 71–77.
- 12  
13 32. C. Ruiz Samblás, J.M. Cadenas, D.A.Pelta, L.Cuadros Rodríguez, *Anal. Bioanal. Chem.*, 2014, **406**, 2591–2601.
- 14  
15 33. T. Galeano Díaz, I. Duran Meras, J. Sánchez Casas, M.F. Alexandre Franco, *Food Control*, 2005, **16**, 339–347.
- 16  
17 34. A. Tres, C Ruiz Samblás, G. Van Der Veer, S.M. Van Ruth, *Food Chem.*, 2013, **137**, 142–150.
- 18  
19 35. C. Ruiz Samblás, C. Arrebola Pascual, A. Tres, S. Van Ruth, L. Cuadros Rodríguez, *Talanta*. 2013, **116**, 788–793.
- 20  
21 36. W. Kamm, F. Dionisi, C. Hischenhuber, H-G. Schmarr, K-H. Engel, *Eur. J. Lipid Sci., Tech.* 2002, **104**, 756–761.
- 22  
23 37. G. Tomasi, F. Savorani, S.B. Engelsen, *J. Chromatogr. A.*, 2011, **1218**, 7832–7840.
- 24  
25 38. K.H. Esbensen, P. Geladi, in *Comprehensive Chemometrics: Chemical and Biochemical Data Analysis. Vol. 2: Data Preprocessing. Linear Soft-Modeling. Unsupervised Data Mining*, eds. S.D. Brown, R. Tauler, B. Walczak, Elsevier, Amsterdam, 2009, ch. 13, pp. 211–226.
- 26  
27 39. C.B.Y. Cordella, in *Analytical Chemistry*, ed. I.S. Krull, InTech, Rijeka, 2012, ch. 1, pp. 1–46.
- 28  
29 40. R. Bro, A.K.Smilde, *Anal. Methods.*, 2014, **6**, 2812–2831.
- 30  
31 41. M. Forina, P. Oliveri, S. Lanteri, M. Casale, *Chemom. Intell. Lab. Syst.*, 2008, **93**, 1320–148.
- 32  
33 42. F. Marini, *Curr. Anal. Chem.*, 2010, **6**, 72–79.
- 34  
35 43. S. Wold, M. Sjostrom, L. Eriksson, *Chemom. Intell. Lab. Syst.*, 2001, **58**, 109–130.
- 36  
37 44. D. Ballabio, V. Consonni, *Anal. Methods.*, 2013, **5**, 3790–3798.
- 38  
39 45. M. Sokolova, N. Japkowicz, S. Szpakowicz, in *Evaluation Methods for Machine Learning: Papers from the 2006 AAAI Workshop (WS-06-06)*, American Association for Artificial Intelligence, 2006.
- 40  
41 46. N. Japkowicz, M. Shah, *Evaluating Learning Algorithms: A Classification Perspective*, Cambridge University Press, New York, 2011, ch. 3-4.
- 42  
43 47. M. Bekkar, H. Kheliouane Djemaa, T. Akrouf Alitouche, *J. Inform. Eng. Appl.*, 2013, **3**, 27-38.
- 44  
45 48. C. Ferri, J. Hernández-Orallo, R. Modroui, *Pattern Recogn. Lett.*, 2009, **30**, 27–38.
- 46  
47 49. M. Felking, in *Quality Measures in Data Mining*, ed. F. Guillet, H.J. Hamilton, Springer-Verlag, Berlin Heidelberg, 2007, ch.12, pp. 290–292.
- 48  
49 50. K.L. Gwet, *Handbook of Inter-Rater Reliability*, Advanced Analytics LLC, Gaithersburg, 3rd ed., 2012, ch. 2, pp. 15–46.
- 50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

**Table 1.** Geographical origin of the 102 palm oil samples.

Asia		Africa		America	
Country	Samples	Country	Samples	Country	Samples
India	5	Cameroon	2	Brazil	16
Indonesia	24	Ghana	20		
Malaysia	19	Guinea	3		
Papua N. Guinea	7	West Africa	5		
Salomon	1				

**Table 2.** Continent distribution of the two set of samples.

Set	Continent	N° samples
Training set 72 samples (70.6 %)	Africa	20
	America	10
	Asia	42
Validation set 30 samples (29.4 %)	Africa	10
	America	6
	Asia	14

**Table 3.** Multiclass contingency table for ternary classification.

	Actual CLASS 1	Actual CLASS 2	Actual CLASS 3	TOTAL
Assigned CLASS 1	$a_1$	$e_{2,1}$	$e_{3,1}$	$AC_1$
Assigned CLASS 2	$e_{1,2}$	$a_2$	$e_{3,2}$	$AC_2$
Assigned CLASS 3	$e_{1,3}$	$e_{2,3}$	$a_3$	$AC_3$
Assigned NO-123	$e_{1,NO}$	$e_{2,NO}$	$e_{3,NO}$	$AC_4$
<b>TOTAL</b>	$TC_1$	$TC_2$	$TC_3$	<b>T</b>

$a_i$  = the number of assignment agreement of the class "i";  $e_i$  = the number of assignment error;  $TC_i$  = the total of actual samples from the class "i";  $AC_i$  = the total of assigned samples to the class "i"; T = the total number of samples. "NO-123" represents a fictitious class where the samples that do not assign to any classes, are allocated.

**Table 4.** Standard contingency table for binary classification.

	LABEL (L) Actual POSITIVE	no-LABEL (nL) Actual NEGATIVE	TOTAL
Assigned POSITIVE	<b>TP</b>	<b>FP</b>	<b>AP = TP+FP</b>
Assigned NEGATIVE	<b>FN</b>	<b>TN</b>	<b>AN = FN+TN</b>
<b>TOTAL</b>	<b>(TL) = TP+FN</b>	<b>(TnL) = FP+TN</b>	<b>T</b>

TP = true positive, the number of positive samples that are correctly identified as positive; FN = false negative, the number of positive samples that are misclassified as negative samples; FP = false positive, the number of negative samples that are incorrectly identified as positive samples; TN = true negative, the number of negative samples that are correctly identified as negative samples; AP and AN = the total of assigned positive and negative samples, respectively; TL and TnL = the number of labelled (actual) samples as positive and negative, respectively; T = the total number of samples.

Calculation example for obtaining the contingency table of the binary classification (class 1 / class n1) from the contingency table shown in Table 3:

$$TP = a_1; FP = e_{2,1} + e_{3,1}; FN = e_{1,2} + e_{1,3} + e_{1,NO}; TN = a_2 + a_3 + e_{2,3} + e_{2,NO} + e_{3,2} + e_{3,NO}$$

**Table 5.** Contingency tables obtained from both HPLC-UV and HPLC-CAD data for the two classification methods (SIMCA and PLS-DA) when the geographical origin is classified by considering three continents: Africa, America and Asia (a three-class classification). In this table, the sample numbers assigned to each class for the validation set (rows) are shown. Between parentheses, the corresponding rates, in %, in relation to the sample total number of each class (columns).

Fingerprint	Assigned class	SIMCA			PLS-DA		
		Actual class			Actual class		
		Africa	America	Asia	Africa	America	Asia
HPLC-UV	Africa	8 (80.0)	2 (33.3)	3 (21.4)	6 (60.0)	0 (0)	2 (14.3)
	America	0 (0)	2 (33.3)	0 (0)	2 (20.0)	2 (33.3)	0 (0)
	Asia	0 (0)	0 (0)	11 (78.6)	2 (20.0)	3 (50.0)	11 (78.6)
	No assigned	2 (20.0)	2 (33.3)	0 (0)	0 (0)	1 (16.7)	1 (7.1)
HPLC-CAD	Africa	10 (100)	2 (33.3)	3 (21.4)	7 (70.0)	1 (16.7)	2 (14.3)
	America	0 (0)	0 (0)	0 (0)	1 (10.0)	3 (50.0)	1 (7.1)
	Asia	0 (0)	2 (33.3)	8 (57.1)	2 (20.0)	1 (16.7)	11 (78.6)
	No assigned	0 (0)	2 (33.3)	3 (21.4)	0 (0)	1 (16.7)	0 (0)

**Table 6.** Values of quality performance features from two classification methods (SIMCA and PLS-DA) by two fingerprint data (HPLC-UV and HPLC-CAD) for the geographical origin between three continents: Africa, America and Asia (a three-class classification). (i): Pooled performance features of the four 3-class classifiers. (ii): overall values of quality performance features.

	HPLC-UV		HPLC-CAD	
	SIMCA	PLS-DA	SIMCA	PLS-DA
<b>(i) Pooled performance features</b>				
Sensibility (or Recall)	0.70	0.63	0.67	0.70
Specificity	0.92	0.80	0.92	0.85
Positive predictive value (or Precision)	0.87	0.67	0.89	0.72
Negative predictive value	0.86	0.81	0.84	0.84
Youden index	0.62	0.44	0.58	0.55
Positive likelihood rate	–	3.97	–	4.71
Negative likelihood rate	0.32	0.44	0.33	0.35
F-measure	0.76	0.69	0.71	0.76
Discriminant power	–	0.52	–	0.63
Efficiency (or Accuracy)	0.85	0.77	0.82	0.81
AUC (or Correctly classified rate)	0.81	0.72	0.79	0.77
Matthews correlation coefficient	0.66	0.46	0.64	0.55
Agreement probability	0.85	0.77	0.82	0.81
Chance agreement probability	0.56	0.57	0.56	0.56
Kappa coefficient	0.63	0.45	0.59	0.55
<b>(ii) Overall performance features</b>				
Overall agreement probability	0.70	0.63	0.60	0.70
Overall chance agreement probability	0.33	0.36	0.32	0.36
Overall KAPPA coefficient	0.55	0.42	0.41	0.53

**FIGURE CAPTIONS**

1  
2  
3  
4  
5  
6  
7  
8  
9  
10 **Figure 1.** Chromatograms of the same sample of palm oil from Africa showing the three  
11 characteristic regions, obtained from the data of the sterolic fraction by: (a)  
12 HPLC-UV, and (b) HPLC-CAD. IS denotes the internal standard. See text for  
13 further descriptions.  
14  
15

16  
17  
18 **Figure 2.** PC1/PC2 scores and PC1 loadings plots obtained from the data of the sterolic  
19 chromatographic data from the palm oil samples of three different continents:  
20 America (green squares); Africa (red rhombus); and Asia (blue triangles), by: (a)  
21 and (b) HPLC-UV; (c) and (d) HPLC-CAD.  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

Figure 1

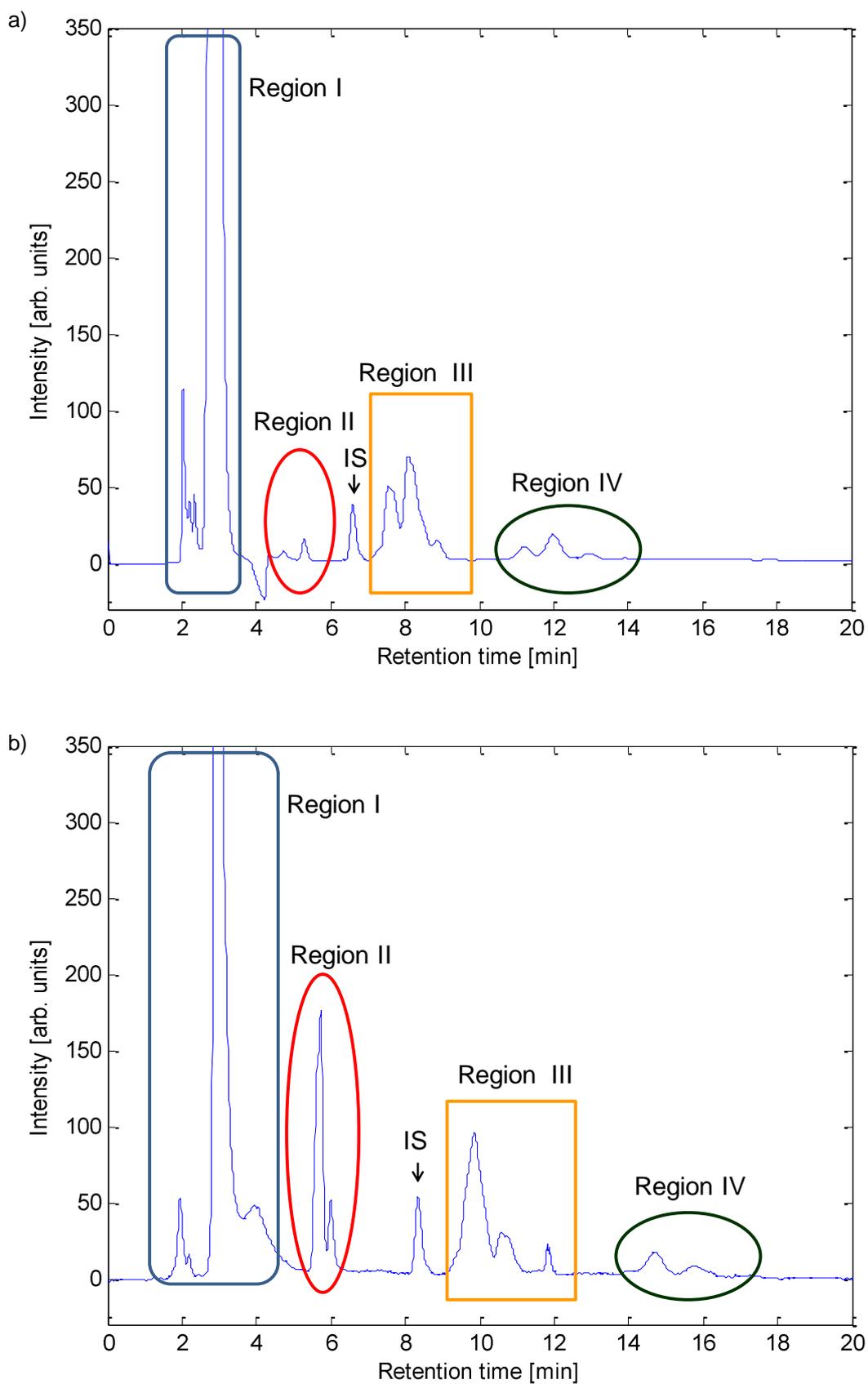


Figure 2

