

ISTITUTO NAZIONALE DI RICERCA METROLOGICA  
Repository Istituzionale

Authentication of cocoa bean shells by near- and mid-infrared spectroscopy and inductively coupled plasma-optical emission spectroscopy

This is the author's submitted version of the contribution published as:

*Original*

Authentication of cocoa bean shells by near- and mid-infrared spectroscopy and inductively coupled plasma-optical emission spectroscopy / Mandrile, Luisa; Barbosa-Pereira, Letricia; Sorensen, Klavs Martin; Giovannozzi, Andrea Mario; Zeppa, Giuseppe; Engelsen, Søren Balling; Rossi, Andrea Mario. - In: FOOD CHEMISTRY. - ISSN 0308-8146. - 292:(2019), pp. 47-57. [10.1016/j.foodchem.2019.04.008]

*Availability:*

This version is available at: 11696/61728 since: 2021-03-09T19:06:42Z

*Publisher:*

Elsevier

*Published*

DOI:10.1016/j.foodchem.2019.04.008

*Terms of use:*

This article is made available under terms and conditions as specified in the corresponding bibliographic description in the repository

*Publisher copyright*

(Article begins on next page)

Manuscript Number: FOODCHEM-D-18-06255R1

Title: Authentication of cocoa bean shells by near-infrared and mid-infrared spectroscopy and inductive coupled plasma-optical emission spectroscopy

Article Type: Research Article (max 7,500 words)

Keywords: cocoa bean shells, food traceability, data fusion, near infrared spectroscopy, mid infrared spectroscopy, inductive coupled plasma.

Corresponding Author: Mrs. Luisa Mandrile,

Corresponding Author's Institution: INRIM

First Author: Luisa Mandrile

Order of Authors: Luisa Mandrile; Letricia Barbosa-Pereira, PhD; Klavs M Sorensen, PhD; Andrea M Giovannozzi, PhD; Giuseppe Zeppa, Prof.; Søren B Engelsen, Prof.; Andrea M Rossi, PhD

Abstract: The aim of this study was to evaluate the efficacy of a multi-analytical approach for origin authentication of cocoa beans shells (CBS). The overall chemical profiles of cocoa bean shells from different origins were collected and measured using diffuse reflectance near-infrared spectroscopy (NIRS) and attenuated total reflectance mid-infrared spectroscopy (ATR-FT-IR) for molecular composition, as well as inductive coupled plasma-optical emission spectroscopy (ICP-OES) for elemental composition. Exploratory chemometric techniques were employed to identify systematic patterns related to the geographical origin of samples based on each technique using Principal Components Analysis (PCA). A combination of the three techniques proved to be the most promising approach to establish classification models. Partial Least Squares-Discriminant Analysis model of the fused PCA scores of three independent models was used and compared with single technique models. CBS samples were better classified by the fused model. Satisfactory classification rates were obtained for Central Africa samples with accuracy of 0.84.

# Highlights

---

- Multi-analytical approach for origin authentication of cocoa beans shells is proposed.
- Principal Component Analysis of NIR, ATR-FT-IR and ICP-OES data was discussed.
- Samples from Ecuador and central Africa were precisely classified by PLS-DA.
- Samples from São Tomé showed more features in common with the American samples than with the African samples.
- CBS samples were better classified by the fused model than by the three single analytical techniques.

1 **Authentication of cocoa bean shells by near-infrared**  
2 **and mid-infrared spectroscopy and inductive**  
3 **coupled plasma-optical emission spectroscopy**

4 Luisa Mandrile<sup>a</sup>, Letricia Barbosa-Pereira<sup>b</sup>, Klavs Martin Sorensen<sup>c</sup>, Andrea Mario Giovannozzi<sup>a</sup>,  
5 Giuseppe Zeppa<sup>b</sup>, Søren Balling Engelsen<sup>c</sup> and Andrea Mario Rossi<sup>a</sup>

6 <sup>a</sup> *Quality of life Division, Food Metrology program, Istituto Nazionale di Ricerca Metrologica, Strada delle Cacce, 91 10135,*  
7 *Torino, Italy*

8 <sup>b</sup> *Department of Agricultural, Forestry, and Food Sciences (DISAFA), University of Turin, Largo Paolo Braccini 2, 10095*  
9 *Grugliasco (TO), Italy*

10 <sup>c</sup> *Department of Food Science, University of Copenhagen, Rolighedsvej 26 DK-1958 Frederiksberg, Denmark*

11

12 <sup>\*</sup> *Corresponding author Luisa Mandrile, tel +39 011 3919329; e-mail [L.mandrile@inrim.it](mailto:L.mandrile@inrim.it)*

13

14 **Keywords:** cocoa bean shell, food traceability, data fusion, near infrared spectroscopy, mid infrared  
15 spectroscopy, inductive coupled plasma.

16 **Abstract**

17 The aim of this study was to evaluate the efficacy of a multi-analytical approach for origin authentication of  
18 cocoa beans shells (CBS). The overall chemical profiles of cocoa bean shells from different origins were  
19 collected and measured using diffuse reflectance near-infrared spectroscopy (NIRS) and attenuated total  
20 reflectance mid-infrared spectroscopy (ATR-FT-IR) for molecular composition, as well as inductive coupled  
21 plasma-optical emission spectroscopy (ICP-OES) for elemental composition. Exploratory chemometric

22 techniques were employed to identify systematic patterns related to the geographical origin of samples based on  
23 each technique using Principal Components Analysis (PCA). A combination of the three techniques proved to be  
24 the most promising approach to establish classification models. Partial Least Squares-Discriminant Analysis  
25 model of the fused PCA scores of three independent models was used and compared with single technique  
26 models. CBS samples were better classified by the fused model. Satisfactory classification rates were obtained  
27 for Central Africa samples with accuracy of 0.84.

## 28 **1. Introduction**

29 Since the 19<sup>th</sup> century cocoa has seen a continuous growth of consumption in a variety of forms, leading to an  
30 outstanding economic interest of chocolate industries for constant innovation and modernization. As many other  
31 agro-food activities, cocoa industry produces large amounts of by-products (<https://www.icco.org/>). Cocoa bean  
32 shells (CBS) is one of the main by-products, which represents the 12 % of weight after husking and grinding of  
33 dried cocoa seeds. CBS represents a non-negligible disposal problem and thus legislation and environmental  
34 issues are forcing industries to define process optimization and recovery/recycling strategies. Recently,  
35 bioconversion of by-products has raised the interest of scientific research and in several countries strategic vision  
36 or dedicated policies are being prepared to manage food industry wastes in the most efficient way – abandoning  
37 the “take, make and dispose” behavior and instead acting out a circular economy paradigm (Sørensen, Aru,  
38 Khakimov, Aunskjær, Engelsen, 2018). The increasing interest for byproducts has certainly an environmental  
39 basis, but an important role is played by the tendency to reduce the use of synthetic additives and replace them  
40 with natural substances in food. Research concerning new natural additives with high quality/costs ratio is  
41 increasing nowadays (Carocho, Morales, Ferreira, 2015). Moreover, the demand of new functional foods, rich in  
42 bio compounds such as polyphenols, fiber, n-3 fatty acids etc., drives interest for rich food wastes, such as seeds  
43 husks (Andrade, Gonçalves, Maraschin, Ribeiro-do-Valle, Martínez, Ferreira, 2012; Jansman, Versteegen,  
44 Huisman, Van den Berg, 1995). Vegetal by-products are rich of nutrients, such as fiber, polyphenols, minerals  
45 and their recycling represent one of the valorization strategies. The development of CBS valorization strategies  
46 is aimed at reducing the environmental impact of the cocoa production and provides information to promote  
47 conversion of a by-product into added-value products with application in food and healthcare sectors. The  
48 definition of the chemical composition of CBS from different countries is meant to evaluate the systematic

49 differences due to their origin. Chemical analysis of CBS has been carried out in several research papers because  
50 of its interesting features related to flavor, phenolic compounds and nutritional values (Barbosa-Pereira,  
51 Guglielmetti, Zeppa, , 2018; Manzano, Hernández, Quijano-Avilés, Barragán, Chóez-Guaranda, Viteri, Valle,  
52 2017; Redgwell, Trovato, Merinat, Curti, Hediger, Manez, 2003; Serra Bonvehí, and Escolá Jordà, 1998;  
53 Martín- Cabrejas, Valiente, Esteban, Mollá, Waldron, 1994;), however a complete characterization, using  
54 different methodologies to highlight similarities and differences in composition of samples from different  
55 countries has not been accomplished yet. In this work, CBS samples from different countries were analyzed with  
56 three different analytical methods. Near infrared spectroscopy (NIRS), mid infrared spectroscopy by attenuated  
57 total reflectance (ATR-FT-IR) and inductively coupled plasma-optical emission spectroscopy (ICP-OES) were  
58 used to collect a wide chemical information, both molecular and elementary. The aim of this study was to  
59 evaluate the validity of simple and rapid analytical techniques, supported by a chemometric approach, for the  
60 identification of differences due to different geographical origin of samples of CBS, with the perspective of a  
61 future application for traceability and origin authentication of CBS as food additive.

62 Nowadays, the exchange of food is realized in a complex and interconnected global net, and food products are  
63 often exposed to frauds, false information, contamination risk and counterfeiting. For this reason, it is extremely  
64 important to protect and valorize authentic products, including regionals specialties. Innovative, reliable  
65 strategies to individuate specific markers of origin, as well as characteristic compositional patterns that can be  
66 associated to a precise origin are needed (Mandrile, Giovannozzi, Zeppa, Rossi, 2016). Geographical origin  
67 indicators should provide an analytical response to the geographical traceability problem and support the  
68 documental certification, which is used today to guarantee food and food-additives provenience. Different  
69 techniques such as NMR and isotope ratio mass spectrometry can play a relevant role to provide origin indicators  
70 (Lee, et al., 2011). Rapid and non-destructive techniques, such as near infrared spectroscopy, are particularly  
71 interesting because of the possibility to obtain an efficient and non-biased overview of the sample chemistry  
72 (Sørensen, Khakimov, Engelsen, 2016). The chemical specificity and ease of sampling of NIR spectroscopy  
73 make it an attractive tool for rapid and comprehensive food analysis. The complex pattern of signals revealed by  
74 IR analysis, both in the near and mid infrared spectral region, is correlated to the content of the different  
75 chemical constituents, such as proteins, fatty acids, carbohydrates, alimentary fibers and phenolic compounds.  
76 Statistics and multivariate data analysis offer powerful tools to identify robust correlations between measured

77 data and geographical origin, and validated models can provide useful methods for the recognition of unknown  
78 samples, with a certain probability (Peres, Barlet, Loiseau, Montet, 2007; Kelly, Heaton, Hoogewerff, 2005). In  
79 this work, chemometrics was used for data analysis to calculate at first explorative, and subsequently predictive,  
80 models. Principal Component Analysis for data exploration and visualization is a well-established strategy to  
81 allow the extraction of useful information from numerous experimental results in food science (Munck,  
82 Nørgaard, Engelsen, Bro, Andersson, 1998). Moreover, data fusion for multi-block analysis was used to improve  
83 models, gaining information from several different analytical techniques (Biancolillo, Bucci, Magrì, Magrì,  
84 Marini, 2014; Skov, Honoré, Hansen, Næs, Engelsen, 2014; Silvestri et. al, 2014; Zakaria, et al, 2010).

85

## 86 **2. Material and Methods**

### 87 *2.1 Samples*

88 Fermented and dried cocoa (*Theobroma cacao* L.) samples were selected and collected within COVALFOOD  
89 project funded by European Union's Seventh Framework, involving five Italian chocolate industries. A complete  
90 list of 78 samples with the associated information about supplier, provenience and variety is reported in table  
91 1S.1 in supplementary information. For an easier exploration of the sample pool, charts of geographical and  
92 | varietal distribution are shown in figure 1S.1. All samples were imported as untreated raw materials, and the  
93 | geographical origin was guaranteed by the supplying industry. All samples were roasted and decorticated in  
94 laboratory in a ventilated oven for 20 min at 130°C. After roasting, the fragile shell of the beans was separated  
95 by mechanical rubbing and removed by hoover suction. The collected cocoa bean shells (CBS) were ground  
96 using an ultra-centrifugal mill Retsch ZM 200 (RetschGmbH, Haan, Germany) and stored as dry fine powders  
97 (250 µm) in a desiccator in closed containers.

### 98 *2.2 Near infrared spectroscopy*

99 NIR spectra of CBS were collected in the spectral range 10000 - 4000 cm<sup>-1</sup> (1000 - 2500 nm) using an Antaris II  
100 FT-NIR spectrometer (Thermo Fisher, Waltham, USA) in diffuse reflectance mode. The integrating sphere  
101 accessorize was used to collect diffuse reflected light. CBS was analyzed without sample pretreatment; 0.1 g of  
102 powder in a quartz glass vial located over the integrating sphere. 32 scans were collected per each sample with

103 spectral resolution of  $8\text{ cm}^{-1}$ . A clean flat golden surface was used for background collection. Three  
104 measurement replicates were collected per sample. All samples were measured in randomized order.

### 105 *2.3 Mid infrared spectroscopy*

106 ATR-FT-IR spectra in the mid infrared region between  $500 - 4000\text{ cm}^{-1}$  were collected using Nicolet FT-IR  
107 spectrometer (Thermo Fisher, Waltham, USA), Germanium crystal ( $n = 5.7$ ) for total reflection was used which  
108 allows a maximum sample penetration of  $1\text{ }\mu\text{m}$ . 64 scans were needed for a good signal to noise with  $4\text{ cm}^{-1}$   
109 resolution. The sample powder was pressed with a conical tip on the crystal, the pressure applied was 15 Bar.  
110 The tip and the crystal were washed with ethanol between one sample analysis and the following. Three spectra  
111 were collected for each sample, resampling at each replicate.

### 112 *2.4 ICP-OES elemental composition*

113 ICP-OES measurements were performed on an Agilent 5100 Synchronous Vertical Dual View (Agilent, Santa  
114 Clara, California, USA), equipped with an EasyFit torch (Agilent P/N G8010-60228). Samples were measured in  
115 radial mode, using a plasma flow of  $12\text{ ml/min}$  and nebulizer flow of  $0.7\text{ ml/min}$ , with a rinse time of 15 seconds  
116 and stabilization time of 15 seconds, in three replicates. Viewing height was set to  $8\text{ mm}$ , and pump speed to 12.  
117 Prior to measurement, the samples were digested in an Antor Paar Multiwave GO microwave oven: 5 mg of  
118 CBS samples were placed in the oven teflon tubes,  $1\text{ ml}$  of  $\text{HNO}_3$  5 % v/v was added, and the tubes were sealed  
119 to manufacturer specifications. The temperature ramp was set to reach  $180^\circ$  in 5 min, then held constant, and the  
120 total treatment lasted 40 min. After digestion the samples were further diluted with  $4\text{ ml}$   $\text{HNO}_3$  5 % v/v to obtain  
121 a clear solution, before being put in tubes and placed in the auto-sampler for the ICP analysis. All glassware,  
122 tubes and equipment were cleansed in  $\text{HNO}_3$  5 % v/v as needed.

### 123 *2.5 Data treatment*

124 Chemometric data analysis was carried out using PLS Toolbox from Eigenvector Research, Inc. (Manson, WA)  
125 for Matlab R2015a (Mathworks, Natick, USA). Principal Components Analysis (PCA) method is a linear  
126 factorization method uniquely suited for data exploration. As an explorative tool, PCA provides visualization of  
127 multivariate data as score points in a model space (Wold, Esbensen, Geladi 1987). PCA scores plot are useful to

128 explore data and to find correlation between measured variables and the information of interest, such as  
129 geographical provenience of CBS, in this case. Then PLS-DA (Barker and Rayens 2003) models were calculated  
130 to compare the classification performances of the three techniques separately with the results obtained by joining  
131 the three datasets and considering all information contemporarily. Ten classes were considered: Central Africa,  
132 Ecuador, Gulf of Mexico, Indonesia, Mexico, Peru, São Tomé, Colombia, Venezuela and Brazil. All the  
133 calculated PLS-DA models were validated using leave-one group-out cross validation. The subsets of samples  
134 used as tests sets in cross validation corresponds to the country of origin. For each technique data preprocessing  
135 details are reported. Leave-one group-out cross validation was performed, using as group vector the country of  
136 origin. Sensitivity (True Positive/(True Positive+False Negative)), Specificity (True Negative/(True  
137 Negative+False Positive)), Accuracy (correctly classified samples/total samples) and Precision (True  
138 Positive/(True Positive +False Positive) were considered as model evaluation parameters for each class in cross  
139 validation to compare classification performances of different techniques.

#### 140 *2.5.1 NIRS data treatment*

141 Preprocessing of NIRS data was applied to extract useful information from the dataset. Absolute absorbance  
142 variations and unwanted light scattering were removed using preprocessing of the NIRS data (Martens et al,  
143 2003). The most effective preprocessing was chosen based on the minimum differences between replicates on  
144 the PCA scores plots relative to the distance between samples. 2<sup>nd</sup> derivative (Savitzky Golay, filter width 15  
145 and polynomial order 2) coupled with standard normal variate (SNV); normalization was useful to remove  
146 random shift of the baseline offset (Barnes, Dhanoa, Lister, 1989). In addition, the derivatives of spectra were  
147 calculated to increase sensitivity to data trends changings. Processed spectra were shown in figure 2S.1.  
148 Unwanted variability was successfully removed as demonstrated by the narrow grouping of the replicates  
149 obtained after processing shown in figure 2S.2 in supplementary information. PCA was applied to visualize data  
150 and to investigate systematic differences among samples, and variables with peculiar relevance were identified.  
151 4LVs PLS-DA classification model was also calculated to discriminate classes of samples from different  
152 geographical areas. Same spectra preprocessing was used.

#### 153 *2.5.2 MIRS data treatment*

154 Preprocessing of data was performed to suppress useless variability associated to unwanted noise. The selection  
155 criterion for data preprocessing was the maximized closeness of the scores of technical replicates on PC1, as  
156 shown in figure 3S.1 in supplementary information. Baseline correction (using asymmetric weighted least  
157 squares algorithm, with basis filter of order 2) (Peng, Peng, Jiang, Wei, Li, Tan, 2010) followed by second  
158 derivative (Savitzky Golay, filter width 15 and polynomial order 2) and mean centering was selected as optimal  
159 preprocessing. PCA model for data visualization and exploration was calculated; PLS-DA classification model  
160 using 4 LVs of the same preprocessed data was also calculated to compare MIRS classification capabilities with  
161 the other techniques.

### 162 2.5.3 ICP-OES data treatment

163 ICP emission spectra were evaluated for quantification using a calibration curve per element. The calibration  
164 curves were estimated using two series of standards prepared by dilution of a certified standard mix (ICP Multi-  
165 element standard solution IV, Sigma Aldrich, Germany) containing known concentration of 21 elements (Al, B,  
166 Ba, Bi, Ca, Cd, Co, Cr, Cu, Fe, K, Li, Mg, Mn, Mo, Na, Ni, Pb, Sr, Tl, Zn). Standard concentrations were 0, 0.2,  
167 0.4, 0.6, 0.8, 1, 2, 4, 6, 8, 10, 20, 30, 40, 60, 80 100 mg/100g of the certified standard concentration, which was 5  
168 mg/l for all elements, out of Potassium that was 50 mg/l in the standard solution. Three emission wavelengths  
169 were monitored per each element, then the intensity revealed for only one  $\lambda$  was selected per each element based  
170 on the best correlation coefficient of the corresponding calibration curve and trying to avoid interferences  
171 between different elements:  $\lambda_{Al}=237.3$  nm;  $\lambda_B=249.7$  nm;  $\lambda_{Ba}=455.4$  nm;  $\lambda_{Bi}=190.2$  nm;  $\lambda_{Ca}=396.8$  nm;  $\lambda_{Cd}=228.8$   
172 nm;  $\lambda_{Co}=230.8$  nm;  $\lambda_{Cr}=206.2$  nm;  $\lambda_{Cu}=324.8$  nm;  $\lambda_{Fe}=234.4$  nm;  $\lambda_K=766.5$  nm;  $\lambda_{Li}=670.8$  nm;  $\lambda_{Mg}=285.2$  nm;  
173  $\lambda_{Mn}=259.4$  nm;  $\lambda_{Mo}=203.8$  nm;  $\lambda_{Na}=589.0$  nm;  $\lambda_{Ni}=221.6$  nm;  $\lambda_{Pb}=217.0$  nm;  $\lambda_{Sr}=421.6$  nm;  $\lambda_{Tl}=351.9$  nm;  
174  $\lambda_{Zn}=202.5$  nm.

175 The table of results was then imported in Matlab (Mathworks, Natick, USA) and processed with the PLS  
176 Toolbox for PCA model calculation and PLS-DA classification. Autoscaling was performed on the data. Three  
177 LVs were considered for PLS-DA classification model. Cross validation was used to evaluate the classification  
178 capabilities of the model, leaving one country out at each validation step, as described for the other techniques.

### 179 2.5.4 Data fusion

180 The multi-block tool of PLS toolbox by Eigenvector was used to fuse the PCA scores from the three single PCA  
181 models of the different analytical techniques. A joined model exploiting mid-level data fusion was obtained (;  
182 Borràs, et al, 2014). To make the interpretation clearer, the measurement replicates were averaged, and one  
183 matrix line per each sample was maintained for the three different original datasets (NIRS, MIR-ATR and ICP).  
184 Each block was first decomposed by PCA, and the resulting scores were fused into a new dataset. The samples'  
185 scores for the most relevant PCs were considered to calculate a new fused model. Seven PCs were considered for  
186 MIRS and ICP, and six PCs were considered for NIRS. Thus, twenty initial variables were used to build the new  
187 joined PCA model. Default autoscale was applied before joining data. PLS-DA method was then performed with  
188 autoscaled data to obtain a classification model (Ballabio, Consonni, 2013). The class vector was represented by  
189 the area of origin. It was composed of 10 classes i.e. Central Africa, Colombia, Ecuador, Gulf of Mexico,  
190 Indonesia, Mexico, Peru, São Tomé, Venezuela, Brazil. Unfortunately the number of samples per each class was  
191 not balanced, due to sample availability. Five latent variables were considered for the PLS-DA model, based on  
192 the minimum average classification error in cross validation, using leave-one country-out cross validation  
193 strategy.

### 194 **3. Results and Discussion**

#### 195 *3.1 NIRS spectroscopy characterization of CBS samples*

196 The NIRS profiles show the typical broad bands of overtones and combination bands of vibrational modes  
197 associated to the main constituents of vegetal origin materials. The assignment of the most bands of the NIR  
198 spectrum are reported in table 2S.1 in the supplementary information (Jacobsen, et. al. 2011). The mean NIR  
199 spectra of all CBS samples is shown in figure 1 a, together with the standard deviation profiles. Similar spectral  
200 shape was obtained for all samples, the same bands are present in all spectra with slight differences in mutual  
201 intensities.

#### 202 **Figure 1**

203 Vibrational spectroscopy represents a rapid strategy to gather chemical information of a complex matrix,  
204 reducing costs, time and environmental impact of analysis. NIR spectra can be effectively correlated to the main

205 alimentary components as widely reported in literature (De Oliveira, Roque, de Maia, Stringheta, Teófilo, 2018;  
206 Dong, Sørensen, He, Engelsen, 2017; Mandrile, Fusaro, Amato, Marchis, Martra, Rossi, 2018).

207 The sensitivity of NIRS to the botanical variety was tested at first, since it has been previously demonstrated in  
208 literature that differences in the chemical composition of different varieties of *Theobroma Cacao L.* were present  
209 (Elwers, Zambrano, Rohsius, Lieberei, 2009). The outcome of the PCA on the NIR spectra is shown in figure 1  
210 b. In contrast with expectations, different botanical varieties did not cause evident systematic clustering of NIR  
211 spectra. The scores of NIR spectra of *Forastero* and *Trinitario* samples were overlapped in the scores plot  
212 (figure 1 b), no separation occurred neither in the PC2/PC1 plot, nor in the later PCs (plots not shown). This can  
213 be probably attributed to the complexity of the samples' set, that introduces a lot of confusing variability.  
214 However, Arriba samples, a specific variety cultivated in Ecuador only (green squares on the scores plot in  
215 figure 1b), was specifically, even though not selectively, characterized by negative scores on PC1 and positive  
216 scores on PC2 attesting the capability of NIR spectra to catch common chemical features of Arriba samples. The  
217 loadings profiles (figure 2S.3 a) and the variance captured (figure 2S.4) allows to define what spectral regions  
218 are involved in each relevant PC. PC1, is mainly characterized by fatty acids bands as  $5670\text{-}5780\text{ cm}^{-1}$  ( $1^{\text{st}}$  C-H  
219 str) and  $4325\text{ cm}^{-1}$  ( $1^{\text{st}}$  C-H str +  $1^{\text{st}}$  C-H def  $\text{CH}_2$ ),  $4250\text{ cm}^{-1}$  ( $1^{\text{st}}$  C-H str +  $1^{\text{st}}$  C-H def). In addition PC1 captures  
220 also some regions related to proteins such as  $5170\text{-}5190\text{ cm}^{-1}$  ( $2^{\text{nd}}$  C=O of CONH),  $5269\text{ cm}^{-1}$  ( $2^{\text{nd}}$  C=O of  
221 COOH),  $6320\text{ cm}^{-1}$  ( $1^{\text{st}}$  N-H str of CONH) and  $6535\text{ cm}^{-1}$  ( $1^{\text{st}}$  N-H str of  $\text{RNH}_2$ ) and  $6950\text{ cm}^{-1}$ . PC2, instead,  
222 shows three maxima at  $4400\text{ cm}^{-1}$  ( $1^{\text{st}}$  O-H str +  $1^{\text{st}}$  C-C str, associated to starch),  $4763\text{ cm}^{-1}$  ( $2^{\text{nd}}$  O-H def +  $2^{\text{nd}}$   
223 C-O str of starch) and  $5000\text{ cm}^{-1}$  ( $2^{\text{nd}}$  O-H def +  $1^{\text{st}}$  C-O def of starch), this means that PC2 mostly represents the  
224 starch content into the samples. PCA highlighted a major content of fatty acids and vegetal proteins in the  
225 examined Arriba samples as shown in figure 1 c, d, whereas lower intensity in the spectral regions associable to  
226 polysaccharides, such as starch, was measured (corresponding enlarged spectral region not shown for brevity  
227 reasons).

228 As far as correlations between the geographical origin and NIR spectra are concerned, the information provided  
229 by the scores plot seems confused at a first look, however some interesting considerations can be underlined.  
230 Common features of all samples coming from central Africa were noticed in the scores plot (figure 2 a) when  
231 considering PC2. On average, central Africa samples (red rhombus in figure 2 a) show positive scores on PC2,

232 related to polysaccharides and starch bands mainly (figures 2S.3, 2S.4 can be consulted for all attributions of  
233 spectral bands to the PCs). Moreover other common features were noticed in further PCs, such as negative scores  
234 on PC3 (figure 2S.6 b) (where the main contributions are  $5218\text{ cm}^{-1}$ , 1<sup>st</sup> O-H str of phenols,  $5878\text{ cm}^{-1}$  1<sup>st</sup> C-H str  
235  $\text{CH}_3$ ,  $6075\text{ cm}^{-1}$  1<sup>st</sup> C-H str of R-CH-CH,  $7062\text{ cm}^{-1}$ , 2<sup>nd</sup> C-H str + 1<sup>st</sup> C-H def of aromatic compounds) and  
236 positive again on PC4 (Figure 2S.6 c) which is related mainly to carbohydrates ( $4790\text{ cm}^{-1}$  1<sup>st</sup> O-H str + 1<sup>st</sup> O-H  
237 def ROH o sucrose and starch,  $6264\text{ cm}^{-1}$ , 1<sup>st</sup> O-H str intramolecular H-bond of starch or glucose). Although the  
238 separation of the examined groups is not sufficient for selective discrimination, it was confirmed that the  
239 geographical origin information is captured by NIRS. As shown in figure 2 a, African samples from São Tomé (a  
240 little island in Guinea Gulf, at latitude  $0^\circ$ ) show features in common with samples coming from America, which  
241 on average showed negative scores on PC2. The scores of São Tomé samples (light blue rhombus in figure 2 a)  
242 are mixed with Gulf of Mexico Samples, this can be attributed to similar environmental and climatic conditions  
243 of the little islands, that influences the chemical composition of Cocoa fruits, and therefore of CBS (see also  
244 figures 2S.6 a to appreciate similitudes of São Tomé with samples from the islands and coasts of Gulf of  
245 Mexico). Moreover, Ecuador samples seemed more similar to the African samples than to the American, indeed,  
246 in figure 2 a, orange circles corresponding to Ecuador samples are mixed with red rhombus corresponding to  
247 samples from Central Africa. In figure 2 b the average NIR spectra of the macro classes, Africa and America, are  
248 compared with the spectra of São Tomé and Ecuador, that show peculiar behavior in contrast with the general  
249 trend.

250 The Asian samples are separated from the others (blue triangles in figure 2 b), because of high values on PCs 4,  
251 5 and 6. PC4 is characterized by a peak around  $4530\text{ cm}^{-1}$ . This spectral region, represented in figure 2 d is  
252 assigned to ROH combination modes, so it can be hypothesized that sugars' content differs for Asian samples  
253 with respect to all the others. The most represented spectral region in PC5 (which is relevant for the clustering of  
254 Asian samples) is the side of the peak at  $6300\text{ cm}^{-1}$ . This region, represented in figure 2 e, highlights that the  
255 bands' shape is relevant, more than its intensity in this case. PC6 is also responsible for the following spectral  
256 regions:  $4466\text{ cm}^{-1}$  (beta-glucan),  $5114\text{ cm}^{-1}$  (2<sup>nd</sup> C=O of esters) and  $7147\text{ cm}^{-1}$  typical of R-OH (as already  
257 mentioned figures 2S.3, 2S.4 can be consulted for all attributions of spectral bands to the PCs).

258 **Figure 2**

259 The definition of rules to correlate the NIR spectra variability with the geographic area of origin based on the  
260 PCA scores plot of NIR spectra is not immediate. However, some common trends were noticed for samples from  
261 the same area, and NIR spectra demonstrated to contain useful information for geographical provenience  
262 analysis.

### 263 3.2 ATR-FT-IR spectra

264 Spectral profiles in the mid infrared region are shown in figure 3 a. As well as for NIRS, ATR-FT-IR  
265 spectroscopy is expected to deliver information about the chemical composition of CBS samples including most  
266 of biochemical species present in the matrix. Although absorption bands in the mid infrared region are more  
267 defined and narrower because primary vibration modes absorb in this spectral region, the visual interpretation of  
268 spectra is difficult, especially in the so-called fingerprint region, between  $1750\text{ cm}^{-1}$  and  $500\text{ cm}^{-1}$ . Main bands  
269 interpretation is reported in table 3S.1 in supplementary information. (Socrates, 2001; Rubio- Diaz,  
270 Rodriguez- Saona, 2010; Li-Chan, Chalmers, Griffiths, 2011). The region between  $2260\text{-}2440\text{ cm}^{-1}$ , where  $\text{CO}_2$   
271 band is present, was excluded.

### 272 **Figure 3**

273 MIRS spectra provided information in agreement with NIRS investigation. Signals are more defined and spectral  
274 specificity is increased compared to NIRS, and PCA scores plots investigation resulted an effective strategy to  
275 explore spectra similarities. Similarities and differences between samples are ruled by PC1, 2 and 3. The  
276 correspondence between PCs and MIR spectral regions was evaluated analyzing figure 3S.4, where the MIR  
277 spectrum was superimposed over the histogram of the percentage of variance captured by each PC, to understand  
278 what bands drive the scores distribution on the scores plot. PC1 is mainly dominated by  $\text{CH}_x$  vibrations in the  
279  $3000\text{-}2800\text{ cm}^{-1}$  and  $1460\text{-}1420\text{ cm}^{-1}$  region (samples with high intensity of signals at  $2920\text{ cm}^{-1}$  and  $1463\text{ cm}^{-1}$   
280 present lower values of PC1), moreover  $1730\text{ cm}^{-1}$  peak (C=O stretching) that showed increased intensity in  
281 Arriba samples is also represented in PC1; PC2 captures variance in  $1700\text{-}1650\text{ cm}^{-1}$  region (high values of PC2  
282 mean lower intensity at  $1560\text{ cm}^{-1}$  and  $1525\text{ cm}^{-1}$  of amide I-II and lower intensity of the  $1690\text{ cm}^{-1}$  shoulder).  
283 Several peaks associated to carbohydrates are also relevant, for example  $763\text{ cm}^{-1}$  related to pyranose compounds  
284 is modeled by PC5. Variety information reveals a certain grouping of Arriba sample that show high PC2 scores

285 and lower intensity of PC5 in Arriba samples, in agreement with NIRS results. The scores plot colored by variety  
286 information is shown in figure 3S.5.

287 The different geographical provenience drives a differentiation between samples and some general  
288 considerations can be extracted from the scores plot (figure 3 b,c). PC2 certainly explains interesting  
289 characteristics of Central Africa samples, that show positive scores on PC2. Samples from São Tomé showed  
290 more similarities with samples from Gulf of Mexico, Venezuela and Colombia, as attested also by NIRS data  
291 shown in the previous paragraph. This confirms that similar climatic and environmental conditions are crucial in  
292 determining the chemical composition captured by spectroscopic techniques, as previously reported in literature  
293 for cocoa samples (Marseglia, et al, 2017). African samples show higher intensity at  $2954\text{ cm}^{-1}$  and  $2870\text{ cm}^{-1}$  in  
294 the  $\text{CH}_x$  stretching vibrations (Figure 3 d). Moreover, PC5 and PC6 were relevant to identify features in common  
295 between Ecuadorian samples. 87% of Ecuador samples were placed to the left of the left diagonal of the  
296 PC6/PC5 plot (figure 3 c). This is due to the ratio between  $1280\text{ cm}^{-1}$  (Amide III of  $\beta$ -sheet proteins) and  $1320$   
297  $\text{cm}^{-1}$  or  $1440\text{ cm}^{-1}$  that allows to separate samples from Ecuador from other American samples, as shown in  
298 figure 3 e. Moreover low values in PC5 reflect low intensities at  $673\text{ cm}^{-1}$  and  $1600\text{ cm}^{-1}$  (ring breathing modes  
299 of polysaccharides) as already noticed for Arriba samples (enlarged spectral regions not shown for brevity  
300 reasons).

301 The ATR-FT-IR spectrum represents the sum of numerous bands of several functional groups, which are  
302 contemporarily present in more than one biochemical compound. Beyond the hypothesized interpretation, it  
303 should be stressed that an accurate understanding of what peaks and bands drive the scores distribution should  
304 be managed carefully to avoid misinterpretation. To univocally associate the relevant spectral regions to specific  
305 classes of compounds remains complicated when a whole complex matrix such as food is analyzed. However,  
306 the possibility to identify spectral features that precisely, characterize samples from the same origin is an  
307 indication that a correlation between geographical origin and vibrational spectra can be modeled.

### 308 *3.3. ICP-OES elemental characterization of CBS samples*

309 The raw ICP-OES results are shown in Table 4S.1 in supplementary information. The most abundant elements  
310 are by far Ca, Mg, K which have a concentration at least one order of magnitude higher compared to all other

311 elements. Among the secondary elements, particularly relevant were Al, Fe and Li (Barker and Rayens 2003).  
312 Relevant amounts of lead were revealed in all samples (around 0.3 mg/kg), which is a high value compared with  
313 the average content of lead in foods reported in 2007 by the Agency for Toxic Substances and Disease Registry  
314 (Abadin H., et al. 2007). All other elements were revealed in concentration lower than 0.2 mg/kg, particularly  
315 low concentrations were determined for Ni and Cr. PCA was used to identify major variance directions that can  
316 be related to geographical origin. Five samples were identified as very different from the others. They were SB3,  
317 SB4 from Brazil, ICAM10 from Congo, FER8 from Uganda and FER13 from Côte d'Yvoire. These samples  
318 were excluded as outliers because of their very low K content. Boron, Potassium, Magnesium and Calcium are  
319 responsible of the most variance captured by PC1, which resulted not to be particularly correlated to provenience  
320 of samples. Aluminum, Chromium, Iron, Sodium and Nickel are particularly relevant for PC2, whereas  
321 Cadmium, Cobalt and Molybdenum together with Calcium and Manganese are mostly represented in PC3, as  
322 shown in figure 4 d.

323 Examining the PC2/PC3 loadings and scores plot (figure 4 a, b), high levels of Fe and Al resulted to be  
324 characteristic for African continent for most of Central Africa Samples, moreover a general deficiency of Ca, K,  
325 Mg, Ni, was revealed. Interestingly some similitudes of São Tomé samples with American samples were  
326 captured by PC2. Precisely a relatively higher content of Fe, Al, Cu and Ni was revealed for this samples, this  
327 trend makes São Tomé samples more like American than to African samples. Moreover, São Tomé samples are  
328 characterized by high content of Ba with respect to others. Conversely, Ecuador samples did not show any  
329 specific elemental profile.

#### 330 **Figure 4**

#### 331 *3.4 Data fusion to merge chemical information provided by the different analytical techniques*

332 The idea of data fusion is to merge information, provided by different analytical determinations, in one single  
333 data set, to enhance the quality of the results. The obtained joined PCA model clearly shows that all the three  
334 datasets provide useful information for the final model. It was noticed that the three most represented variables  
335 in PC1 were one from MIR-ATR, one from ICP and one from NIRS (figure 5S.1 in supplementary information).  
336 The scores plot and the loadings projected on the PC2/PC1 space are shown in figure 5. The grouping of samples

337 based on the geographical origin was improved by the multi analytical model. Proximity, and hence common  
338 features, were appreciated for samples from the same geographical area.

339 Classification models were calculated to quantify the grouping performances of the joined model compared to  
340 the three single models, based on the geographical origin. Even though interesting observations were previously  
341 discussed for the three techniques separately, and some correlation between geographical origin and the  
342 composition was defined, single technique outputs were not accurate and precise for the recognition of the  
343 geographical origin of samples in predictive classification models. In table 1 the most classification figure of  
344 merit (sensitivity, specificity, error rate, accuracy, precision) relative to PLS-DA classification models for the  
345 geographical discrimination were reported. The classification performances for the samples' classes composed of  
346 more than 5 samples were shown. Classification results were higher for the joined model compared to each of  
347 the three single models for Central Africa, Ecuador and Gulf of Mexico classes. This experimental evidence was  
348 in agreement with literature findings corroborating mid-level or high-level data fusion to increase predictive  
349 performance of classification models (Doeswijk, Smilde, Hageman, Westerhuis, Van Eeuwijk, 2011). Single  
350 techniques provide null accuracy and precision for most classes, out of Central Africa. Moreover, merging  
351 information from the three techniques, the accuracy (correctly classified samples rate) increased.

#### 352 **Table 1**

353 NIRS, MIRS and ICP profiles together deliver sufficiently accurate information to capture the common features  
354 of African samples, and to distinguish them from all the others. Unfortunately, the same is not confirmed for the  
355 other classes. Low stability emerged during cross validation for Ecuador, Gulf of Mexico and Venezuela classes.  
356 Classification results for classes composed of less than 10 samples were not considered statistically valid.

#### 357 **4. Conclusions**

358 Because of the low price and interesting features of CBS, such as the extraordinary similarity to cocoa powder in  
359 terms of color, taste and texture, and the potential beneficial effects on human health, research is needed to assist  
360 the valorization of this food by-product, and to prevent fraud in cocoa powder market. The present work  
361 demonstrates the existence of correlations between the geographical origin and the composition of CBS samples,  
362 even though low specificity for the single country or restricted areas emerged. Some information about what  
363 samples from the same macro-area have in common was described. The selected techniques provided significant

364 criteria to distinguish sample classes, such as Central Africa and Ecuador samples with adequate accuracy and  
365 precision, however it is very difficult to precisely determine what chemical species drive this separation only  
366 using vibrational spectroscopy for chemical composition analysis. Nevertheless, estimates and trends were  
367 determined. The geographical traceability of food based on chemical analysis remains complicated and always  
368 valid rules are rarely identified. The natural variability of most food materials is huge, climatic conditions and  
369 process variables represent an intrinsic limit of this field of study. However, the capability to identify leading  
370 variables, common trends and general indications using rapid and simple techniques is an encouraging result in  
371 this domain. More sensitive and accurate techniques should be used for an exhaustive investigation. Easy-to-use  
372 instrumental analysis still needs the support of heavier analytical strategies for comparison and calibration.

### 373 **Acknowledgements**

374 The present work has been supported by COVALFOOD “Valorisation of high added-value compounds from  
375 cocoa industry by-products as food ingredients and additives” project funded by European Union’s Seventh  
376 Framework programme for research and innovation under the Marie Skłodowska-Curie grant agreement No  
377 609402 - 2020 researchers: Train to Move (T2M).

### 378 **References**

379 Andrade, K. S., Gonçalves, R. T., Maraschin, M., Ribeiro-do-Valle, R. M., Martínez, J., Ferreira, S. R. (2012).  
380 Supercritical fluid extraction from spent coffee grounds and coffee husks: Antioxidant activity and effect of  
381 operational variables on extract composition. *Talanta*, 88, 544-552.

382 Abadin, H., Ashizawa, A., Stevens Y.W., Lladós, F., Diamond, G., Sage, G., Citra, M., Quinones, A., Bosch S.  
383 J., and Swarts, S. G., ATSDR, U. (2007). Toxicological profile for lead. US Department of Health and Human  
384 Services, 1, 582.

385 Ballabio, D., Consonni, V. (2013). Classification tools in chemistry. Part 1: linear models. PLS-DA. *Analytical*  
386 *Methods*, 5(16), 3790-3798.

387 Biancolillo, A., Bucci, R., Magrì, A. L., Magrì, A. D., & Marini, F. (2014). Data-fusion for multiplatform  
388 characterization of an Italian craft beer aimed at its authentication. *Analytica chimica acta*, 820, 23-31.

389 Barbosa-Pereira, L., Guglielmetti, A., & Zeppa, G., 2018. Pulsed Electric Field Assisted Extraction of Bioactive  
390 Compounds from Cocoa Bean Shell and Coffee Silverskin. *Food and Bioprocess Technology*, 11(4), 818-835.

391 Barker, M. and Rayens, W. (2003). Partial least squares for discrimination. *Journal of Chemometrics*, 17(3):  
392 166–173.

393 Barnes, R. J., Dhanoa, M. S., Lister, S. J. (1989). Standard normal variate transformation and de-trending of  
394 near-infrared diffuse reflectance spectra. *Applied spectroscopy*, 43(5), 772-777.

395 Borràs, E., Ferré, J., Boqué, R., Mestres, M., Aceña, L., Busto, O. (2015). Data fusion methodologies for food  
396 and beverage authentication and quality assessment—A review. *Analytica Chimica Acta*, 891, 1-14.

397 Caroch, M., Morales, P., Ferreira, I. C. (2015). Natural food additives: Quo vadis?. *Trends in Food Science &*  
398 *Technology*, 45(2), 284-295.

399 De Oliveira, I. R., Roque, J. V., de Maia, M. P., Stringheta, P. C., & Teófilo, R. F. (2018). New strategy for  
400 determination of anthocyanins, polyphenols and antioxidant capacity of Brassica oleracea liquid extract using  
401 infrared spectroscopies and multivariate regression. *Spectrochimica Acta Part A: Molecular and Biomolecular*  
402 *Spectroscopy*, 194, 172-180.

403 Doeswijk, T. G., Smilde, A. K., Hageman, J. A., Westerhuis, J. A., & Van Eeuwijk, F. A. (2011). On the  
404 increase of predictive performance with high-level data fusion. *Analytica chimica acta*, 705(1-2), 41-47.

405 Dong, Y., Sørensen, K. M., He, S., & Engelsen, S. B. (2017). Gum Arabic authentication and mixture  
406 quantification by near infrared spectroscopy. *Food Control*, 78, 144-149.

407 Elwers, S., Zambrano, A., Rohsius, C., Lieberei, R. (2009), Differences between the content of phenolic  
408 compounds in Criollo, Forastero and Trinitario cocoa seed (*Theobroma cacao* L.), *European Food Research and*  
409 *Technology*, 229(6), 937-948.

410 Jacobsen, S., Søndergaard, I., Møller, B., Desler, T., Munck, L. (2005). A chemometric evaluation of the  
411 underlying physical and chemical patterns that support near infrared spectroscopy of barley seeds as a tool for

412 explorative classification of endosperm genes and gene combinations. *Journal of Cereal Science*, 42(3), 281-  
413 299.

414 Jansman, A. J., Verstegen, M. W., Huisman, J., Van den Berg, J. W. (1995). Effects of hulls of fava beans (*Vicia*  
415 *fabu* L.) with a low or high content of condensed tannins on the apparent ileal and fecal digestibility of nutrients  
416 and the excretion of endogenous protein in ileal digesta and feces of pigs. *Journal of Animal Science*, 73(1), 118-  
417 127.

418 Kelly, S., Heaton, K., Hoogewerff, J. (2005). Tracing the geographical origin of food: The application of multi-  
419 element and multi-isotope analysis. *Trends in Food Science & Technology*, 16(12), 555-567.

420 Lee, A. R., Gautam, M., Kim, J., Shin, W. J., Choi, M. S., Bong, Y. S., Hwang G.S., Lee, K. S. (2011). A  
421 multianalytical approach for determining the geographical origin of ginseng using strontium isotopes,  
422 multielements, and <sup>1</sup>H NMR analysis. *Journal of agricultural and food chemistry*, 59(16), 8560-8567.

423 Li-Chan, E., Chalmers, J., Griffiths, P. (Eds.). (2011). Applications of vibrational spectroscopy in Food Science.  
424 John Wiley & Sons.

425 Luykx, D. M., Van Ruth, S. M. (2008). An overview of analytical methods for determining the geographical  
426 origin of food products. *Food Chemistry*, 107(2), 897-911.

427 Magagna, F., Guglielmetti, A., Liberto, E., Reichenbach, S. E., Allegrucci, E., Gobino, G., ... & Cordero, C.  
428 (2017). Comprehensive Chemical Fingerprinting of High-Quality Cocoa at Early Stages of Processing:  
429 Effectiveness of Combined Untargeted and Targeted Approaches for Classification and Discrimination. *Journal*  
430 *of agricultural and food chemistry*, 65(30), 6329-6341.

431 Mandrile, L., Fusaro, I., Amato, G., Marchis, D., Martra, G., & Rossi, A. M. (2018). Detection of insect's meal  
432 in compound feed by Near Infrared spectral imaging. *Food Chemistry*.

433 Mandrile, L., Zeppa, G., Giovannozzi, A. M., & Rossi, A. M. (2016). Controlling protected designation of origin  
434 of wine by Raman spectroscopy. *Food chemistry*, 211, 260-267.

435 Manzano, P., Hernández, J., Quijano-Avilés, M., Barragán, A., Chóez-Guaranda, I., Viteri, R., Valle, O. (2017).  
436 Polyphenols extracted from Theobroma cacao waste and its utility as antioxidant. *Emirates Journal of Food and*  
437 *Agriculture*, 29(1), 45.

438 Marseglia, A., Acquotti, D., Consonni, R., Cagliani, L. R., Palla, G., & Caligiani, A. (2016). HR MAS 1H NMR  
439 and chemometrics as useful tool to assess the geographical origin of cocoa beans—Comparison with HR 1H  
440 NMR. *Food Research International*, 85, 273-281.

441 Martens, H., Nielsen, J. P., & Engelsen, S. B. (2003). Light scattering and light absorbance separated by  
442 extended multiplicative signal correction. Application to near-infrared transmission analysis of powder mixtures.  
443 *Analytical Chemistry*, 75(3), 394-404.

444 Martín- Cabrejas, M. A., Valiente, C., Esteban, R. M., Mollá, E., Waldron, K. (1994). Cocoa hull: a potential  
445 source of dietary fibre. *Journal of the Science of Food and Agriculture*, 66(3), 307-311.

446 Munck, L., Nørgaard, L., Engelsen, S. B., Bro, R., Andersson, C. A. (1998). Chemometrics in food science—a  
447 demonstration of the feasibility of a highly exploratory, inductive evaluation strategy of fundamental scientific  
448 significance. *Chemometrics and Intelligent Laboratory Systems*, 44(1), 31-60.

449 Peng, J., Peng, S., Jiang, A., Wei, J., Li, C., & Tan, J. (2010). Asymmetric least squares for multiple spectra  
450 baseline correction. *Analytica chimica acta*, 683(1), 63-68.

451 Peres, B., Barlet, N., Loiseau, G., Montet, D. (2007). Review of the current methods of analytical traceability  
452 allowing determination of the origin of foodstuffs. *Food Control*, 18(3), 228-235.

453 Redgwell, R., Trovato, V., Merinat, S., Curti, D., Hediger, S., Manez, A. (2003). Dietary fibre in cocoa shell:  
454 characterisation of component polysaccharides. *Food Chemistry*, 81(1), 103-112.

455 Rubio- Diaz, D. E., Rodriguez- Saona, L. E. (2010). Application of Vibrational Spectroscopy for the Study of  
456 Heat- Induced Changes in Food Components. *Handbook of Vibrational Spectroscopy*.

457 Schwanninger, M., Rodrigues, J. C., Fackler, K. (2011). A review of band assignments in near infrared spectra of  
458 wood and wood components. *Journal of Near Infrared Spectroscopy*, 19(5), 287-308.

- 459 Serra Bonvehí, J., and Escolá Jordà, R. (1998). Constituents of Cocoa Husks, *Z. Naturforsch.* 53c, 785-792.
- 460 Silvestri, M., Elia, A., Bertelli, D., Salvatore, E., Durante, C., Vigni, M. L., ... & Cocchi, M. (2014). A mid level  
461 data fusion strategy for the Varietal Classification of Lambrusco PDO wines. *Chemometrics and Intelligent*  
462 *Laboratory Systems*, 137, 181-189.
- 463 Skov T., Honoré A.H., Hansen H.M., Næs T., S.B. Engelsen,, (2014). Chemometrics in Foodomics: Handling  
464 data structures from multiple analytical platforms, TRAC-Trends. *Analytical Chemistry*, 60, 71-79.
- 465 Socrates, G. (2001). Infrared and Raman characteristic group frequencies: tables and charts. John Wiley & Sons.
- 466 Sørensen, K. M., Khakimov, B., & Engelsen, S. B. (2016). The use of rapid spectroscopic screening methods to  
467 detect adulteration of food raw materials and ingredients. *Current Opinion in Food Science*, 10, 45-51.;
- 468 Sørensen, K. M., Aru, V., Khakimov, B., Aunskjær, U., & Engelsen, S. B. (2018). Biogenic Amines: a key  
469 freshness parameter of animal protein products in the coming circular economy. *Current Opinion in Food*  
470 *Science*.
- 471 Wold S., Esbensen K., Geladi P. (1987). Principal component analysis. *Chemometrics and Intelligent*  
472 *Laboratory Systems*, 2, 37-52.
- 473 Zakaria, A., Shakaff, A.Y.M., Adom, A.H., Ahmad, M., Masnan, M.J., Aziz, A.H.A., Fikri, N.A., Abdullah,  
474 A.H. and Kamarudin, L.M. (2010). Improved classification of *Orthosiphon stamineus* by data fusion of  
475 electronic nose and tongue sensors. *Sensors*, 10(10), 8782-8796.

## 476 **FIGURE CAPTIONS**

477 **Figure 1**–a) Mean NIR spectrum of all CBS samples (green) and standard deviation limits (blue); b) Scores plot  
478 of NIRS data PCA colored in accordance with variety; c, d) Zoom of average spectrum of Arriba samples  
479 compared with the mean spectrum calculated considering all other NIR spectra.

480 **Figure 2**– a) PC2/PC1 scores plot of NIR spectra of CBS sample colored by geographical origin. b)  
481 PC4/PC5/PC6 scores plot of NIR spectra of CBS sample colored by geographical origin. c) Average NIR spectra  
482 of CBS from Africa and America as macro-classes (red and green respectively) and mean spectra of São Tomé

483 and Ecuador groups (light blue and orange respectively); d, e) Zoom on the spectral regions which make Asian  
484 samples different from all other CBS samples;

485 **Figure 3**– a) ATR-FT-IR average spectrum of all CBS samples (green) and standard deviation limits (blue); b)  
486 PC2 scores plot which highlight common behavior of African samples; c) PC5/PC6 scores plot that allow to  
487 highlight characteristic trend for Ecuador samples; d) MIR average spectra of CH<sub>x</sub> stretching bands of samples  
488 different geographical origin; e) MIR average spectra of Ecuador sample compared with Americans in the  
489 spectral region where Ecuador samples show distinct characteristics with respect to American samples.

490 **Figure 4**– PCA model of ICP-OES data outputs, 2D a) loading and b) scores plots; c) Histogram of mean data  
491 for the considered macro-classes (Africa and America) and São Tomé samples that show peculiar feature with  
492 respect to others; d) Variance captured per each principal component.

493 **Figure 5**– Joined PCA model of NIRS+ICP+MIRS, a) loadings and b) scores plot on PC1 and PC2.

494 **Table 1**–Cross Validation outputs of PLS-Discriminant Analysis classification models for geographical origin  
495 discrimination: a) Joined classification model with 5 LVs, classification performances in leave-one origin-out  
496 cross validation; b) NIRS PLS-DA model with 4 LVs classification performances in leave-one origin-out cross  
497 validation; c) MIRS PLS-DA model with 4 LVs classification performances in leave-one origin-out cross  
498 validation; d) ICP-OES PLS-DA model with 3 LVs classification performances in leave-one origin-out cross  
499 validation.

500

501

502

503

504

505

**Declaration of interests**

The authors declare that they have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper.

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

Table1

Class	Technique	N	Sensitivity	Specificity	Accuracy	Precision
			(true positive ratio)	(true negative ratio)		
Central Africa	a) <b>Joined</b>	22	0.68	0.92	0.84	0.79
	b) NIRS	19	0.68	0.86	0.81	0.00
	c) MIRS	19	0.32	0.70	0.59	0.29
	d) ICP-OES	19	0.50	0.83	0.75	0.50
Gulf of Mexico	a) <b>Joined</b>	9	0.33	0.82	0.76	0.21
	b) NIRS	9	0.00	0.87	0.75	0.00
	c) MIRS	9	0.00	0.82	0.71	0.00
	d) ICP-OES	9	0.00	0.87	0.75	0.00
São Tomé	a) <b>Joined</b>	6	0.33	0.95	0.9	0.40
	b) NIRS	6	0.00	0.91	0.86	0.00
	c) MIRS	6	0.00	0.90	0.83	0.00
	d) ICP-OES	6	0.00	0.92	0.86	0.00
Venezuela	a) <b>Joined</b>	10	0.10	0.87	0.76	0.11
	b) NIRS	12	0.00	0.89	0.74	0.00
	c) MIRS	4	0.00	0.89	0.84	0.00
	d) ICP-OES	12	0.00	0.85	0.69	0.00
Ecuador	a) <b>Joined</b>	10	0.00	0.87	0.74	0.00
	b) NIRS	10	0.00	0.85	0.72	0.00
	c) MIRS	10	0.00	0.81	0.70	0.00
	d) ICP-OES	10	0.00	0.87	0.73	0.00
Indonesia	a) <b>Joined</b>	1	0.00	0.96	0.94	0.00
	b) NIRS	1	0.00	1.00	0.99	0.00
	c) MIRS	1	0.00	1.00	0.99	0.00
	d) ICP-OES	1	0.00	0.98	0.97	0.00
Mexico	a) <b>Joined</b>	2	0.00	0.99	0.96	0.00
	b) NIRS	2	0.00	0.96	0.93	0.00
	c) MIRS	2	0.00	0.94	0.91	0.00
	d) ICP-OES	2	0.00	0.97	0.94	0.00
Peru	a) <b>Joined</b>	4	0.00	0.89	0.84	0.00
	b) NIRS	4	0.00	0.92	0.87	0.00
	c) MIRS	4	0.00	0.97	0.91	0.00
	d) ICP-OES	4	0.00	0.87	0.81	0.00
Colombia	a) <b>Joined</b>	4	0.00	0.95	0.90	0.00
	b) NIRS	4	0.00	0.92	0.87	0.00
	c) MIRS	12	0.00	0.93	0.77	0.00
	d) ICP-OES	4	0.00	0.85	0.80	0.00

**Table 1:** Cross Validation outputs of PLS-Discriminant Analysis classification models for geographical origin discrimination: a) Joined classification model with 5 LVs, classification performances in leave-one origin-out cross validation; b) NIRS PLS-DA model with 4 LVs classification performances in leave-one origin-out cross validation; c) MIRS PLS-DA model with 4 LVs classification performances in leave-one origin-out cross validation; d) ICP-OES PLS-DA model with 3 LVs classification performances in leave-one origin-out cross validation.

Figure 1  
[Click here to download high resolution image](#)

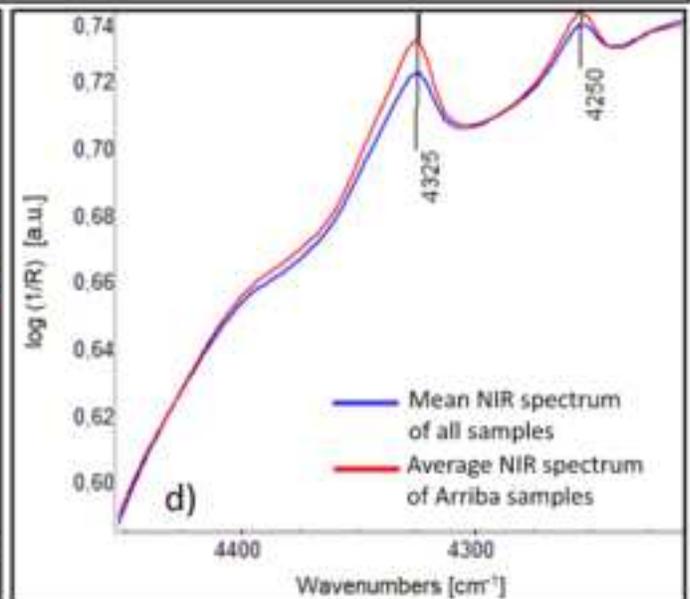
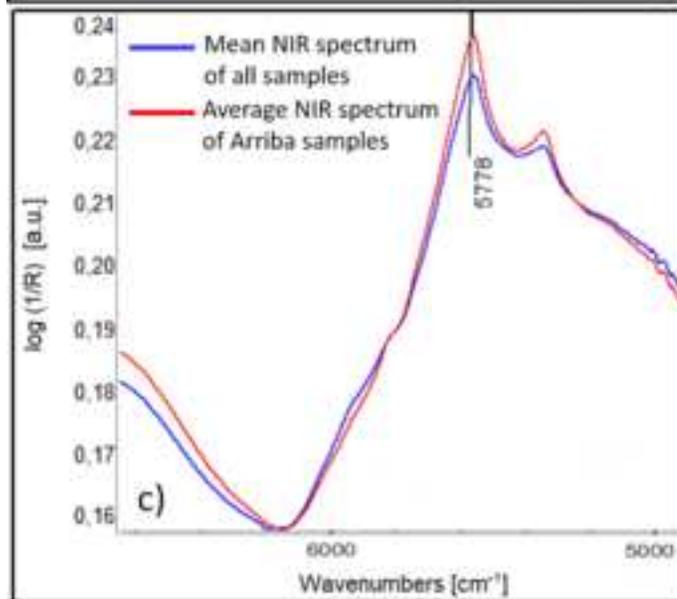
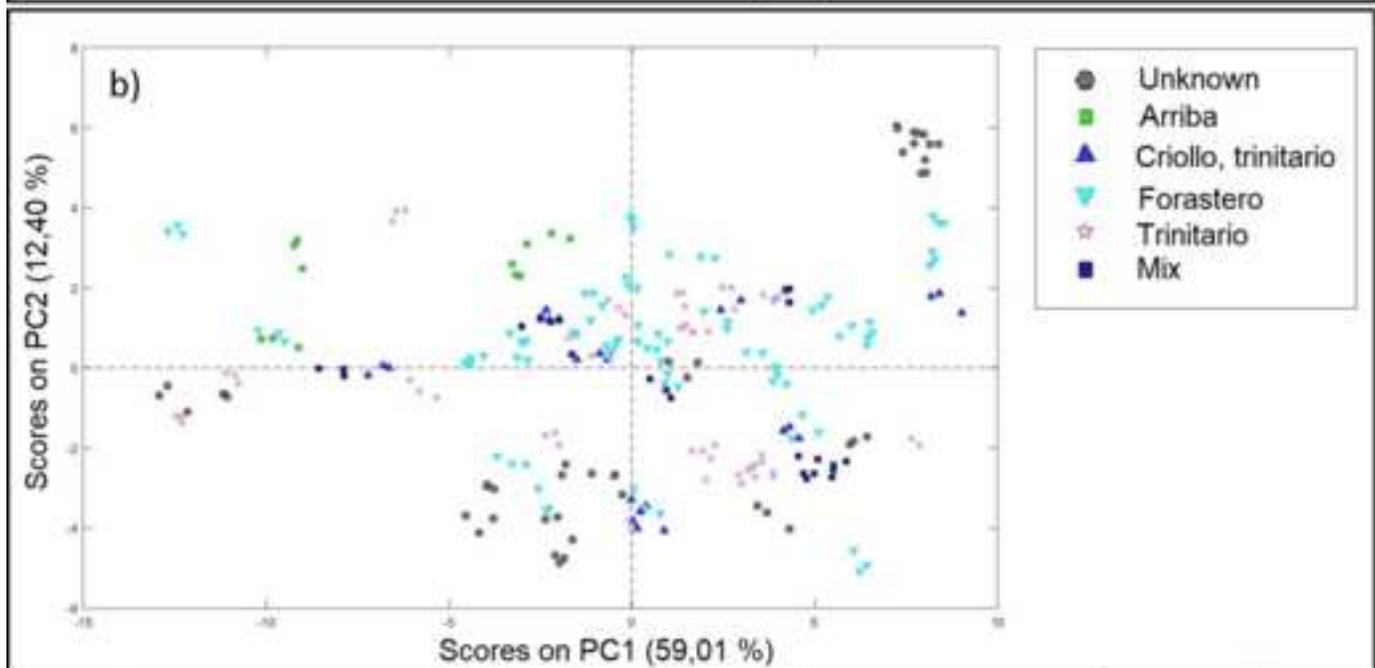
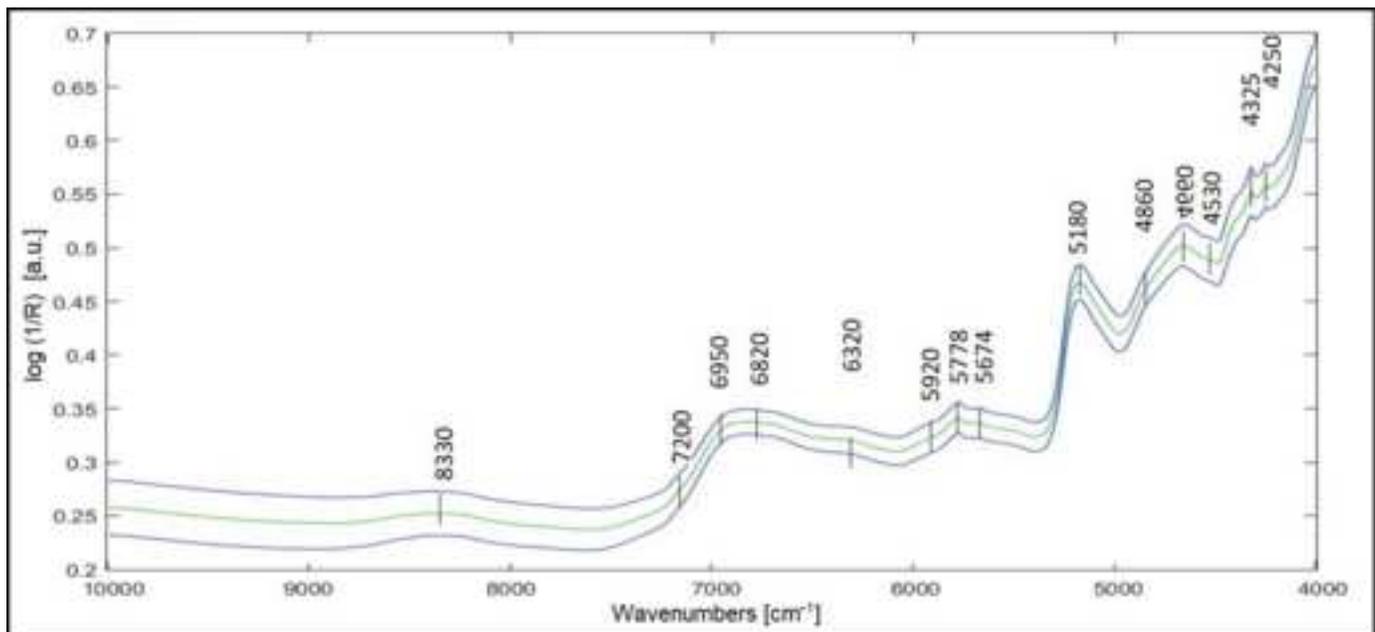


Figure 2

[Click here to download high resolution image](#)

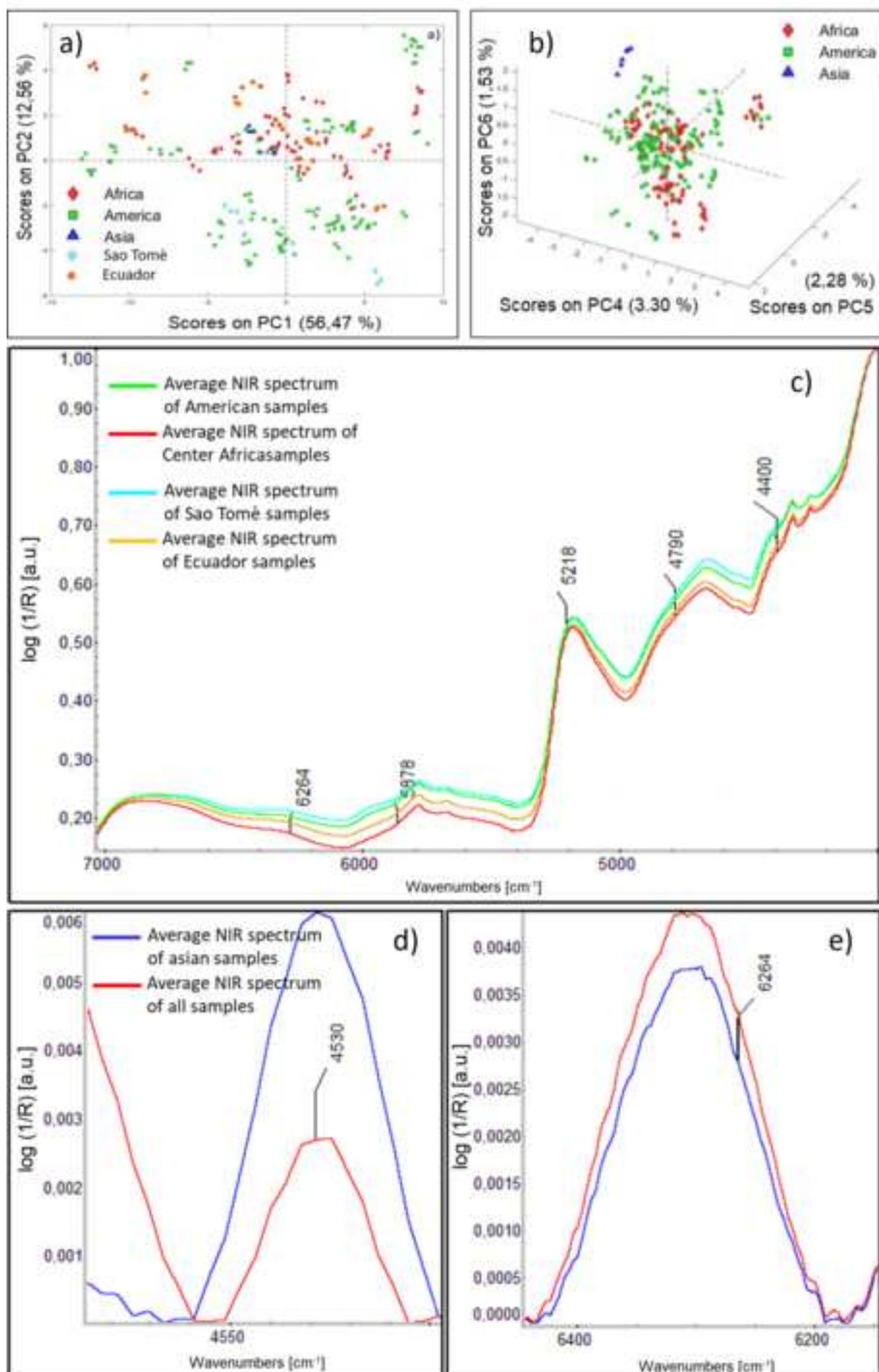


Figure 3  
[Click here to download high resolution image](#)

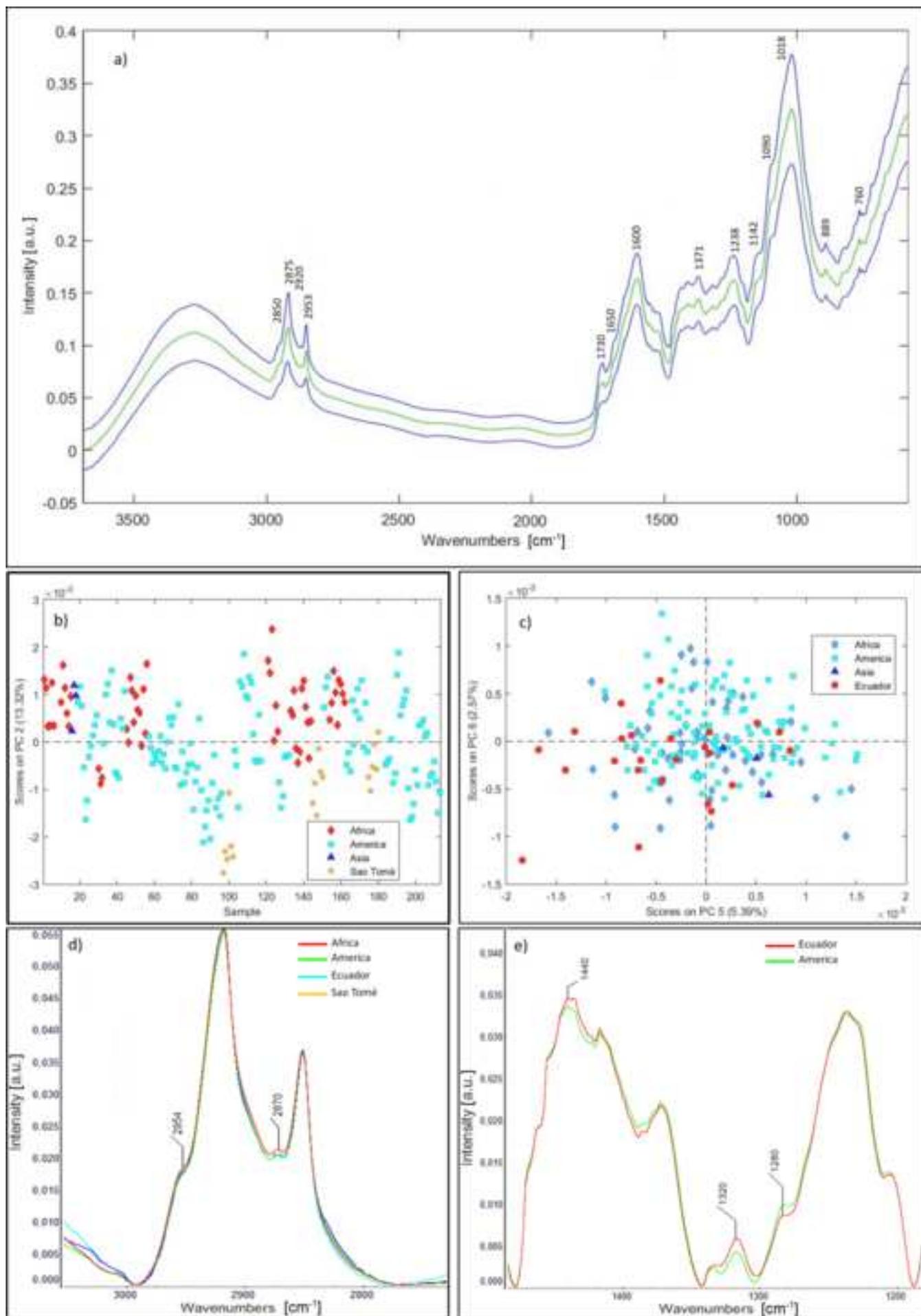


Figure 4  
[Click here to download high resolution image](#)

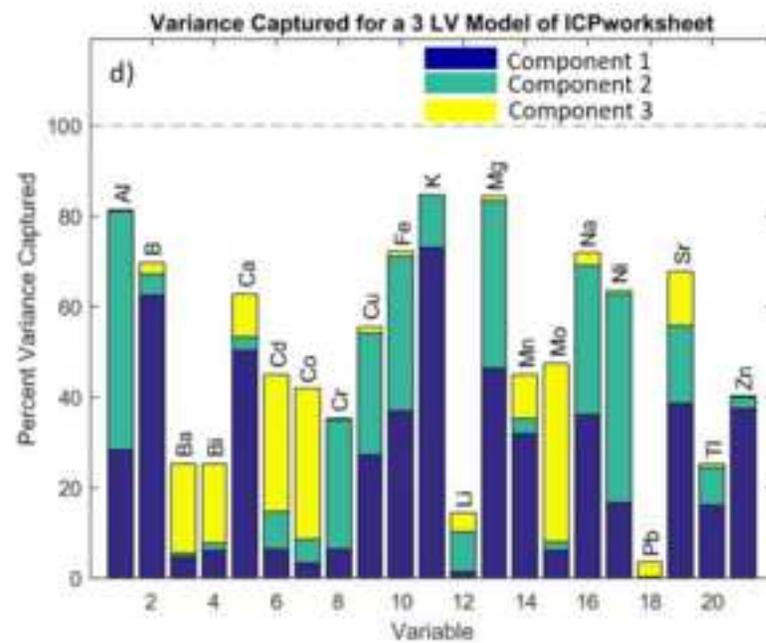
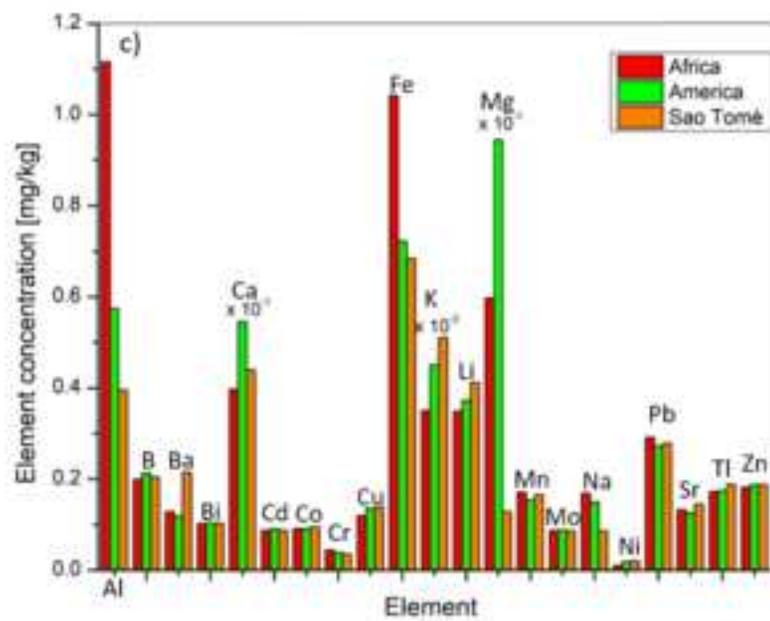
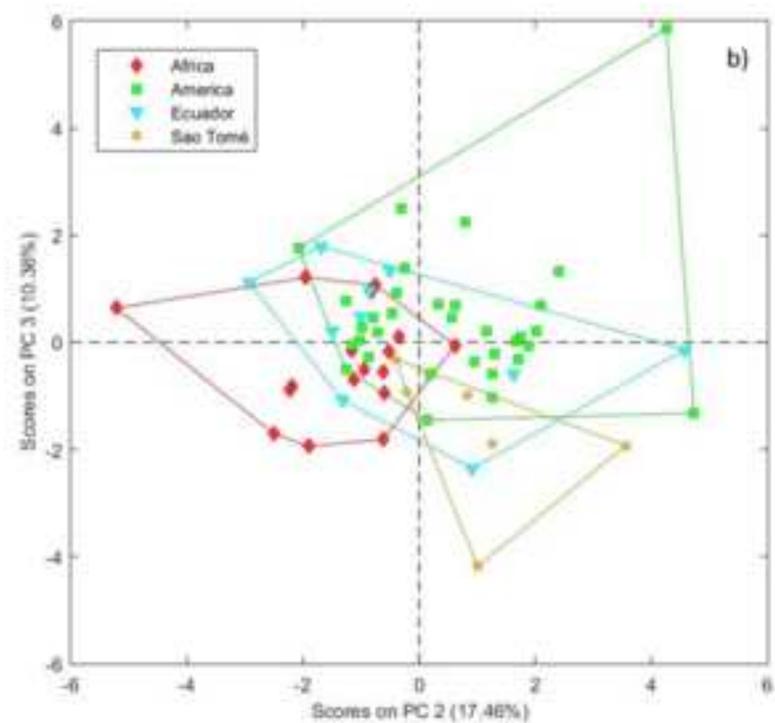
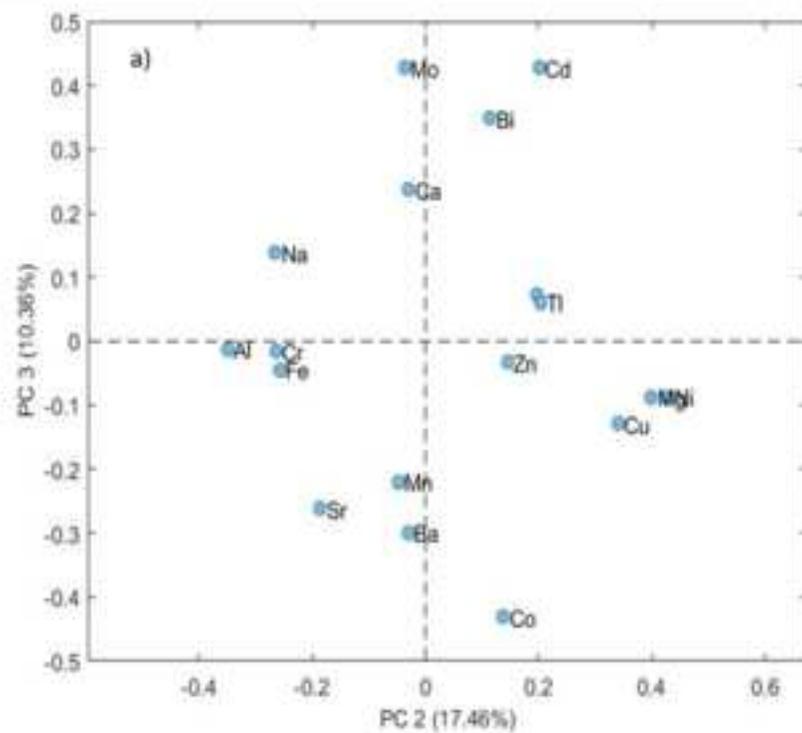
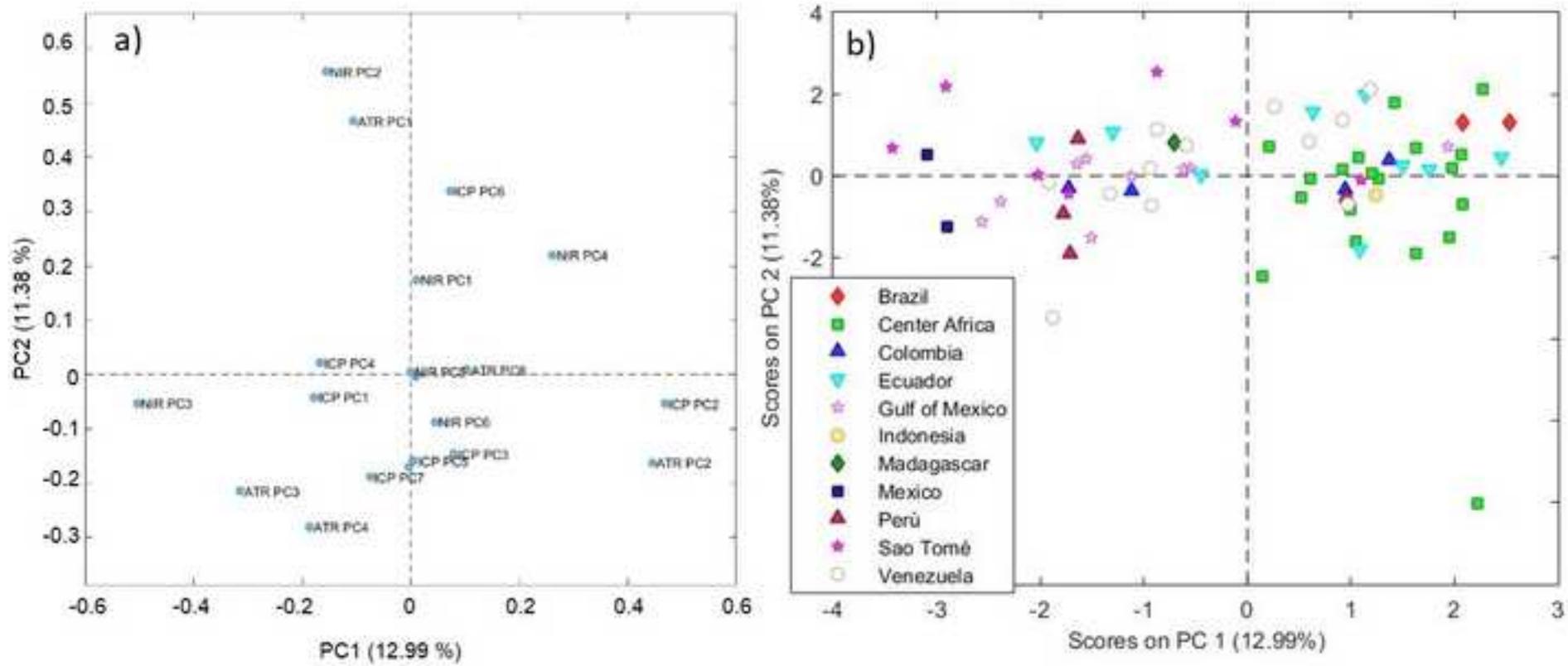


Figure 5  
[Click here to download high resolution image](#)



**Supplementary Material**

[Click here to download Supplementary Material: Supplementary information\\_NEW.docx](#)