

# Determination of the geographical origin of green coffee beans using NIR spectroscopy and multivariate data analysis

*A. Girauda<sup>a</sup>, S. Grassi<sup>b</sup>, F. Savorani<sup>a\*</sup>, G. Gavoci<sup>a</sup>, E. Casiraghi<sup>b</sup>, F. Geobaldo<sup>a</sup>*

<sup>a</sup> Department of Applied Science and Technology, Polytechnic of Turin, Corso Duca degli Abruzzi 24, I-10129, Turin,  
Italy.

<sup>b</sup> Department of Food, Environmental and Nutritional Sciences, University of Milan, Via G. Celoria 2, I-20133, Milan,  
Italy.

\*Corresponding author. E-mail: [francesco.savorani@polito.it](mailto:francesco.savorani@polito.it); Phone: 011 0904562

Authors' e-mail:

Alessandro Girauda, [alessandro.girauda@polito.it](mailto:alessandro.girauda@polito.it)

Silvia Grassi, [silvia.grassi@unimi.it](mailto:silvia.grassi@unimi.it)

Gentian Gavoci, [gentian.gavoci@polito.it](mailto:gentian.gavoci@polito.it)

Ernestina Casiraghi, [ernestina.casiraghi@unimi.it](mailto:ernestina.casiraghi@unimi.it)

Francesco Geobaldo, [francesco.geobaldo@polito.it](mailto:francesco.geobaldo@polito.it)

1 **Abstract**

2 In this work, near infrared (NIR) spectroscopy and multivariate data analysis were investigated as a  
3 fast and non-disruptive method to classify green coffee beans on continents and countries bases.  
4 FT-NIR spectra of 191 coffee samples, origin from 2 continents and 9 countries, were acquired by  
5 two different laboratories.

6 Laboratory-independent Partial Least Square-Discriminant Analysis and interval PLS-DA models  
7 were developed by following a hierarchical approach, i.e. considering at first the continent and then  
8 the country of origin as discrimination rule.

9 The best continent-based classification model was able to identify correctly more than 98% in  
10 prediction, whereas 100% of them were correctly predicted by the best country-based classification  
11 model. The inter-laboratory reliability of the proposed method was confirmed by McNemar test,  
12 since no significant differences ( $P>0.05$ ) were found. Furthermore, a validation was performed  
13 predicting the spectral test set of a laboratory using the model developed by the other one.

14

15

16 **Keywords**

17 Geographical origin; Green coffee beans; NIR Spectroscopy; Chemometrics; Classification;  
18 Variable selection

19

20

21

## 22 1. Introduction

23  
24 The sensory properties of coffee, one of the most popular beverages in the world, are deeply  
25 affected by the chemical composition of the raw coffee beans that is, in turn, highly related to their  
26 geographical growing regions (Esteban-Diez et al., 2007). Coffee plantations are primarily found in  
27 the latitude between 20°N and 20°S, a geographical band encircling the world that is commonly  
28 known as “coffee belt”. This band encompasses regions (belonging to Centre and South-America,  
29 Africa and Asia) roughly bounded by the tropics of Cancer and Capricorn, which offer optimal  
30 climate conditions for growing coffee plants. However, a great variance in coffee quality, taste and  
31 body is found within these world’s regions (Anderson and Smith, 2002). Unavoidably, this  
32 variability aspect also implies considerable differences in terms of commercial values of the product  
33 (Teuber, 2010) and has led, in many cases, to common forms of frauds like mislabelling, i.e.  
34 disguising (not declaring) the right geographical or botanical origin of the coffee beans, and  
35 adulteration, i.e. selling cheaper coffee, or mix of coffees, as pure expensive species/categories  
36 (Alonso-Salces et al., 2009). As a consequence, coffee producers, as well as industrial  
37 manufacturers, are increasingly interested in protecting the market reputation from the  
38 aforementioned socioeconomic issues and have highly encouraged the development of efficient  
39 analytical methods able to evaluate the coffee authenticity (Huck et al., 2005).

40 In this context, during the last couple of decades, several analytical techniques, such as gas and  
41 liquid chromatography (Frega et al., 1995; Bicchi et al., 1997; Gonzalez et al., 2001; Casal et al.,  
42 2003; Martin et al., 2001; Carrera et al., 1998), mass spectrometry (Gil-Agusti et al., 2005) and  
43 nuclear magnetic resonance spectroscopy (Cagliani et al., 2013; Arana et al., 2015) were tested  
44 allowing gathering accurate information about the coffee composition, aiming at assessing its  
45 authenticity.

46 Unfortunately, these methods appear quite expensive, complex to use, often time consuming and  
47 requiring a sample preparation step before analysis that involves, in turn, the use of different kinds  
48 of chemical solvents. In that respect, alternative spectroscopic techniques, such as Raman and near-  
49 infrared (NIR) spectroscopy, provide a valid solution to overcome some of the abovementioned  
50 drawbacks, since they allow performing green, simple, fast and non-invasive analysis directly *in*  
51 *situ* (Marquetti et al., 2016). This is the reason why, over the years, they have found increasing  
52 application as routine analysis in different fields, such as biology, medicine, food science and  
53 technology (Ozaki et al., 2006).

54 For what it concerns the coffee production chain, Raman spectroscopy has been investigated for the  
55 chemical discrimination of Arabica and Robusta species (Rubayiza et al., 2005; El-Abassy et al.,

2011), while NIR spectroscopy has been successfully applied with different purposes, from the coffee varietal differentiation (Kemsley et al. 1995, Esteban-Díez et al., 2007, Buratti et al., 2015) to the prediction of sensory properties of coffee beverage (Esteban-Díez et al., 2004; Ribeiro et al., 2011) and of the coffee roasting degree (Alessandrini et al., 2008; Bertone et al., 2016). Surprisingly, only a few research studies, regarding the use of NIR spectroscopy for the direct determination of the geographic origin of coffee, were found in literature. One of these studies focused on evaluating the variability of coffee samples cultivated in different regions of a single country (Marquetti et al., 2016). Another relevant study, performed by Medina et al., 2017, compared the ability of different spectroscopic techniques for the determination of the geographical origin of roasted coffee samples and was aimed at protecting the Colombian coffee origin. However, none of the proposed studies considered inter-laboratory results comparison, which is instead a relevant point when assessing an analytical method to be a reliable and accessible approach. Boix et al., 2012 evaluated the transfer of a NIR microscopy method for the detection of animal products in feedingstuffs performing Analysis of Variance (ANOVA) on the laboratories' duplicate analyses. A robust Ring test was performed involving 16 laboratories for the prediction of a significant number of wine quality parameters by Fourier transform-middle infrared spectrometry (Patz et al., 2004). However, to the best of our knowledge, this important aspect is usually underestimated in research studies dealing with food classification models based on NIR spectroscopy.

Furthermore, the NIR application in coffee field demands for the development of simplified devices (Pizarro et al., 2007; Calvini et al., 2017) reducing the costs in terms of money and time but leading to high quality results. The realization of *ad hoc* instruments passes through the identification of relevant spectral region, which can be achieved by variable selection strategies.

In the present work, NIR spectroscopy and multivariate data analysis were investigated as a fast and non-invasive method to classify green coffee beans cultivated in different world's continents and countries and to select relevant spectral features in vision of simplified instrument development.

In details, after data exploration by Principal Component Analysis (PCA), Partial Least Square-Discriminant Analysis (PLS-DA) algorithm was applied to develop a series of classification models. In this context, a hierarchical approach has been followed, considering at first the continent of origin and then the country of origin of the coffee samples as discrimination criteria. Moreover, interval Partial Least Square-Discriminant Analysis (iPLS-DA) algorithm was investigated to select the most informative spectral regions, and McNemar test was performed to compare the inter-laboratory model performance.

89

90  
91  
92  
93  
94  
95  
96  
97  
98  
99  
100  
101  
102  
103  
104  
105  
106  
107  
108  
109  
110  
111  
112  
113  
114  
115  
116  
117  
118  
119  
120  
121  
122

**2. Materials and methods**

*2.1 Coffee samples*

One hundred and ninety-one (191) samples of green coffee beans, representative of different geographical origins, were considered in the present study. Out of these, 88 coffee samples came from Centre-South America and 103 samples were grown in Asian countries; the details about the origin are provided in Table 1.

The coffee samples were representative of the main botanical species, i.e. *Robusta* and natural or washed *Arabica*. The producing countries and the species were chosen according to their relevance for the Italian coffee market.

The preliminary post-harvest processing operations on the green coffee beans, i.e. washing of cherry fruits and drying, were performed directly in the producing country.

Two replicates of sample, each one consisting of three hundred grams of coffee beans, were packed in vacuum in light barrier packaging material. The samples were prepared over a period of two months and delivered concurrently, within seven days from the preparation, in eight subsequent shipping to two different laboratories, i.e. the Department of Food, Environmental and Nutritional Sciences (UNIMI, University of Milan) and the Department of Applied Science and Technology (POLITO, Polytechnic of Turin).

**Table 1**

*2.2 NIR spectra acquisition*

NIR analysis was performed within 48 hours from the delivery of the coffee samples in both laboratories independently. The samples were analysed using two different Fourier Transform (FT)-NIR spectrometer (MPA, Bruker Optics, Ettlingen, Germany) provided with integrating sphere. The NIR spectra were collected at room temperature (~20° C) in diffuse reflectance mode directly on the entire coffee beans, hence without any kind of sample pre-treatment. Each sample was always carefully mixed before sampling eighty grams of coffee beans for each of the three replicates. During measurement, the sample holder was kept in rotation to collect more representative spectra on the whole sample. Spectral data were collected over the range 12500 to 3600 cm<sup>-1</sup> (resolution, 8 cm<sup>-1</sup>; scanner velocity, 10kHz; background, 256 scans; sample, 64 scans). Background scans were

123 performed using instrumental internal reference standard. OPUS software (v. 6.5, Bruker Optics,  
124 Ettlingen, Germany) was used for instrumental control and for spectra acquisition.

125

### 126 *2.3 Multivariate data analysis*

127

128 Chemometric elaboration of the NIR spectra was performed under Matlab environment (ver. 9.2,  
129 The Mathworks, Inc., Natick, USA) using the PLS toolbox (ver. 8.5, Eigenvector Research, Inc.,  
130 USA) software package. The three replicates of NIR spectra collected on each coffee sample were  
131 averaged. Therefore, a global dataset, containing a unique, representative spectrum per sample, was  
132 obtained separately from both the laboratories and used independently for all the subsequent data  
133 elaboration steps.

134 Initially, different pre-processing methods were investigated aiming at eliminating potential  
135 artefacts from NIR spectra and correcting their nonlinear behaviour: Standard Normal Variate  
136 (SNV), Multiple Scattering Correction (MSC), as well as 1<sup>st</sup> or the 2<sup>nd</sup> derivative transformation  
137 using the Savitzky-Golay algorithm, applied alone or in combination with SNV or MSC. No matter  
138 the pre-treatment applied, all the data were subjected to mean centering prior to any multivariate  
139 data analysis.

140

#### 141 *2.3.1 Exploratory data analysis*

142

143 Explorative data analysis was carried out by Principal Component Analysis (PCA), based on the  
144 singular value decomposition (SVD) algorithm (Wall et al., 2003), on all the datasets obtained by  
145 applying the different abovementioned pre-processing techniques. The optimal number of Principal  
146 Components (PCs) to be retained in the PCA model was chosen from the analysis of the  
147 corresponding scree plot (Bro and Smilde, 2014).

148

#### 149 *2.3.2 PLS-DA classification models*

150

151 The classification models were developed using Partial Least Square – Discriminant Analysis (PLS-  
152 DA) algorithm (Wold et al., 1984).

153 The samples were divided into the proper number of classes according to the investigated  
154 discrimination criteria. In the first case, i.e. the continent-based discrimination, a vector of two *a-*  
155 *priori* classes was considered, where the first class referred to the American coffee samples and the  
156 other class referred to the Asian samples.

157 In the second case, i.e. the country-based discrimination, five classes were considered. Four of them  
158 corresponded to the four most representative producing countries with respect to the total number of  
159 samples reported in Table 1, i.e. Brazil (38), Honduras (27), India (65) and Vietnam (31). Class 5  
160 groups 30 samples belonging to different countries, beings each country samples' number not  
161 enough to build an independent class maintaining a balanced design. In detail, this class includes  
162 samples produced in Guatemala (11), Colombia (8), Costa Rica (2), Nicaragua (2) and Indonesia  
163 (7), respectively.

164 Approximately the 75% of each sample dataset (144 samples, belonging to six deliveries) was used  
165 as calibration set and the remaining 25% (47 samples, belonging to two deliveries) was used as  
166 external validation set.

167 Models were tested internally by cross validating the calibration set using six cancellation groups,  
168 each group corresponding to one of the six deliveries. Each time, the samples included in the same  
169 cancellation group formed the prediction set, while all the other ones were used to build the model.  
170 Subsequently, the prediction capability of the achieved models was validated on the external test  
171 set.

172 The classification performance of the PLS-DA models was evaluated by comparing the actual (or  
173 reference) class to the class predicted by the model. In this context, four different conditions may  
174 occur based on the model prediction, since the samples can result “true positive” (TP), “true  
175 negative” (TN), “false positive” (FP) or “false negative” (FN) (Szymanska et al., 2012; Ballabio  
176 and Consonni, 2013). A series of statistical parameters were then calculated for all the class  
177 separately. They are:

- 178 - Sensitivity (*SENS*, Equation 1), which expresses the model capability to correctly recognize  
179 samples belonging to the considered class;
- 180 - Specificity (*SPEC*, Equation 2), which describes the model capability to correctly reject  
181 samples belonging to all the other classes;
- 182 - Efficiency (*EFF*), calculated as the geometric mean of *SPEC* and *SENS* (Equation 3).

183

$$184 \quad SENS = \frac{TP}{(TP+FN)} \quad (\text{Eq. 1})$$

185

$$186 \quad SPEC = \frac{TN}{(TN+FP)} \quad (\text{Eq. 2})$$

187

$$188 \quad EFF = \sqrt{(SPEC * SENS)} \quad (\text{Eq. 3})$$

189

190 All these parameters can assume values between 0 (0%) and 1 (100%) and were calculated referring  
191 to the calibration (TRN, CAL), to the cross-validation of the training set (TRN, CV), and to the  
192 prediction of the external test set (TST, PRED).

193 Eventually, the classification results obtained were compared in terms of prediction performance by  
194 applying the ‘testcholdout’ function implemented in Matlab (v. 2016a, The Mathworks, Inc.,  
195 Natick, USA), which performs one-tailed, mid P-value McNemar test (Fagerlan et al., 2013), a  
196 particular case of Fisher’s sign test that verifies if two models have the same error rate (Roggo et  
197 al., 2003). For further details about the statistical test at hand, the reader is referred to Grassi et al.  
198 (2018a). In our case, the algorithm has been implemented to assess which data pre-treatment  
199 technique and/or combination of them gave the best results. Moreover, the test allowed to assess the  
200 inter-laboratory accuracy; that is whether the accuracy of the classification models obtained using  
201 the spectra collected with two different instruments by the same brand, but located in two different  
202 laboratories and analysed by two independent operators, were different.

203

### 204 2.3.3 *Variable selection*

205

206 In addition to the different pre-treatments techniques investigated, a variable selection procedure  
207 was performed using interval Partial Least Square-Discriminant Analysis (iPLS-DA) algorithm  
208 (Nørgaard et al., 2000), in order to identify the best subset of spectral variables to consider for the  
209 elaboration of the PLS-DA classification models. The iPLS-DA algorithm, as well as other interval-  
210 based chemometric techniques available in literature (Savorani et al., 2013), consists essentially in  
211 subdividing the whole spectral range in equidistant intervals (with arbitrary-defined length) and  
212 building a series of ‘local’ models, first using one interval at a time and then adding the intervals in  
213 an iterative way. The best model in terms of selected variables is defined by the lowest Root Mean  
214 Square Error value achieved in cross-validation (RMSECV) (Leardi and Nørgaard, 2004). In the  
215 present work, four different interval size values ( $i$ ) were considered, with length of 80, 40, 20 or 10  
216 spectral variables, respectively.

217

### 218 2.3.4 *Cross-laboratory model validation*

219

220 As a final step of data elaboration, a validation of the classification models was performed by  
221 predicting the coffee origin, always considering both continent and country, using at the same time  
222 the NIR spectra acquired from UNIMI and POLITO laboratories. At first, the best iPLS-DA model  
223 developed with the POLITO calibration set was validated on the UNIMI external test set.



224 Subsequently, the model build using the UNIMI calibration set and the same spectral variables  
225 selected for the abovementioned iPLS-DA model was used to predict the POLITO external test set.  
226 Weighted average efficiency in prediction ( $EFF_{TS}$ ) was calculated to assess models' reliability.

227  
228

### 229 **3. Results and discussion**

230

#### 231 *3.1 Spectral features*

232

233 FT-NIR spectra collected in both laboratories are characterised by absorption bands related to  
234 different vibrations of the chemical bonds. No matter the country of origin, the main observable  
235 absorption bands are related to different organic chemical compounds of coffee, mainly caffeine,  
236 lignin and fatty acids (Figure 1a). In particular, absorption associated with the C- H bonds of  
237 caffeine can be distinguished between 5600 and 4000  $cm^{-1}$ ; in the same area contributes the  
238 absorption of C- H bonds of fatty acids, amino acids and lignin, which also show absorption bands  
239 around 8300-8200  $cm^{-1}$ . O-H absorptions are observable around 10000  $cm^{-1}$ , 6800  $cm^{-1}$  and 4800  
240  $cm^{-1}$  and are mainly linked to cellulose. Moreover, cellulose signal around 4400  $cm^{-1}$  is ascribable  
241 to the C- H and O- H bonds. Absorptions related to N- H bonds of polyamides (6800  $cm^{-1}$ ) and  
242 proteins (6400, 5700 and 4800  $cm^{-1}$ ) are also distinguishable (Buratti et al., 2015; Weyer and  
243 Workman, 2007).

244 Even if the spectral features of coffee beans are the same, few differences exist between the SNV  
245 transformed average spectrum of the American samples and the one of the Asian samples (Figure  
246 1b). In particular, the average spectrum of the American samples shows higher absorption for the  
247 broad band at 8450-8000  $cm^{-1}$ , for the band at 7000-6500  $cm^{-1}$  and the consecutive shoulder up to  
248 6100  $cm^{-1}$  and again between 6000-5400  $cm^{-1}$ , in agreement with what observed by Marquetti et al.,  
249 2016. Conversely, the average spectrum of the Asian samples is characterized by higher absorbance  
250 values from 12500 to 10500  $cm^{-1}$ , which can be related to colour effect from the adjacent visible  
251 region. Moreover, higher absorbance for these samples can be observed at around 5200  $cm^{-1}$  and for  
252 the peak at 4200  $cm^{-1}$ .

253

### 253 **Figure 1**

254

#### 255 *3.2 Exploratory data analysis*

256

257 PCA was applied on all the pre-processed NIR dataset to assess which kind of pre-treatments better  
258 separated the coffee samples based on both the continent and the country of origin.

259 Both the PCA models calculated independently from UNIMI and POLITO laboratories led to a  
260 rather satisfactory discrimination after the combined application of SNV and mean centering on the  
261 raw NIR spectra. The plots reported in Figure 2 refer to the spectra collected from POLITO. It is  
262 noteworthy that the same clustering trends shown in Figure 2 can be observed also for the UNIMI  
263 data (Figure 1S). The scores plot of the first two Principal Components (Figure 2a) shows that PC1  
264 (82.61% of captured variance), is the main direction along which the samples separated according  
265 to the continent of origin. The greatest amount of American coffee samples assumed negative score  
266 values of PC1; thus, they were found in the left side of the scores plot. Conversely, almost all the  
267 Asian coffee samples were found to assume positive score values of PC1, thus they were gathered  
268 in the right side of the score plot. The observed score distribution is in accordance with Medina et  
269 al. (2017), who identified a clear separation between Arabica (American coffee beans) and Robusta  
270 (Indian coffee beans) species using NIR spectroscopy. However, in our study, some of the Asian  
271 coffee samples appear overlapped to the American ones, at least looking at the score plot of the first  
272 two components.

273 The distribution of the scores along PC2 (12.69% of explained variance) could be essentially  
274 ascribed to the potential intrinsic variability existing among coffee samples belonging to the same  
275 continent but to different countries. In particular, a great variability was found among the Asian  
276 samples, which even formed two well-defined clusters. The first one, which included the greatest  
277 part of the samples, is characterized by positive score values on PC2, while the second cluster is  
278 characterized by negative score values on the same component.

279 Figure 2b represents the score plot (PC1 vs. PC2) of the same PCA model presented in Figure 2a,  
280 but where each coffee sample was labelled according to the respective country of origin. As before,  
281 a separation trend can be found mainly along PC1. The Honduran samples assumed the lower score  
282 values, followed by Brazilian samples with score values approximately near to zero. The majority  
283 of Indian samples, as well as the Vietnamese samples, assumed instead positive score values, but no  
284 trend of separation between these two countries was observed along PC1, which was instead found  
285 along PC2. Moreover, some Indian samples appeared to be quite overlapped both to Honduran and  
286 Brazilian samples. This is in accordance with the findings of Bona et al., 2017, who tested within-  
287 country samples and found a clear grouping of coffee beans according to the considered Brazilian  
288 cities along PC1 (84% of variance) but a relevant group scattering along PC2 (12% of variance),  
289 probably related to the genetic variability and the different cultivation conditions.

290 The samples belonging to the class ‘Other’ were found to be spread over all the plot quarters  
291 without grouping in a defined cluster, neither along PC1 nor along PC2. This could be expected as  
292 they belong to different countries located in different continents. Indeed, the 7 samples located in  
293 the PC1 positive quarter have Indonesian origin, whereas the samples from Guatemala (11),  
294 Colombia (8), Costa Rica (2) and Nicaragua (2) assumed negative PC1 values except for one  
295 sample. Due to the lack of homogenous distribution of this class, its implementation in the further  
296 classification models could add errors more than benefits; thus, it was excluded.

297 To relate the samples distribution observed in the scores plot to spectral features, the loadings plot  
298 of PC1 and PC2 are represented in Figure 2c. It was found that at many spectral regions commented  
299 in Section 3.1 the corresponding loading values are different from zero. This means that these  
300 regions contribute in explaining the variability of the coffee samples on both the PCs considered.  
301 By looking at the loading values on PC1 it is possible to state that negative loadings values are  
302 associated with spectral regions where the NIR signal is higher for the American samples, indeed  
303 having negative score values on PC1; whereas the spectral regions with positive loading values  
304 denote a higher NIR signal for the Asian samples; actually they have positive score values on PC1.  
305 In detail, variables negatively influencing the samples distribution along PC1 are the ones  
306 corresponding to the maximum of broad peak at  $8450\text{-}8000\text{ cm}^{-1}$ , the beginning of the peak at  $7000\text{ cm}^{-1}$   
307 and the two small peaks at  $5780$  and  $5680\text{ cm}^{-1}$ . The variable positively influencing samples  
308 distribution along PC1 is the one corresponding to  $4020\text{ cm}^{-1}$ . Concerning PC2 loading values,  
309 spectral variables at  $5200$  and  $4020\text{ cm}^{-1}$  pull samples to negative values; whereas the change in  
310 spectral slope at  $7000\text{ cm}^{-1}$  and the peaks at  $6000$  and  $5400\text{ cm}^{-1}$  influenced samples distribution  
311 towards positive PC1 score values.

## 312 **Figure 2**

### 313 *3.3 PLS-DA classification models*

314  
315  
316 The continent-based classification models were built for the spectra collected in each laboratory  
317 after subjecting them to all the considered pre-treatment techniques.

318 The models showed a good classification capability, leading in any case to EFF values in prediction  
319 of the external test set ( $\text{EFF}_{\text{TS}}$ ) higher than 93.0% (Table 2). Moreover, the model robustness and  
320 reproducibility were proved by the high similarity of the results in calibration, cross-validation and  
321 validation of the external test set, no matter the pre-treatment applied. The model performances  
322 agree with those reported by Medina et al., 2017 for the PLS-DA classification of coffee beans by  
323 species (Arabica vs. Robusta), which led to 100% (model) accuracy using a 7-fold cross-validation.

324 However, they did not perform an external validation as it was done in the present work for making  
325 the results even more robust. Indeed, the results of the McNemar test demonstrated that all the  
326 models obtained, no matter the pre-treatment or the laboratory performing the analysis, were  
327 comparable in terms of classification error rate in prediction ( $P > 0.05$ , data not shown). Thus, the  
328 best PLS-DA classification results could be considered the ones obtained after SNV transformation  
329 for both laboratories, since they combine a soft mathematical data pre-treatment and a high EFF in  
330 prediction, giving results analogous to those calculated with any other pre-treated data.

### 331 **Table 2**

332

333 Regarding the country-based classification models (Table 3), the highest performance in cross-  
334 validation was obtained for the Brazilian samples, with 98.1-95.9% of  $EFF_{CV}$  for all the models  
335 calculated, except for those developed after second derivative transformation. Honduras, India and  
336 Vietnam classes were better predicted when soft pre-processing strategies (namely SNV and MSC)  
337 were applied, leading to an  $EFF_{CV}$  ranging from 97.8 to 90.0%. The results obtained from the  
338 prediction of the external test set confirmed what already observed for the internal validation of the  
339 model. The PLS-DA results reported by Marquetti et al., 2016 showed a better performance,  
340 reaching up to 100 % of sensitivity and specificity. Nevertheless, it should be mentioned that, in  
341 that case, the model was validated using 18 samples whereas, in the present study, 47 samples  
342 belonging to two deliveries, were considered as validation set.

### 343 **Table 3**

344

345 To further assess the models' robustness, the classification performance in prediction was evaluated  
346 by McNemar test giving an objective comparison of the results.

347 From the comparison of the results developed with the FT-NIR data collected at POLITO, it  
348 resulted that all the models gave comparable performances ( $P > 0.05$ ), except the ones obtained  
349 after second derivative transformation alone or combined with SNV and MSC (Table 4).

350 The best model developed with FT-NIR spectra collected at UNIMI was the one combining SNV  
351 and first derivative transformation, which led to  $EFF_{TS}$  of 98.5, 98.7, 93.5 and 83.7% for Brazil,  
352 Honduras, India and Vietnam, respectively (Table 3). However, the McNemar results demonstrated  
353 that this model has no significant differences ( $P > 0.05$ ) if compared with the performances in  
354 prediction of raw, SNV and MSC models always developed with the data collected in the UNIMI  
355 laboratory (Table 4).

356 Comparing POLITO and UNIMI results by McNemar test, it was observed that the models  
357 developed with no pre-treatment (raw data) or with SNV, SNV combined with first derivative, and  
358 MSC pre-processing gave comparable prediction performances ( $P > 0.05$ , Table 4).

#### 359 **Table 4**

### 361 *3.4 Variable selection*

362  
363 In order to improve the model performances and in vision of a simplified instrument with reduced  
364 spectral range, a variable selection strategy (iPLS-DA) was applied for the best models obtained  
365 with the whole spectral range.

366 Since a soft mathematical pre-treatment could be considered the better solution for practical  
367 implementation (Grassi et al., 2018a), the models developed with SNV transformed data could be  
368 considered the most convenient strategy to be further optimized by variable selection for both  
369 laboratories.

370 Therefore, the continent-based models developed with SNV data by both laboratories were  
371 subjected to iPLS-DA variable selection, as reported in Section 2.3.3. The  $EFF_{TS}$  of the iPLS-DA  
372 models was always 100% for America-class and 96.5% for Asia-class, no matter the data points  
373 interval size ( $i-80$ ,  $i-40$ ,  $i-20$  or  $i-10$ ) (data not shown). These results are in agreement with the  
374 variable reduction approach proposed by Calvinini et al., 2017, who obtained classification models  
375 with  $EFF_{TS}$  of 98.3% and 100.0% for Arabica and Robusta, respectively, by applying different filter  
376 simulations on NIR-hyperspectral data.

377 Moreover, the McNemar test was performed to compare the results obtained with the whole spectra  
378 range and after the variable selection. P-values higher than 0.05 confirmed the rejection of the null  
379 hypothesis, meaning that no significant difference exists between the models. Thus, a model  
380 developed with just 40 spectral variables out of 1154 (i.e. the whole spectral range) selected in four  
381 different spectral regions ( $12258-12188\text{ cm}^{-1}$ ,  $5855-5786\text{ cm}^{-1}$ ,  $5315-5246\text{ cm}^{-1}$  and  $4852-4783\text{ cm}^{-1}$ )  
382 <sup>1</sup>) was found to ensure an efficient continent-based discrimination, decreasing considerably the NIR  
383 instrument complexity at the same time.

384 In the case of country-based classification, the iPLS-DA models gave  $EFF_{TS}$  ranging from 100% to  
385 94.9%, depending on the class and the interval size considered (Table 5).

#### 386 **Table 5**

387  
388 The excellent performances of the country-based models were confirmed by McNemar results that  
389 reported no significant differences ( $P > 0.05$ ) among all the iPLS-DA models considered. Thus, in

390 this case, 90 spectral variables selected in three different spectral regions (9018-8871  $\text{cm}^{-1}$ , 8632-  
391 8177  $\text{cm}^{-1}$  and 6009-5940  $\text{cm}^{-1}$ ) guarantee discriminant performances comparable to the whole  
392 spectral range SNV-models. This is in accordance with the use of portable devices with spectral  
393 range reduced from 12500-4000  $\text{cm}^{-1}$  to 10500-6000  $\text{cm}^{-1}$  and only 125 variables, which  
394 demonstrated to give reliable results for both food authentication (Grassi et al, 2018a) and process  
395 monitoring (Grassi et al., 2018b).

396

### 397 *3.5 Cross-laboratory model validation*

398

399 The iPLS-DA continent-based classification model developed using the 40 abovementioned spectral  
400 variables was validated cross-laboratory. Both combinations, i.e. calibration set by POLITO and  
401 external set by UNIMI and vice versa, gave a weighted average  $\text{EFF}_{\text{TS}}$  of 100%

402 Concerning the country-based discrimination, the iPLS-DA model developed with the POLITO  
403 training set gave a weighted  $\text{EFF}_{\text{TS}}$  of 95.9% when validated on the UNIMI external test set, while  
404 the prediction of the POLITO external test set from the iPLS-DA model calibrated on the UNIMI  
405 training dataset led to a weighted  $\text{EFF}_{\text{TS}}$  equal to 94.6%. Therefore, in all cases, the prediction  
406 performances were comparable to those achieved by the iPLS-DA classification models calibrated  
407 and validated using the NIR spectra acquired by the same laboratory.

408

409

## 410 **4. Conclusions**

411

412 In the present study, the feasibility of implementing an automated system for the determination of  
413 the geographical origin of green coffee beans, based on NIR spectroscopy and multivariate data  
414 analysis, has been demonstrated.

415 Besides the automation, the main benefits related to this kind of analytical approach are the  
416 objectivity, the non-destructive nature and its rapidity, leading to a cost-effective improvement of  
417 the quality assurance of such a key worldwide food product. Actually, with appropriate  
418 industrialization and once the chemometric model has been properly calibrated, the time elapsed  
419 from the acquisition of NIR spectra on unknown samples and their subsequent classification would  
420 require just a few seconds. Therefore, this method could represent a concrete and effective answer  
421 to the need, claimed by coffee producers, industrial manufacturers, as well as by the Food Control  
422 Authority, of affordable, rapid and efficient technologies for the evaluation of food quality and  
423 authenticity (in this case applied to coffee beans).

424 Moreover, the variable selection results establish the groundwork for the development of a portable  
425 and cost-effective handheld NIR device, customised for origin discrimination of green coffee beans  
426 directly “in-field” to certify authenticity and counteract frauds.

427 Last but not least, the promising results achieved by the cross-laboratory model validation  
428 demonstrate the potential transferability of a NIR spectroscopy-based method among different  
429 production sites or industries, where the availability of more instruments and different operators is  
430 required to perform routine quality control analyses.

431

432

### 433 **References**

434

435 - Alessandrini, L., Romani, S., Pinnavaia, G., & Dalla Rosa, M. (2008). Near infrared  
436 spectroscopy: An analytical tool to predict coffee roasting degree. *Analytica Chimica Acta*,  
437 *625(1)*, 95-102.

438 - Alonso-Salces, R. M., Serra, F., Reniero, F., & Héberger, K. (2009). Botanical and  
439 Geographical Characterization of Green Coffee (*Coffea arabica* and *Coffea canephora*):  
440 Chemometric Evaluation of Phenolic and Methylxanthine Contents. *Journal of Agricultural*  
441 *and Food Chemistry*, *57(10)*, 4224-4235.

442 - Anderson, K. A., & Smith, B. W. (2002). Chemical Profiling to Differentiate Geographic  
443 Growing Origins of Coffee. *Journal of Agricultural and Food Chemistry*, *50(7)*, 2068-2075.

444 - Arana, V. A., Medina, J., Alarcon, R., Moreno, E., Heintz, L., Schäfer, H., & Wist, J.  
445 (2015). Coffee’s country of origin determined by NMR: The Colombian case. *Food*  
446 *Chemistry*, *175*, 500-506.

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

449 - Bertone, E., Venturello, A., Giraud, A., Pellegrino, G., & Geobaldo, F. (2016).  
450 Simultaneous determination by NIR spectroscopy of the roasting degree and  
451 Arabica/Robusta ratio in roasted and ground coffee. *Food Control*, *59*, 683-689.

452 - Bicchi, C. P., Panero, O. M., Pellegrino, G. M., & Vanni, A. C. (1997). Characterization of  
453 roasted coffee and coffee beverages by solid phase microextraction-gas chromatography and  
454 principal component analysis. *Journal of Agricultural and Food Chemistry*, *45(12)*, 4680-  
455 4686.

456 - Bro, R., & Smilde, A. K. (2014). Principal component analysis. *Analytical Methods*, *6*,  
457 2812-2831.

- 458 - Boix, A., Fernández Pierna, J. A., von Holst, C., & Baeten, V. (2012). Validation of a near  
459 infrared microscopy method for the detection of animal products in feedingstuffs: results of  
460 a collaborative study. *Food Additives & Contaminants: Part A*, *29*(12), 1872-1880.
- 461 - Bona, E., Marquetti, I., Link, J. V., Makimori, G. Y. F., da Costa Arca, V., Guimarães  
462 Lemes, A. L., Ferreira, J. M. G., dos Santos Scholz, M. B., Valderrama, P., & Poppi, R. J.  
463 (2017). Support vector machines in tandem with infrared spectroscopy for geographical  
464 classification of green arabica coffee. *LWT-Food Science and Technology*, *76*, 330-336.
- 465 - Buratti, S., Sinelli, N., Bertone, E., Venturello, A., Casiraghi, E., & Geobaldo, F. (2015).  
466 Discrimination between washed Arabica, natural Arabica and Robusta coffees by using near  
467 infrared spectroscopy, electronic nose and electronic tongue analysis. *Journal of the Science  
468 of Food and Agriculture*, *95*(11), 2192-2200.
- 469 - Cagliani, L. R., Pellegrino, G., Giugno, G., & Consonni, R. (2013). Quantification of *Coffea  
470 arabica* and *Coffea canephora* var. Robusta in roasted and ground coffee blends. *Talanta*,  
471 *106*, 169-173.
- 472 - Calvini, R., Amigo, J. M., & Ulrici, A. (2017). Transferring results from NIR-hyperspectral  
473 to NIR-multispectral imaging systems: a filter-based simulation applied to the classification  
474 of Arabica and Robusta green coffee. *Analytica Chimica Acta*, *967*, 33-41.
- 475 - Carrera, F., Leon-Camacho, M., Pablos, F., & Gonzalez, A. G. (1998). Authentication of  
476 green coffee varieties according to their sterolic profile. *Analytica Chimica Acta*, *370*, 131-  
477 139.
- 478 - Casal, S., Alves, M. R., Mendes, E., Oliveira, M. B. P. P., & Ferreira, M. A. (2003).  
479 Discrimination between Arabica and Robusta coffee species on the basis of their amino acid  
480 enantiomers. *Journal of Agricultural and Food Chemistry*, *51*(22), 6495-6501.
- 481 - El-Abassy, R. M., Donfack, P., & Materny, A. (2011). Discrimination between Arabica and  
482 Robusta green coffee using visible micro Raman spectroscopy and chemometric analysis.  
483 *Food Chemistry*, *126*(3), 1443-1448.
- 484 - Esteban-Díez, I., González-Sáiz, J. M., & Pizarro, C. (2004). Prediction of sensory  
485 properties of espresso from roasted coffee samples by near-infrared spectroscopy. *Analytica  
486 Chimica Acta*, *525*(2), 171-182.
- 487 - Esteban-Díez, I., González-Sáiz, J. M., Sáenz-González, C., & Pizarro, C. (2007). Coffee  
488 varietal differentiation based on near infrared spectroscopy. *Talanta*, *71*(1), 221-229.
- 489 - Fagerlan, M. W., Lydersen, S., & Laake, P. (2013). The McNemar test for binary matched-  
490 pairs data: mid-p and asymptotic are better than exact conditional. *BMC Medical Research  
491 Methodology*, *13*, 1-8.



- 492 - Frega, N., Bocci, F., & Lercker, G. (1995). Determination of Robusta in commercial coffee  
493 blends with Arabica coffee. *Industria Alimentari*, 34(339), 705-708.
- 494 - Gil-Agusti, M. T., Campostrini, N., Zolla, L., Ciambella, C., Invernizzi, C., & Righetti, P.  
495 G. (2005). Two-dimensional mapping as a tool for classification of green coffee bean  
496 species. *Proteomics*, 5(3), 710-718.
- 497 - González, A. G., Pablos, F., Martín, M. J., León-Camacho, M., & Valdenebro, M. S. (2001).  
498 HPLC analysis of tocopherols and triglycerides in coffee and their use as authentication  
499 parameters. *Food Chemistry*, 73(1), 93-101.
- 500 - Grassi, S., Casiraghi, E., & Alamprese, C. (2018a). Handheld NIR device: A non-targeted  
501 approach to assess authenticity of fish fillets and patties. *Food Chemistry*, 243, 382-388.
- 502 - Grassi, S., Cardone, G., Bigagnoli, D., & Marti, A. (2018b). Monitoring the sprouting  
503 process of wheat by non-conventional approaches. *Journal of Cereal Science*, 83, 180-187.
- 504 - Huck, C. W., Guggenbichler, W., & Bonn, G. K. (2005). Analysis of caffeine, theobromine  
505 and theophylline in coffee by near infrared spectroscopy (NIRS) compared to high-  
506 performance liquid chromatography (HPLC) coupled to mass spectrometry. *Analytica  
507 Chimica Acta*, 538, 195-203.
- 508 - Kemsley, E. K., Ruault, S., & Wilson, R. H. (1995). Discrimination between *Coffea arabica*  
509 and *Coffea canephora* variant Robusta beans using infrared-spectroscopy. *Food Chemistry*,  
510 54(3), 321-326.
- 511 - Leardi, R., & Nørgaard, L. (2004). Sequential application of backward interval partial least  
512 squares and genetic algorithms for the selection of relevant spectral regions. *Journal of  
513 Chemometrics*, 18(11), 486-497.
- 514 - Marquetti, I., Link, J. V., Guimarães Lemes, A. L., dos Santos Scholz, M. B., Valderrama,  
515 P., & Bona, E. (2016). Partial least square with discriminant analysis and near infrared  
516 spectroscopy for evaluation of geographic and genotypic origin of Arabica coffee.  
517 *Computers and Electronics in Agriculture*, 121, 313-319.
- 518 - Martín, M. J., Pablos, F., González, A. G., Valdenebro, M. S., & León-Camacho, M. (2001).  
519 Fatty acid profiles as discriminant parameters for coffee varieties differentiation. *Talanta*,  
520 54(2), 291-297.
- 521 - Medina, J., Caro Rodríguez, D., Arana, V. A., Bernal, A., Esseiva, P., & Wist, J. (2017).  
522 Comparison of Attenuated Total Reflectance Mid-Infrared, Near Infrared, and <sup>1</sup>H-Nuclear  
523 Magnetic Resonance Spectroscopies for the Determination of Coffee's Geographical Origin.  
524 *International Journal of Analytical Chemistry*. <https://doi.org/10.1155/2017/7210463>.

- 525 - Nørgaard, L., Saudland, A., Wagner, J., Nielsen, J. P., Munck, L., & Engelsen, S. B. (2000).  
526 Interval partial least-squares regression (iPLS): a comparative chemometric study with an  
527 example from near-infrared spectroscopy. *Applied Spectroscopy*, 54 (3), 413-419.
- 528 - Ozaki, Y., McClure W. F., & Christy, A. A. (2006). *Near-Infrared Spectroscopy in Food*  
529 *Science and Technology*. United States of America: Wiley & Sons, Inc.
- 530 - Patz, C. D., Blieke, A., Ristow, R., & Dietrich, H. (2004). Application of FT-MIR  
531 spectrometry in wine analysis. *Analytica Chimica Acta*, 513(1), 81-89.
- 532 - Pizarro, C., Esteban-Díez, I., González-Sáiz, J. M., & Forina, M. (2007). Use of near-  
533 infrared spectroscopy and feature selection techniques for predicting the caffeine content  
534 and roasting color in roasted coffees. *Journal of Agricultural and Food Chemistry*, 55 (18),  
535 7477-7488.
- 536 - Ribeiro, J. S., Ferreira, M. M. C., & Salva, T. J. G. (2011). Chemometric models for the  
537 quantitative descriptive sensory analysis of Arabica coffee beverages using near infrared  
538 spectroscopy. *Talanta*, 83(5), 1352-1358.
- 539 - Roggo, Y., Duponchel, L., & Huvenne, J. P. (2003). Comparison of supervised pattern  
540 recognition methods with McNemar's statistical test: Application to qualitative analysis of  
541 sugar beet by near-infrared spectroscopy. *Analytica Chimica Acta*, 477(2), 187-200.
- 542 - Rubayiza, A. B., & Meurens, M. (2005). Chemical discrimination of Arabica and Robusta  
543 coffees by Fourier transform Raman spectroscopy. *Journal of Agricultural and Food*  
544 *Chemistry*, 53(12), 4654-4659.
- 545 - Savorani, F., Rasmussen, M. A., Rinnan, Å., & Engelsen, S. B. (2013). Chapter 12 -  
546 Interval-based chemometric methods in NMR foodomics. *Data Handling in Science and*  
547 *Technology*, 28, 449-486.
- 548 - Szymańska, E., Saccenti, E., Smilde, A. K., & Westerhuis, J. A. (2012). Double-check:  
549 validation of diagnostic statistics for PLS-DA models in metabolomics studies.  
550 *Metabolomics*, 8, 3-16.
- 551 - Teuber, R. (2010). Geographical Indications of Origin as a Tool of Product Differentiation:  
552 The Case of Coffee. *Journal of International Food & Agribusiness Marketing*, 22 (3-4),  
553 277-298.
- 554 - Wall, M. E., Rechtsteiner, A., & Rocha, L. M. (2003). Singular value decomposition and  
555 principal component analysis. In: *A Practical Approach to Microarray Data Analysis* (pp.  
556 91-109). Boston: Springer.

- 557 - Wold, S., Ruhe, A., Wold, H., & Dunn, W. J. (1984). The Collinearity Problem in Linear  
558 Regression. The Partial Least Squares (PLS) Approach to Generalized Inverses. *Journal of*  
559 *Scientific and Statistical Computing*, 5 (3), 735-743.
- 560 - Weyer, L., & Workman Jr, J. (2007). *Practical guide to interpretive near-infrared*  
561 *spectroscopy*. Boca Raton: CRC press.
- 562

563 **Figure captions**

564 **Figure 1** - Average NIR spectrum of American coffee samples (blue line) and Asian coffee samples  
565 (red line): **(a)** from raw data and **(b)** after SNV pre-treatment. (For interpretation of the references  
566 to colour in this figure legend, the reader is referred to the web version of this article.)

567 **Figure 2** - PCA results of spectra collected by POLITO and pre-treated by SNV and mean  
568 centering: **(a)** Scores (PC1 vs. PC2) of coffee samples labelled according to the continent of origin;  
569 **(b)** scores (PC1 vs. PC2) of coffee samples labelled according to the country of origin; **(c)** loadings  
570 plot for PC1 and PC2. (For interpretation of the references to colour in this figure legend, the reader  
571 is referred to the web version of this article.)

572 **Figure 1S** - PCA results of spectra collected by UNIMI and pre-treated by SNV and mean  
573 centering: **(a)** Scores (PC1 vs. PC2) of coffee samples labelled according to the continent of origin;  
574 **(b)** scores (PC1 vs. PC2) of coffee samples labelled according to the country of origin; **(c)** loadings  
575 plot for PC1 and PC2. (For interpretation of the references to colour in this figure legend, the reader  
576 is referred to the web version of this article.)

577

## Highlights

- NIR Spectroscopy investigated to identify the geographical origin of green coffee
- Partial Least Square-Discriminant Analysis applied to build classification models
- Coffee origin for both continent and country was tested as discrimination parameter
- Variable selection applied to NIR spectra allowed improving the model performance
- Inter-laboratory comparison of the classification results was made by McNemar test

**Table 1.** Continents and countries of origin of green coffee samples.

<i>Continent</i>	<i>Country</i>	<i>Sample no.</i>
AMERICA	BRAZIL	38
	HONDURAS	27
	GUATEMALA	11
	COLOMBIA	8
	COSTA RICA	2
	NICARAGUA	2
ASIA	INDIA	65
	VIETNAM	31
	INDONESIA	7

**Table 2.** Results of continent-based Partial Least Squares-Discriminant Analysis (PLS-DA) on NIR spectra of green coffee samples after different mathematical pre-treatments: Efficiency percentage obtained in calibration (CAL), cross-validation (CV) and prediction of the external test set (TS).

		AMERICA			ASIA		
		<i>CAL</i>	<i>CV</i>	<i>TS</i>	<i>CAL</i>	<i>CV</i>	<i>TS</i>
<b>POLITO</b>	Raw	98.6	98.6	100.0	94.6	94.6	100.0
	SNV	98.6	98.6	100.0	93.2	93.2	93.1
	SNV + d1	98.6	98.6	100.0	93.2	93.2	96.6
	SNV + d2	98.6	98.6	100.0	95.9	93.2	96.6
	MSC	98.6	98.6	100.0	93.2	93.2	93.1
	MSC + d1	98.6	98.6	100.0	93.2	93.2	96.6
	MSC + d2	98.6	98.6	100.0	95.9	93.2	96.6
	d1	98.6	98.6	100.0	95.9	95.9	100.0
	d2	98.6	98.6	100.0	97.3	94.6	96.6
<b>UNIMI</b>	Raw	98.6	98.6	100.0	97.3	97.3	100.0
	SNV	98.6	98.6	100.0	93.2	91.9	96.5
	SNV + d1	98.6	98.6	100.0	94.6	93.2	93.1
	SNV + d2	97.2	97.2	100.0	94.6	93.2	96.5
	MSC	98.6	98.6	100.0	93.2	91.9	96.5
	MSC + d1	98.6	98.6	100.0	94.6	93.2	93.1
	MSC + d2	97.2	97.2	100.0	94.6	93.2	96.5
	d1	98.6	98.6	100.0	91.9	91.9	96.5
	d2	97.2	97.2	100.0	95.9	93.2	96.5

**Note.** POLITO: models developed from the NIR data collected at the Department of Applied Science and Technology (Polytechnic of Turin); UNIMI: models developed from the NIR data collected at the Department of Food, Environmental and Nutritional Sciences (University of Milan). Raw: data elaborated without any pre-treatment except mean centering; SNV: standard normal variate; MSC: multiplicative scatter correction; d1: first derivative; d2: second derivative.

**Table 3.** Results of country-based Partial Least Squares-Discriminant Analysis (PLS-DA) on NIR spectra of green coffee samples after different mathematical pre-treatments: Efficiency percentage obtained in calibration (CAL), cross-validation (CV) and prediction of the external test set (TS).

	BRAZIL			HONDURAS			INDIA			VIETNAM			
	CAL	CV	TS	CAL	CV	TS	CAL	CV	TS	CAL	CV	TS	
<b>POLITO</b>	Raw	98.1	97.6	100.0	97.3	92.0	100.0	94.1	92.3	98.1	95.1	92.1	94.9
	SNV	97.6	97.6	100.0	99.5	95.5	100.0	97.5	91.6	93.5	95.1	92.1	88.2
	SNV + d1	97.6	97.6	98.5	95.0	95.0	98.7	93.4	92.3	89.0	92.6	92.1	83.7
	SNV + d2	97.0	97.0	98.5	92.8	91.0	93.5	86.5	86.4	76.1	92.1	89.5	77.5
	MSC	97.6	97.6	100.0	99.5	95.5	100.0	97.5	92.3	93.5	95.1	94.6	88.2
	MSC + d1	97.6	97.6	98.5	95.0	95.0	98.7	93.4	92.3	89.0	92.6	92.1	83.7
	MSC + d2	97.0	97.0	98.5	92.8	91.0	93.5	86.5	86.4	76.1	92.1	89.5	77.5
	d1	97.6	97.0	100.0	95.0	92.6	98.7	93.4	90.9	89.0	92.6	86.8	82.5
	d2	95.9	90.4	97.0	91.0	85.6	94.9	84.0	82.1	76.1	92.6	89.5	77.5
	<b>UNIMI</b>	Raw	97.6	97.6	100.0	96.8	94.0	100.0	90.5	89.4	92.0	90.0	89.5
SNV		98.1	98.1	100.0	98.4	97.8	100.0	95.0	93.9	92.0	90.0	90.0	77.5
SNV + d1		98.1	96.5	98.5	98.4	90.4	98.7	95.0	93.9	93.5	90.0	87.3	83.7
SNV + d2		97.6	94.0	98.5	92.9	86.5	93.5	84.6	85.8	74.3	86.8	86.4	70.7
MSC		98.1	98.1	100.0	98.4	97.8	100.0	95.0	93.9	92.0	90.0	90.0	77.5
MSC + d1		97.6	95.9	98.5	92.9	83.6	93.5	85.8	85.8	74.3	87.3	86.8	70.7
MSC + d2		97.6	94.0	98.5	92.9	86.0	93.5	84.6	84.6	74.3	86.8	86.4	70.7
d1		97.0	96.5	98.5	90.7	86.5	93.5	85.1	86.2	74.3	84.5	84.1	70.7
d2		95.1	94.0	98.5	90.1	86.0	94.9	84.6	86.9	74.3	86.8	87.3	69.7

**Note.** POLITO: models developed from the NIR data collected at the Department of Applied Science and Technology (Polytechnic of Turin); UNIMI: models developed from the NIR data collected at the Department of Food, Environmental and Nutritional Sciences (University of Milan). Raw: data elaborated without any pre-treatment except mean centering; SNV: standard normal variate; MSC: multiplicative scatter correction; d1: first derivative; d2: second derivative.



**Table 4.** Results of McNemar test: P-value resulting from pair comparison of the country-based Partial Least Squares-Discriminant Analysis (PLS-DA) on NIR spectra of green coffee samples after different mathematical pre-treatments. In bold, the p-values <0.05 denote a significant difference between the considered models.

		POLITO								UNIMI									
		Raw	SNV	SNV + d1	SNV + d2	MSC	MSC + d1	MSC + d2	d1	d2	Raw	SNV	SNV + d1	SNV + d2	MSC	MSC + d1	MSC + d2	d1	d2
<b>POLITO</b>	Raw	1.000	0.250	0.063	0.002	0.250	0.063	<b>0.002</b>	0.063	<b>0.002</b>	0.125	0.125	0.125	<b>0.001</b>	0.125	<b>0.001</b>	<b>0.001</b>	<b>0.001</b>	<b>0.001</b>
	SNV	0.250	1.000	0.250	<b>0.021</b>	1.000	0.250	<b>0.021</b>	0.250	<b>0.021</b>	0.625	0.625	0.500	<b>0.012</b>	0.625	<b>0.012</b>	<b>0.012</b>	<b>0.012</b>	<b>0.012</b>
	SNV + d1	0.063	0.250	1.000	0.070	0.250	1.000	0.070	1.000	0.070	0.625	0.688	0.625	<b>0.039</b>	0.688	<b>0.039</b>	<b>0.039</b>	<b>0.039</b>	<b>0.039</b>
	SNV + d2	<b>0.002</b>	<b>0.021</b>	0.070	1.000	<b>0.021</b>	0.070	1.000	0.070	1.000	<b>0.039</b>	<b>0.039</b>	0.065	0.500	<b>0.039</b>	0.500	0.500	0.500	0.500
	MSC	0.250	1.000	0.250	<b>0.021</b>	1.000	0.250	<b>0.021</b>	0.250	<b>0.021</b>	0.625	0.625	0.500	<b>0.012</b>	0.625	<b>0.012</b>	<b>0.012</b>	<b>0.012</b>	<b>0.012</b>
	MSC + d1	0.063	0.250	1.000	0.070	0.250	1.000	0.070	1.000	0.070	0.625	0.688	0.625	<b>0.039</b>	0.688	<b>0.039</b>	<b>0.039</b>	<b>0.039</b>	<b>0.039</b>
	MSC + d2	<b>0.002</b>	0.021	0.070	1.000	0.021	0.070	1.000	0.070	1.000	<b>0.039</b>	<b>0.039</b>	0.065	0.500	<b>0.039</b>	0.500	0.500	0.500	0.500
	d1	0.063	0.250	1.000	0.070	0.250	1.000	0.070	1.000	0.070	0.625	0.688	0.625	<b>0.039</b>	0.688	<b>0.039</b>	<b>0.039</b>	<b>0.039</b>	<b>0.039</b>
	d2	<b>0.002</b>	<b>0.021</b>	0.070	1.000	<b>0.021</b>	0.070	1.000	0.070	1.000	<b>0.039</b>	<b>0.039</b>	0.065	0.500	<b>0.039</b>	0.500	0.500	0.500	0.500
	<b>UNIMI</b>	Raw	0.125	0.625	0.625	<b>0.039</b>	0.625	0.625	0.039	0.625	<b>0.039</b>	1.000	1.000	1.000	<b>0.008</b>	1.000	<b>0.008</b>	<b>0.008</b>	<b>0.008</b>
SNV		0.125	0.625	0.688	<b>0.039</b>	0.625	0.688	0.039	0.688	<b>0.039</b>	1.000	1.000	1.000	<b>0.008</b>	1.000	<b>0.008</b>	<b>0.008</b>	<b>0.008</b>	<b>0.008</b>
SNV + d1		0.125	0.500	0.625	0.065	0.500	0.625	0.065	0.625	0.065	1.000	1.000	1.000	<b>0.021</b>	1.000	<b>0.021</b>	<b>0.021</b>	<b>0.021</b>	<b>0.021</b>
SNV + d2		<b>0.001</b>	<b>0.012</b>	<b>0.039</b>	0.500	<b>0.012</b>	<b>0.039</b>	0.500	<b>0.039</b>	0.500	<b>0.008</b>	<b>0.008</b>	<b>0.021</b>	1.000	<b>0.008</b>	1.000	1.000	1.000	1.000
MSC		0.125	0.625	0.688	0.039	0.625	0.688	0.039	0.688	<b>0.039</b>	1.000	1.000	1.000	0.008	1.000	<b>0.008</b>	<b>0.008</b>	<b>0.008</b>	<b>0.008</b>
MSC + d1		<b>0.001</b>	<b>0.012</b>	<b>0.039</b>	0.500	<b>0.012</b>	<b>0.039</b>	0.500	<b>0.039</b>	0.500	<b>0.008</b>	<b>0.008</b>	<b>0.021</b>	1.000	<b>0.008</b>	1.000	1.000	1.000	1.000
MSC + d2		<b>0.001</b>	<b>0.012</b>	<b>0.039</b>	0.500	<b>0.012</b>	<b>0.039</b>	0.500	<b>0.039</b>	0.500	<b>0.008</b>	<b>0.008</b>	<b>0.021</b>	1.000	<b>0.008</b>	1.000	1.000	1.000	1.000
d1		<b>0.001</b>	<b>0.012</b>	<b>0.039</b>	0.500	<b>0.012</b>	<b>0.039</b>	0.500	<b>0.039</b>	0.500	<b>0.008</b>	<b>0.008</b>	<b>0.021</b>	1.000	<b>0.008</b>	1.000	1.000	1.000	1.000
d2		<b>0.001</b>	<b>0.012</b>	<b>0.039</b>	0.500	<b>0.012</b>	<b>0.039</b>	0.500	<b>0.039</b>	0.500	<b>0.008</b>	<b>0.008</b>	<b>0.021</b>	1.000	<b>0.008</b>	1.000	1.000	1.000	1.000

**Note.** POLITO: models developed from the NIR data collected in the Department of Applied Science and Technology (Polytechnic of Turin); UNIMI: models developed from the NIR data collected in the Department of Food, Environmental and Nutritional Sciences (University of Milan); Raw: data elaborated without any pre-treatment except mean centering; SNV: standard normal variate; MSC: multiplicative scatter correction; d1: first derivative; d2: second derivative.

**Table 5.** Results of iPLS-DA variable selection strategy applied on the country-based Partial Least Squares-Discriminant Analysis (PLS-DA) on the SNV-NIR spectra of green coffee samples: Efficiency percentage obtained in calibration (CAL), cross-validation (CV) and prediction of the external test set (TS).

		BRAZIL			HONDURAS			INDIA			VIETNAM		
		<i>CAL</i>	<i>CV</i>	<i>TS</i>	<i>CAL</i>	<i>CV</i>	<i>TS</i>	<i>CAL</i>	<i>CV</i>	<i>TS</i>	<i>CAL</i>	<i>CV</i>	<i>TS</i>
<b>POLITO</b>	<i>i</i> -80	98.1	98.1	100.0	99.5	98.9	100.0	100.0	97.5	100.0	100.0	95.1	100.0
	<i>i</i> -40	97.6	97.6	100.0	99.5	98.9	100.0	98.9	96.4	100.0	100.0	95.1	100.0
	<i>i</i> -20	98.1	97.6	100.0	99.5	98.4	100.0	100.0	96.0	98.1	100.0	97.6	94.9
	<i>i</i> -10	98.1	98.1	100.0	98.4	97.8	100.0	97.8	95.3	100.0	100.0	95.1	100.0
<b>UNIMI</b>	<i>i</i> -80	98.1	98.1	100.0	97.8	97.8	100.0	94.9	94.2	97.3	97.1	94.6	98.6
	<i>i</i> -40	98.1	97.6	100.0	98.9	96.1	100.0	98.2	96.4	98.1	97.6	95.1	94.9
	<i>i</i> -20	97.6	97.0	100.0	98.9	98.9	100.0	96.7	94.2	100.0	99.5	94.6	100.0
	<i>i</i> -10	97.6	97.6	100.0	98.9	98.9	100.0	96.7	95.3	100.0	99.5	94.6	100.0

**Note.** POLITO: models developed from the NIR data collected in the Department of Applied Science and Technology (Polytechnic of Turin); UNIMI: models developed from the NIR data collected in the Department of Food, Environmental and Nutritional Sciences (University of Milan).

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

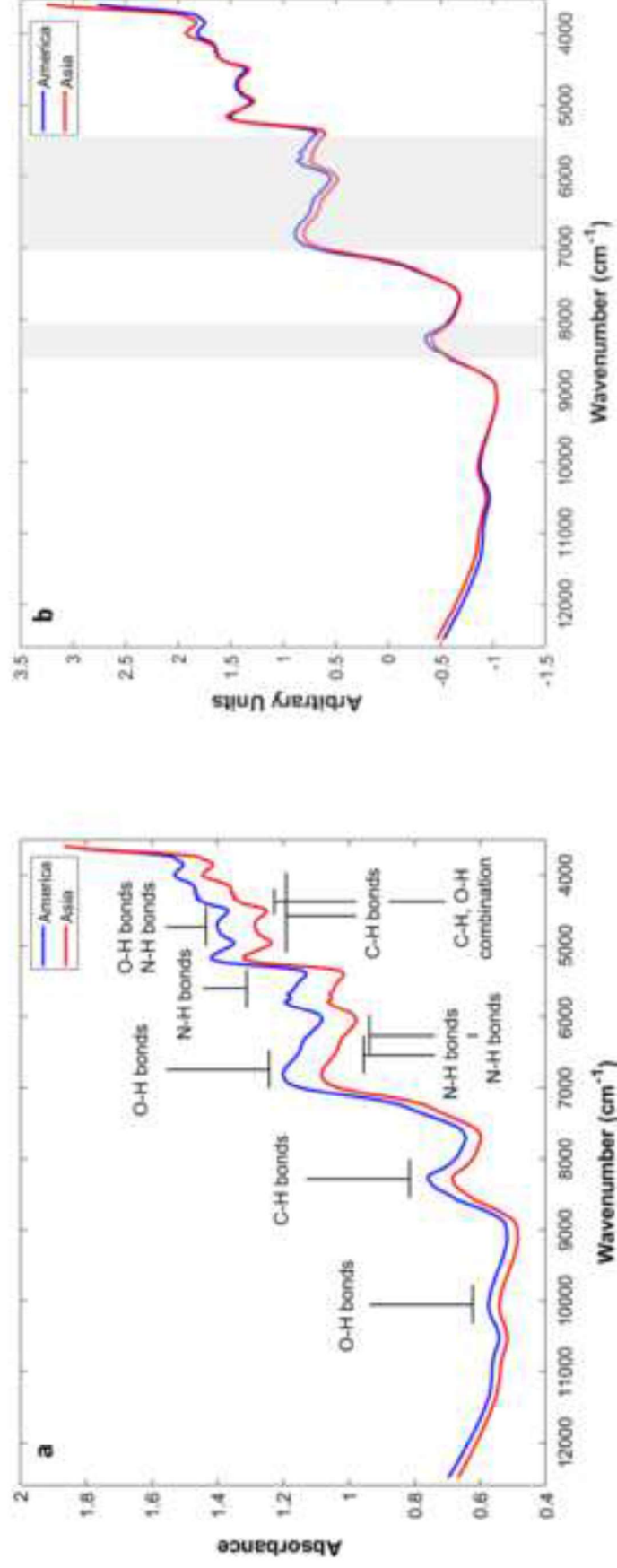


Figure 2  
[Click here to download high resolution image](#)

