

Accepted Manuscript

Classification of cowpea beans using multielemental fingerprinting combined with supervised learning

Michael Pérez-Rodríguez, José E. Gaiad, Melisa J. Hidalgo, María V. Avanza, Roberto G. Pellerano



PII: S0956-7135(18)30395-5
DOI: 10.1016/j.foodcont.2018.08.001
Reference: JFCO 6265
To appear in: *Food Control*
Received Date: 18 June 2018
Accepted Date: 01 August 2018

Please cite this article as: Michael Pérez-Rodríguez, José E. Gaiad, Melisa J. Hidalgo, María V. Avanza, Roberto G. Pellerano, Classification of cowpea beans using multielemental fingerprinting combined with supervised learning, *Food Control* (2018), doi: 10.1016/j.foodcont.2018.08.001

This is a PDF file of an unedited manuscript that has been accepted for publication. As a service to our customers we are providing this early version of the manuscript. The manuscript will undergo copyediting, typesetting, and review of the resulting proof before it is published in its final form. Please note that during the production process errors may be discovered which could affect the content, and all legal disclaimers that apply to the journal pertain.

1 **Classification of cowpea beans using multielemental fingerprinting**
2 **combined with supervised learning**

3

4 Michael Pérez-Rodríguez*, José E. Gaiad, Melisa J. Hidalgo, María V. Avanza, Roberto G.
5 Pellerano

6

7 Institute of Basic and Applied Chemistry of the Northeast of Argentina (IQUIBA-NEA),
8 National Scientific and Technical Research Council (CONICET), Faculty of Natural and Exact
9 Sciences and Surveying, National University of the Northeast – UNNE, Av. Libertad 5470,
10 Corrientes 3400, Argentina.

11

12

13

14

15

16

17

18

19

20

21

22

23 * Corresponding author. Phone: +54 379 445 7996

24 E-mail address: michaelpr1984@gmail.com (M. Pérez-Rodríguez).

25

26 **ABSTRACT**

27
28 Multielemental compositions (Ag, As, Ba, Be, Cd, Cs, Co, Cr, Cu, Mo, Ni, Pb, Sb, Se, Sn, Sr,
29 Tl, Rb, V, and Zn) of 106 cowpea bean samples belonging to different varieties collected from
30 the province of Corrientes in Argentina were determined using inductively coupled plasma mass
31 spectrometry (ICP-MS). Based on the multielemental data, five supervised learning techniques,
32 namely, linear discriminant analysis (LDA), partial least square discriminant analysis (PLS-DA),
33 k nearest neighbors (k-NN), random forest (RF), and support vector machine (SVM) with radial
34 basis function Kernel, were computed aiming at building classification models that allow one to
35 predict the botanical variety of the samples based on their element profiles. The best
36 classification performance was obtained by SVM with 93% accuracy rate. The model developed
37 through this method enabled the correct separation of the samples into the five cowpea varieties
38 investigated, where 100% sensitivity was achieved for most of the predicted classes. Thus,
39 SVM was the algorithm selected for the classification of the cowpea beans according to their
40 botanical variety. Multielemental determination coupled with supervised pattern recognition
41 techniques have proved to be an interesting approach for differentiating a diverse range of
42 cowpea genotypes. This study has contributed toward generalizing the use of multielemental
43 fingerprinting as a promising tool for testing the authenticity of cowpea beans on a global scale.

44
45
46 **Keywords:** Cowpea bean; multielemental fingerprinting; genotype; supervised learning; ICP-
47 MS; authenticity.

48

49

50 1. Introduction

51

52 Grain legumes are an important source of essential nutrients in several developing
53 countries owing to nutritional and socioeconomic reasons. These edible plants are regarded as
54 the main protein source for populations in Asia, Africa and Latin America (Kato, Fernandes,
55 Bacchi, Sarriés, & Reyes, 2015; Santos, Santos, Fernandes, Castro, & Korn, 2013). Within
56 the Leguminosae family, cowpea beans (*Vigna unguiculata* L. Walp.) contain good contents
57 of vitamins, fibers, minerals, and other nutrients that are vital for the normal metabolic
58 functioning of the human body (Alayande, Mustapha, Dabak, & Ubom, 2012). In addition,
59 phenolic compounds, which are known for their use in the treatment and preventive action
60 against various pathologies, such as diabetes (Asgar, 2013), arterial hypertension (Souza,
61 Marcadenti, & Portal, 2017), and cardiovascular diseases (Scolaro, Kim, & Castro, 2018), are
62 found to be present in the composition of cowpea beans (Moreira-Araújo et al., 2018).

63 The presence of multielemental contents in cowpea products may stem from several
64 factors. The contents might come from the soil, environment, botanical variety (genotype) or
65 may be introduced during the production process (crop cultivation, transport, storage, or
66 preparation) (Kato et al., 2015; Santos, Gramacho, & Teixeira, 2008). Trace elements
67 quantification in legumes may also reveal the presence of inorganic contaminants which are
68 dangerous to human health (Zhou et al., 2016). In view of that, several international
69 organizations have set tolerance levels for toxic metals in leguminous seeds (MERCOSUR
70 2012; Codex Alimentarius, 2011; European Community, 2006). Considering the presence of
71 such harmful substances in grain legumes, performing an exhaustive quality control is
72 essentially important to ensure human food safety.

73 A huge number of authors have published studies in the literature regarding the
74 nutritional composition of different cowpea genotypes grown in South Africa (Belane &

75 Dakora, 2011), West Africa (Alayande, Mustapha, Dabak, & Ubom, 2012; Inobeme,
76 Nlemadim, Obigwa, Ikechukwu, & Ajai, 2014), Argentina (Avanza, Acevedo, Chaves, &
77 Añón, 2013), and Brazil (Santos et al., 2008, 2013, 2009, 2008). Predominantly,
78 multielemental analysis techniques with high sensitivity and ability to determine the isotope
79 composition of a sample were used in these studies. In this context, inductively coupled
80 plasma mass spectrometry (ICP-MS) has proved to be a rapid and high-accuracy technique
81 for simultaneous determination of several elements. The outstanding properties of this
82 technique are reflected in its conceivably superior performance when it comes to conducting
83 reliable detection and quantification of trace elements in a wide variety of samples (Ammann,
84 2007). Furthermore, the information obtained about the multielemental composition of
85 samples can be analyzed by pattern recognition techniques that provide one with powerful
86 tools for identifying geographical origin, quality grade, or genetic characteristics of foods
87 (Liu, Xue, Wang, Li, Xue, & Xu, 2012).

88 The present work focuses on characterizing cowpeas using ICP-MS based multielemental
89 fingerprinting (Ag, As, Ba, Be, Cd, Cs, Co, Cr, Cu, Mo, Ni, Pb, Sb, Se, Sn, Sr, Tl, Rb, V, and
90 Zn) aiming at identifying differences among the cultivated varieties and check the
91 authenticity of the beans through the application of machine learning techniques. To this end,
92 the study applied some pattern recognition techniques, such as, principal component analysis
93 (PCA), linear discriminant analysis (LDA), partial least squares discriminant analysis (PLS-
94 DA), k nearest neighbors (k-NN), random forest (RF), and support vector machines (SVM),
95 on multielemental data measured by ICP-MS for bean seeds produced in the northeastern
96 region of Argentina (Corrientes Province). The combined benefits of analytical and
97 chemometric techniques studied in this work were tested as a strategic tool for the
98 differentiation of a wide range of cowpea genotypes, which exhibited the potential to meet
99 the quality regulatory requirements for this food.

100

101 **2. Materials and methods**

102

103 *2.1. Instrumentation*

104

105 A Spex 6750 (Metuchen, NJ, USA) cryogenic mill was used to reduce the particle sizes of
106 the bean samples, while decomposition was performed using a microwave digestion system,
107 Ethos One (Milestone, Chicago, USA) equipped with programmable power control (1600 W
108 maximum power) and HPR 1000/10s segmented rotor (operating conditions of 35 bar
109 maximum pressure and 260 °C maximum temperature) with 10 reaction vessels.

110 Multielemental fingerprinting was determined using an Agilent 7700 Series (Agilent
111 Technologies, Japan) inductively coupled plasma mass spectrometer. The instrument was
112 equipped with a cooled double-pass quartz spray and a MicroMist glass concentric nebulizer.
113 The Fassel-type ICP torch was constituted by three-cylinder assembly with injector diameter
114 of 2.5 mm. Ni sampler and skimmer cones of 1.0 mm and 0.4 mm were used. The ICP-MS
115 operating conditions are summarized in Table 1. Internal standards were used in all
116 determinations in order to correct interferences stemming from the sample matrices. In order
117 of mass number, the selected isotopes for measurement included the following: ^9Be , ^{51}V ,
118 ^{53}Cr , ^{59}Co , ^{60}Ni , ^{63}Cu , ^{66}Zn , ^{75}As , ^{78}Se , ^{85}Rb , ^{88}Sr , ^{98}Mo , ^{107}Ag , ^{111}Cd , ^{118}Sn , ^{121}Sb , ^{133}Cs ,
119 ^{137}Ba , ^{205}Tl , and ^{208}Pb .

120

121 Insert Table 1

122

123 *2.2. Chemicals and standard solutions*

124

125 All chemical reagents employed were of ultrapure analytical grade. Nitric acid (65%
126 m/m) and hydrogen peroxide (30% m/m) were purchased from Merck (Darmstadt, Germany).
127 Nitric acid was further purified by sub-boiling distillation for later use. All standard and
128 working solutions were prepared in deionized water (18.0 M Ω cm at 25 °C) obtained from a
129 Milli-Q Plus Water purification system (Millipore Corp., Molsheim, France). Multielemental
130 calibration solutions were prepared from multi-element standard solution, TraceCERT®
131 CRM, purchased from Sigma-Aldrich (St. Louis, MO, USA). For assessing the accuracy of
132 analytical method, a plant based SRM tomato leaves (NIST® SRM 1573a) purchased from
133 Sigma-Aldrich (St. Louis, MO, USA) was used.

134 Materials, including plastic containers, polyethylene flasks, pipette tips, and PFA Teflon
135 digestion vessels, were constantly checked for contaminations by blank quality control tests,
136 and glassware were strictly avoided.

137

138 *2.3. Bean samples*

139

140 A total of 106 cowpea bean samples were analyzed so as to determine their mineral
141 contents. Different agricultural cooperatives from the province of Corrientes (Argentina)
142 provided us with cowpeas belonging to five botanical varieties, namely: Alarcon (ALA),
143 California Black (CBK), Carnevalia Ensiformis (CES), Moro Riach (MRH), and San Miguel
144 (SML). The samples were obtained between September and October 2016–2017. After their
145 collection, the samples were immediately stored in zipped bags in a vacuum desiccator
146 cabinet until analysis.

147

148 *2.4. Sample preparation*

149

150 The bean samples were homogenized in a cryogenic mill, and approximately 500 mg of
151 dry samples were placed in closed vessels in a microwave oven, where 5 mL HNO₃ and 2 mL
152 H₂O₂, both in concentrated solution, were added later. The procedure was completed by using
153 the following temperature program: first stage: 25–200 °C for 10 min; second stage: 200 °C
154 for 15 min; and third stage: 200–110 °C for 10 min; followed immediately by ventilation at
155 room temperature for 10 min. Finally, the digested samples were diluted to 10 mL with
156 deionized water and prepared for analysis by ICP-MS. Blank solutions were prepared based
157 on the same procedure applied for the samples. To correct any instrument signal drift during
158 analysis and suppress mistakes in the analyte quantification, an internal standard solution
159 (Sigma-Aldrich, St. Louis, MO, USA) containing ⁷²Ge and ¹¹⁵In at 100 µg L⁻¹ was added to
160 all samples prior to the digestion stage. All measurements were performed in triplicate.

161

162 *2.5. Analytical performance parameters*

163

164 The analytical features of the proposed ICP-MS method were measured through
165 performance parameters, such as coefficient of determination, limit of quantification (LOQ),
166 accuracy, precision and overall recovery.

167 To construct the multi-element calibration curves, five different concentration levels in
168 triplicate were used, and the calibration ranges were modified according to the expected
169 mineral concentration ranges. Linear regression analysis by the least squares method was
170 used to calculate the coefficients of determination (r^2). The limits of quantification (LOQ)
171 were determined as ten times the standard deviation of measurements of 10 blank solutions,
172 divided by the slope of the calibration curve, according to the IUPAC recommendations
173 (Thompson, Ellison, & Wood, 2002).

174 Precision assays were conducted under conditions of repeatability and intermediate
175 precision over a period of three weeks, as such the samples under study were digested and
176 analyzed in triplicate in different days. The variability of the measurement was expressed as
177 the relative standard deviation (%RSD).

178 The accuracy of the method was assessed using replicate analysis of 5 different samples of
179 standard reference materials (SRM 1573a). In addition, recovery studies were carried out by
180 fortifying randomly selected bean samples with multi-element standard solutions at
181 concentrations of 10 and 100 $\mu\text{g kg}^{-1}$. The solutions, including the blank solutions, the
182 fortified solutions, and the certified material solutions, were all prepared following the same
183 procedure as the samples. **One limitation of this study is that the contents of some elements**
184 **such as, Ag, Ba, Be, Cs, Mo, Pb, Sn, Sr and Tl were not certified in the NIST SRM 1573a.**
185 **Therefore, only fortification results for these elements are available.**

186

187 *2.6. Multivariate data analysis*

188

189 The data matrix consisted of 106 rows, which corresponded to different cowpea bean
190 samples, and 20 columns containing the mineral concentration values (Ag, As, Ba, Be, Cd,
191 Cs, Co, Cr, Cu, Mo, Ni, Pb, Sb, Se, Sn, Sr, Tl, Rb, V, and Zn). Multivariate data processing
192 was carried out via machine learning techniques using caret package in R-project software
193 version 3.3.3 (R Core Team, 2017).

194 Principal component analysis (PCA) was used for exploratory data analysis aiming at
195 visualizing the natural distribution of samples in a reduced dimensional space. The dimension
196 reduction in PCA allows one to accurately represent high-dimensional data in lower
197 dimensional space. It is worth noting that this unsupervised learning method also verifies the
198 existence of relationships between the variables in a multidimensional space (Bro & Smilde,

199 2014; Varmuza & Filzmoser, 2009). Unsupervised learning focuses on describing the
200 associations and patterns among a set of input measures, revealing properties of the data
201 density (Hastie, Tibshirani, & Friedman, 2008).

202 After that, five supervised learning techniques were applied to the dataset for the
203 predictive modeling of botanical varieties of cowpea beans. To achieve this goal, two linear
204 methods, including linear discriminant analysis (LDA) and partial least squares discriminant
205 analysis (PLS-DA), and three non-linear models, namely, k nearest neighbors (k-NN),
206 support vector machine (SVM) and random forest (RF) were compared according to their
207 classification performance. Supervised learning focuses on predicting the value of an
208 outcome measure based on a number of input measures. The predictions are based on the
209 training samples of previously solved cases (samples), where the joint values of all of the
210 variables are known. Thus, a model is trained to generate reasonable classifications for new
211 data (Hastie et al., 2008).

212 During the training step, the parameters needed to build optimal classifiers were optimized
213 by ten-fold cross-validation so as to avoid the occurrence of bias. For this purpose, the data
214 matrix was randomly split into 10 mutually exclusive subsets, and each classifier was trained
215 and tested 10 times, so that the cross-validation estimate of metrics is over-tested 10 times.
216 This procedure is repeated five times, so each subset is used to test the model developed at
217 least one time. The parameters that were subjected to optimization included the number of
218 significant components (ncomp) for PLS-DA, number of k neighbors for k-NN, number of
219 trees (nt) and number of variables tried at each split (mtry) for RF, in addition to penalty
220 factor (C), γ of the γ -insensitive loss function and kernel type for SVM.

221 After selecting the optimal values for each model, some basic measures derived from the
222 confusion matrix were considered in order to evaluate the classification obtained using the
223 supervised algorithms under investigation (LDA, PLS-DA, k-NN, RF and SVM). The metrics

224 calculated per class included the following: sensitivity (correct positive predictions divided
225 by the number of positive cases), specificity (correct negative predictions divided by the
226 number of negative cases), and overall accuracy (all correct predictions divided by the total
227 number of cases examined) (Lantz, 2015).

228 LDA classification model is used for differentiating classes of samples by minimizing the
229 variance within classes and maximizing the variance among classes. This method is based on
230 the estimation of several canonical or discriminant functions, which are linear combinations
231 of the original variables in order to optimize the separation (Moncayo, Manzoor, & Caceres,
232 2015).

233 PLS-DA is a linear classification algorithm based on partial least squares regression used
234 for conducting predictive modeling. This technique allows class distinction, which enables
235 one to find latent variables from dependent variables with maximum covariance (Brereton &
236 Lloyd, 2014).

237 K-NN is a distance-based non-parametric discriminant technique. This method uses the
238 distance between objects to assign one of them to the most common class among the k-
239 nearest neighbors. The optimal size of neighbor k must be optimized (Bevilacqua et al.,
240 2013).

241 RF is a supervised learning method based on a set of tree decision predictors for
242 classification, which was created as an extension of bagging. In this method, a random vector
243 is generated independently on the input vector, and each tree casts a vote meant for
244 classifying an input vector (Archer & Kimes, 2008; Hernández-Pereira, Álvarez-Estévez, &
245 Moret-Bonillo, 2015).

246 SVM is a machine learning technique that performs a mapping of the training data in a
247 high-dimensional space directed toward building a classifier in that space (Batista et al.,

248 2012). This method classifies the data by constructing a separate hyperplane in n-dimensional
249 space, which maximizes the margin between classes (Bona et al., 2017).

250 Aside from the calculated metrics, Kappa coefficient of agreement was likewise
251 considered to enable us to compare the performance of the proposed classification models,
252 since the cowpea bean dataset is unbalanced in terms of the sample number per group
253 (variety). This robust statistic is commonly used either for evaluating inter-classifier or intra-
254 classifier reliability (McHugh, 2012). In addition, Kappa statistical measure is a form of
255 correlation coefficient, so that squaring the correlation value facilitates its interpretation.
256 Squared Kappa is referred to as coefficient of determination; it is defined as the amount of
257 variation in the dependent variable that can be explained by the independent variable
258 (Stephens & Diesing, 2014; Tang, Hu, Zhang, Wu, & He, 2015). A detailed description
259 regarding the calculation of the kappa statistic can be found in the work published by
260 McHugh, 2012 (McHugh, 2012).

261

262 **3. Results and discussion**

263

264 *3.1. Analytical features of the ICP-MS method under investigation*

265

266 Table 2 summarizes the analytical features of the method studied in this work. The
267 calibration curves obtained exhibited good linearity in the selected concentration range for
268 each element, with coefficients of determination (r^2) greater than 0.9985. The LOQs ranged
269 from 2 to 545 $\mu\text{g kg}^{-1}$. These values indicated that the method proposed in this work was
270 sufficiently sensitive for determining trace elements in cowpea beans.

271

272 Insert Table 2

273

274 The intra-day and inter-day precisions for the elements measured were within the ranges
275 of 0.2–4.8% and 1.5–6.9%, respectively. Accuracy assays performed by analyzing SRM
276 1573a showed good agreement between the results obtained by the present method and the
277 certified values. The average recoveries of the elements ranged from 93% to 107%, with RSD
278 values being less than 12% (Table 2). Furthermore, overall recoveries in fortified samples
279 were in the range of 86.2–109.2%. According to the Codex Alimentarius recommendations
280 (Codex Alimentarius, 2009), these results are within the acceptable criteria for recovery (70–
281 110 %) and intermediate precision ($\leq 15\%$ RSD). In effect, the results clearly demonstrate that
282 the performance of the method proposed here was satisfactory and that the microwave
283 digestion step was essentially effective in contributing toward obtaining reliable results with
284 adequate accuracy and precision.

285

286 3.2. Multielemental fingerprinting of cowpea bean samples

287

288 The applicability of the present ICP-MS method was checked through the determination of
289 multi-elements in *Vigna unguiculata* L. Walp. samples. Table 3 shows the element
290 concentrations obtained for the different cowpea bean varieties under study. The results are
291 expressed as average values of three measurements with the corresponding standard deviation
292 (SD).

293

294 Insert Table 3

295

296 Based on the results obtained, Ba and Zn were found to be the most abundant elements in
297 the cowpea beans investigated, with values ranging from 19.80–96.60 and 17.34–60.27 mg

298 kg^{-1} , respectively. Ba concentrations were higher than the reported values for cowpea beans
299 ($0.2\text{--}12 \text{ mg kg}^{-1}$) (Santos et al., 2013). Zn presented values similar to those found in different
300 cowpea genotypes ($44\text{--}65 \text{ mg kg}^{-1}$) (Belane & Dakora, 2011).

301 In terms of abundance, the elements that came next included Cu, Sr, and Mo. Cu content
302 ranged from 2.60 to 7.90 mg kg^{-1} , whereas Sr exhibited concentrations ranging from 1.50 to
303 8.12 mg kg^{-1} . These elements showed relatively greater variability compared to those found
304 in other cowpea varieties ($3.1\text{--}5.8 \text{ mg kg}^{-1}$ for Cu and $3.3\text{--}4.8 \text{ mg kg}^{-1}$ for Sr) (Santos et al.,
305 2013). With regard to Mo, MRH presented the highest average content followed by CBK,
306 SML, and CES; while ALA exhibited an average value that was farthest from the rest of the
307 varieties. In general, Mo contents in the beans investigated here were relatively lower than
308 those observed in common beans (*Phaseolus vulgaris* species) by other authors ($1.3\text{--}5.4 \text{ mg}$
309 kg^{-1}) (Kato et al., 2015).

310 The third group that came next in terms of abundance consisted of Ni, Pb, Sn, Be, Rb, Cr,
311 V, and Tl, in decreasing concentration order, with average values ranging from $0.09\text{--}0.67 \text{ mg}$
312 kg^{-1} . The highest Ni concentration was found in the SML variety, whereas the lowest was
313 obtained in the ALA variety. The concentrations of this element were lower compared to
314 those reported by Santos et al. in common beans ($0.75\text{--}6.7 \text{ mg kg}^{-1}$) and cowpeas ($2.9\text{--}3.4$
315 mg kg^{-1}) (Santos et al., 2013). Pb concentrations ranged from 0.14 to 0.52 mg kg^{-1} for the
316 cowpea varieties investigated here, with the highest values observed in beans from CBK and
317 lowest values in beans from MRH. The overall Pb concentration of 0.28 mg kg^{-1} obtained
318 was slightly above the tolerance limit in dried legumes, which is 0.20 mg kg^{-1} , set by the
319 MERCOSUR (MERCOSUR, 2012), FAO/WHO (Codex Alimentarius, 2011), and European
320 Commission (European Community, 2006). The concentrations of Sn and Be were found to
321 be lowest in beans from the ALA variety, while being highest for the MRH variety; the
322 average concentration of Sn (0.26 mg kg^{-1}) was found to be well below the allowable

323 maximum limit in canned foods according to the regulations of MERCOSUR (250 mg kg⁻¹)
324 (MERCOSUR, 2012) and the European Community (200 mg kg⁻¹) (European Community,
325 2006). The overall concentration values found for Cr and V ranged between 0.12–0.21 and
326 0.11–0.15 mg kg⁻¹, respectively. The concentration of Rb ranged from 0.05 to 0.45 mg kg⁻¹,
327 being lower than those determined in different common bean varieties collected from
328 Hungary (2.6–18.2 mg kg⁻¹) (Kato et al., 2015). With regard to Tl, similar concentration
329 profiles were observed between samples from different varieties of cowpeas investigated, and
330 their mean contents were lower than those found in vegetables (0.02–0.3 mg kg⁻¹)
331 (Karbowska, 2016).

332 The trace elements at concentrations lower than 0.1 mg kg⁻¹ can be arranged in the
333 following order: Se > Sb > Co > Ag > As > Cd > Cs. The last two elements were found at
334 ultra-trace levels varying between 5 and 25 µg kg⁻¹. Co concentrations (0.02–0.08 mg kg⁻¹)
335 were below those found in cowpea beans (0.7–2.3 mg kg⁻¹) (Santos et al., 2013). Sb contents
336 were in the range of 0.02–0.14 mg kg⁻¹, with an overall average of 80 µg kg⁻¹. However, this
337 element is not commonly found in legumes. On the other hand, As and Cd elements exhibited
338 concentration values of 30 and 17 µg kg⁻¹, respectively; these values were strictly in
339 compliance with the regulatory limits of 100 µg kg⁻¹ for each element in dried legumes
340 imposed by MERCOSUR (MERCOSUR, 2012) and FAO/WHO (Codex Alimentarius,
341 2011).

342

343 3.3. Exploratory statistical analysis

344

345 Firstly, the data matrix that consisted of concentrations of 20 elements from 106 cowpea
346 bean samples was autoscaled, and PCA was applied for exploratory analysis. This helped
347 ensure an equal contribution of variables to the results. PCA reduces the number of variables

348 used in data description, and is the most frequently applied method for computing
349 components (linear latent variables) (Bro & Smilde, 2014; Varmuza & Filzmoser, 2009). The
350 first three principal components (PCs) extracted according to Kaiser Criterion represented
351 16.2%, 11.1%, and 9.1%, respectively, of the variability of the system.

352 Fig. 1 shows the score-plot for PC2 vs. PC1, where one can observe a clear overlap among
353 the scores corresponding to the different cowpea samples identified according to their
354 botanical variety. Samples from the ALA variety showed negative scores in PC1, and were
355 more easily differentiated from the SML and MRH varieties which exhibited positive scores
356 in PC1. However, samples from the CBK and CES varieties presented a combination of
357 positive and negative scores in PC2; the CBK group exhibited a yet distant sample, which
358 tended to complicate the distinction between them. The orientation of the variables in the
359 PC2–PC1 plane is observed from the loading plot shown in Fig. 2. PC1 was strongly
360 influenced by the values of Sn, Mo, Ba, Ni and Zn with positive contributions and by those of
361 V and Sr with negative contributions. The dominant variables in PC2 included Cr, Ni, Tl, Se,
362 Rb, Be and As.

363

364 Insert Fig. 1

365

366 Insert Fig. 2

367

368 In short, the results obtained by PCA showed that the samples from the ALA variety can
369 be distinguished from those of the SML and MRH varieties according to their multielemental
370 fingerprinting. Nonetheless, the last two groups would be difficult to differentiate from one
371 another due to the overlap of their scores. In the same way, although one could differentiate
372 between the samples from the CBK and CES varieties, these varieties cannot be differentiated

373 from the rest of the samples. Through the application of this unsupervised method of pattern
374 recognition, one notices the natural grouping of the 106 cowpea bean samples in the original
375 data matrix, indicating a slight tendency to group some samples. Remarkably, the systematic
376 separation of samples is not clear. Hence, the application of supervised pattern recognition
377 methods for the development of classification models is required to enable one to distinguish
378 the cowpea bean samples from each other according to their botanical varieties.

379

380 *3.4. Predictive modeling for sample classification*

381

382 Upon the completion of the exploratory analysis, several machine learning techniques
383 were applied to the data matrix aiming at the predictive modeling of the cowpea bean
384 varieties. To construct the different classification models, the dataset was divided into two
385 subsets randomly; these included a training set with known class memberships used to
386 calculate the classifiers, and a test set containing samples that were not included in the
387 training and which also had known class memberships that served to validate the models
388 constructed. Considering the different cowpea varieties, a random sampling of the samples
389 was performed in order to balance the group distributions within the splits (stratified
390 sampling). The samples included in each set were randomly changed for each model that was
391 replicated. The training set was formed by 70% of the total samples ($n = 75$), while the
392 remaining 30% ($n = 31$) constituted the testing set. The parameters requiring optimization
393 were calculated during the training stage via the cross-validation technique described above;
394 the maximum accuracy was selected as the best criterion. The samples included in the test set
395 were used to evaluate the performance of the methods developed here against an unknown set
396 of samples.

397 Firstly, the classical LDA method was applied to the whole set of beans samples, where a
398 centroid was created which represented the mean position of all points in all directions. The
399 prediction of results was carried out by projecting the new samples in accordance with the
400 minimal distance to the centroid of each class. Here, a success rate of 87% was recorded in
401 the test set while better sensitivity was obtained for the ALA, MRH and SML varieties.

402 Fig. 3 shows the distribution patterns of cowpea bean samples according to their botanical
403 variety in the plot defined by the first two canonical discriminant functions (DFs). As can be
404 observed, this figure depicts a good discrimination between the four main groups formed by
405 samples from ALA (negative scores on DF2), MRH (positive scores on DF2), SML (majority
406 of positive scores with a negative score farthest on DF2), and CBK + CES (positive and
407 negative scores on DF2). The latter group appears at the center of the graph, exhibiting some
408 degree of difficulty to successfully distinguish the varieties by which it is constituted.

409

410 Insert Fig. 3

411

412 Secondly, the cowpea data matrix was analyzed using PLS-DA aiming at combining the
413 properties of partial least squares regression with the discriminative ability of a classification
414 technique for predictive modeling (Ballabio & Consonni, 2013). Accordingly, the number of
415 significant components for the PLS regression was optimized ($n_{comp} = 2$). Through the
416 combination of these chemometric tools, one is able to obtain additional information
417 regarding the importance of the variables in the classification model, thus enabling the
418 selection of variables and the reduction of noise. Fig. 4 shows the importance of variables for
419 the PLS-DA model. From what can be noted, Ag and Cr were the variables that exerted the
420 highest influence on the discrimination of the CES and CBK groups, while the ALA and
421 MRH groups were mostly influenced by Sn and Mo. For the SML group, none of the

422 variables were found to have stood out among the lots. Surprisingly though, a global accuracy
423 of 77% was obtained despite the fact that the independent variables of the model did not
424 appear to be highly collinear. In view of that, supervised non-linear methods were needed to
425 solve the prediction of varieties in cowpea bean samples as a result of the difficulties
426 encountered in discriminating some sample groups.

427

428 Insert Fig. 4

429

430 For the construction of the k-NN model, the optimal size of neighbor k was optimized, so
431 that when $k = 21$ was applied, the highest average accuracy was obtained. The parameter k
432 indicates the number of neighbors which are considered to predict the unknown sample
433 classes by majority. This method stores all cases and classifies new samples by projection
434 into the multivariate space and attributing these samples to the class of their closest neighbor
435 in the training set. The balanced accuracies for all 106 samples were 63% for ALA samples,
436 50% for CBK samples, 70% for CES samples, 79% for MRH samples and 50% for SML
437 samples; and the overall classification accuracy was 50%. The K-NN model failed to resolve
438 the problem of sample classification; the linear methods studied presented better results.

439 In the RF model, 500 trees were constructed, and an $mtry = 2$ was obtained during the
440 training stage. Each tree was grown using a bootstrapped sample from the original learning
441 sample. However, at each node of the tree, a set of variables were randomly selected while a
442 size was defined. This random selection of features at each node decreases the correlation
443 between the trees in the forest, thereby resulting in a decline in the forest error rate (Archer &
444 Kimes, 2008). Based on the RF classifier, the samples from ALA, CES and MRH groups
445 were classified correctly according to their element profiles; here, the overall classification
446 accuracy obtained was 86% as the samples belonging to the CBK and SML varieties were

447 difficult to differentiate. It is worth pointing out that some works published in the literature
448 have shown an improved accuracy of RF in comparison to other supervised learning methods
449 (Canizo et al., 2017; Villafaña et al., 2017). Oddly enough, the classification performance
450 obtained here was clearly not as expected.

451 By contrast, the SVM classifies the samples by constructing a separate hyperplane in n -
452 dimensional space, maximizing the margin between the classes. This model uses an iterative
453 algorithm, learning the distribution of samples at the boundaries of each of the classes
454 considered (Bona et al., 2017).

455 From the classification results obtained by the LDA and PLS-DA models, a non-linear
456 relation between classes and objects can be considered. In these cases, the use of Kernel trick
457 is recommended as it enables the transformation of the input data in a linearly separable
458 higher-dimensional feature space (Moncayo et al., 2015). In this work, several types of kernel
459 functions were used to carry out the classification: linear kernel, polynomial kernel, gaussian
460 kernel, radial basis function (RBF) and sigmoid kernel. The best results were obtained when
461 the RBF kernel was used. In general, this function is run first due to its ability to handle
462 multivariate data.

463 The complexity of the SVM algorithm is controlled by a penalty error function (C value),
464 which is meant to improve the prediction result and prevent over-fitting (Lantz, 2015). Here,
465 γ is the RBF kernel free parameter. Upon the completion of training, the C and γ were
466 selected, and the model was verified with a set of tests leading to a prediction result. The
467 hyperparameters, $C = 0.5$ and $\gamma = 0.096$, were the best fit for the classification model. Table 4
468 shows the SVM model achieve 100% of sensitivity for prediction of four varieties of cowpea
469 beans (ALA, CBK, CES and MRH). The SML variety reached only 67% which is still
470 considered a good result considering the great similarity between the different groups.
471 Finally, the overall classification accuracy was 93%. Thus, this was the most suitable

472 supervised learning method for predicting the variety of cowpea beans from their
473 multielemental fingerprinting.

474

475 *3.5. Comparing classification performance obtained by the chemometric methods under*
476 *investigation*

477

478 The classification performance of the supervised learning methods investigated in this
479 work was compared in terms of sensitivity, accuracy and kappa metrics. Clearly, Kappa value
480 becomes greatly important in cases involving classifiers with unbalanced datasets (Lantz,
481 2015), as is the case of this work. This statistic provides a more robust measure of agreement
482 than accuracy, since it takes into account the expected agreement by random chance
483 (Stephens & Diesing, 2014).

484 Table 4 summarizes the results obtained for the different algorithms. The order of
485 successful predictions for the models was as follows: SVM > LDA > RF > PLS-DA > k-NN.
486 In general, the samples from the ALA and MRH varieties can be classified by the different
487 models with 100% sensitivity, with the exception of K-NN. LDA and RF models presented
488 similar performance in terms of overall accuracy and sensitivity per class. A further
489 observation that merits mentioning is that the ALA, SML and MRH groups were correctly
490 classified by the LDA method as predicted in Fig. 3. The classification by PLS-DA exhibited
491 an excellent performance for the ALA, CES and MRH samples, but it failed to separate the
492 samples of CBK and SML. When the K-NN method was used, the problem related to the
493 classification of the samples was not possible to be solved. This algorithm failed to
494 differentiate all the cowpeas varieties investigated here, obtaining a very low overall
495 accuracy. The SVM with RBF kernel function was the model that attained the best success
496 rate in the test set, with 100% sensitivity in most of the groups considered.

497

498 Insert Table 4

499

500 On the other hand, according to Cohen (McHugh, 2012), the Kappa statistic can be more
501 easily interpreted if we rely on Table 5. For instance, any kappa value below 0.60 indicates
502 that the confidence intervals obtained for this statistical measure are wide enough; this
503 implies that about half of the data may be incorrect. Such a situation can be noted in the case
504 of the K-NN model, where only 10% (kappa = 0.32) of the data were analyzed correctly,
505 which may have even included anything from good to poor agreement. Thus, the results
506 obtained by this classifier are found to be unreliable. The PLS-DA model showed an
507 agreement of 46% (kappa = 0.68), indicating a moderate level of agreement, since 54% of the
508 data analyzed were erroneous. The RF (Kappa = 0.81) and LDA (Kappa = 0.82) models were
509 within the data reliability range of 64–81%, which indicates a strong agreement of the
510 classifiers. Finally, the SVM model showed a kappa value equal to 0.91, indicating an almost
511 perfect agreement between the predicted model values and the actual values. Furthermore,
512 this classifier obtained the highest overall accuracy among the lots. The model created by
513 SVM was, thus, selected for predicting the botanical variety of cowpea bean samples
514 according to their element profiles.

515

516 Insert Table 5

517

518 Based on the estimated metrics per class, good classification results were obtained by a
519 nonlinear model probably due to the flexibility and ability of the SVM algorithm to create a
520 generalized model (Gaiad et al., 2016). An SVM classifier generally finds a linear or non-

521 linear decision surface in a feature space that separates the training classes with the largest
522 distance between objects on the boundaries (Moncayo et al., 2015).

523 In general, multivariate analysis applied to the cowpea data matrix demonstrated that the
524 concentration of traces elements investigated here varied among the varieties of the same
525 bean species; thus enabling the distinction and authentication of bean samples of different
526 varieties based on their multielemental fingerprinting.

527

528 **4. Conclusions**

529

530 The quantitative cowpea fingerprints consisting of 20 trace elements were obtained by
531 ICP-MS. This technique has been proved to be a fast and reliable tool for generating highly
532 chemical information-rich multielemental fingerprinting. The essential elements, including
533 Co, Cr, Cu, Mo, V, and especially Zn, exhibited good nutritional contribution in the bean
534 samples. For non-essential elements, the concentration levels of As, Cd and Pb showed no
535 significant danger to human health. By projecting multielemental data on computing platform
536 of five supervised learning techniques, the study demonstrated that there are significant
537 differences among the concentrations of some elements across different cowpea varieties.
538 The most relevant variations observed were associated with Ag, As, Ba, Cr, Mo, Ni, Rb, Sb,
539 Sn and V, which enabled the successful classification of cowpea samples according to their
540 botanical variety. Based on the models proposed using the LDA, PLS-DA, k-NN, SVM, and
541 RF algorithms, the classification performance of SVM was found to be the best, as it attained
542 a success rate of 93% during the test step, with 100% sensitivity for most of the predicted
543 classes. In view of that, the SVM model was regarded the most robust classification
544 algorithm for predicting genotype of bean samples from their element profiles.
545 Multielemental determination combined with supervised pattern recognition techniques have

546 proved to be an essentially promising tool for differentiating cowpea varieties. This work
547 unfolds the great potential of multielemental fingerprinting when it comes to evaluating the
548 authenticity of foods, such as cowpea beans.

549

550 **Conflicts of interest**

551 The authors have no conflicts of interest to declare.

552

553 **Acknowledgements**

554 The authors would like to express their sincerest gratitude to the National University of the
555 Northeast (SGCyT-UNNE) and the National Scientific and Technical Research Council (File
556 LH: 172645 CONICET) for granting a postdoctoral scholarship to M. Pérez-Rodríguez.

557

558

559 **References**

560 Alayande, L. B., Mustapha, K. B., Dabak, J. D., & Ubom, G. A. (2012). Comparison of
561 nutritional values of brown and white beans in Jos North Local Government markets.
562 *African Journal of Biotechnology*, *11*, 10135–10140.

563 Ammann, A. A. (2007). Inductively coupled plasma mass spectrometry (ICP MS): a versatile
564 tool. *Journal of Mass Spectrometry*, *42*, 419–427.

565 Archer, K. J., & Kimes, R. V. (2008). Empirical characterization of random forest variable
566 importance measures. *Computational Statistics & Data Analysis*, *52*, 2249–2260.

567 Asgar, M. A. (2013). Anti-diabetic potential of phenolic compounds: A review. *International*
568 *Journal of Food Properties*, *16*, 91–103.

569 Avanza, M., Acevedo, B., Chaves, M., & Añón, M. (2013). Nutritional and anti-nutritional
570 components of four cowpea varieties under thermal treatments: Principal component

- 571 analysis. *LWT - Food Science and Technology*, 51, 148–157.
- 572 Ballabio, D., & Consonni, V. (2013). Classification tools in chemistry. Part 1: linear models.
573 PLS-DA. *Analytical Methods*, 5, 3790–3798.
- 574 Batista, B. L., Silva, L. R. S., Rocha, B. A., Rodrigues, J. L., Berretta-silva, A. A., Bonates,
575 T. O., Gomes, V. S. D., Barbosa, R. M., & Barbosa, F. (2012). Multi-element
576 determination in Brazilian honey samples by inductively coupled plasma mass
577 spectrometry and estimation of geographic origin with data mining techniques. *Food*
578 *Research International*, 49, 209–215.
- 579 Belane, A. K., & Dakora, F. D. (2011). Levels of nutritionally-important trace elements and
580 macronutrients in edible leaves and grain of 27 nodulated cowpea (*Vigna unguiculata* L.
581 Walp.) genotypes grown in the Upper West Region of Ghana. *Food Chemistry*, 125, 99–
582 105.
- 583 Bevilacqua, M., Bucci, R., Magrì, A. D., Magrì, A. L., Nescatelli, R., & Marini, F. (2013).
584 Data handling in science and technology, Chapter 5 - Classification and class-
585 modelling, Vol. 28, pp. 171–233, Oxford: Elsevier.
- 586 Bona, E., Marquetti, I., Varaschim, J., Yasuo, G., Makimori, F., Arca, C., Guimar, L.,
587 Brígida, M., Valderrama, P., & Poppi, R. J. (2017). Support vector machines in tandem
588 with infrared spectroscopy for geographical classification of green arabica coffee. *LWT*
589 *- Food Science and Technology*, 76, 330–336.
- 590 Brereton, R. G., & Lloyd, G. R. (2014). Partial least squares discriminant analysis : taking the
591 magic away. *Journal of Chemometrics*, 28, 213–225.
- 592 Bro, R., & Smilde, A. K. (2014). Principal component analysis. *Analytical Methods*, 6, 2812–
593 2831.
- 594 Canizo, B. V, Escudero, L. B., Pérez, M. B., Pellerano, R. G., & Wuilloud, R. G. (2018).
595 Intra-regional classification of grape seeds produced in Mendoza province (Argentina)

- 596 by multi- elemental analysis and chemometrics tools. *Food Chemistry*, 242, 272–278.
- 597 Codex Alimentarius. (2009). Guidelines for the design and implementation of national
598 regulatory food safety assurance program associated with the use of veterinary drugs in
599 food producing animals. CAC/GL–71, p. 22.
- 600 Codex Alimentarius. (2011). Joint FAO/WHO Food Standards Programme. Codex
601 Committee on contaminants in foods. Fifth Session, The Hague, The Netherlands, 21-
602 25/03/2011.
- 603 Código Alimentario Argentino. (2012). Reglamento Técnico MERCOSUR sobre “Límites
604 máximos de contaminantes inorgánicos en alimentos”. Resolución Conjunta 116/2012
605 y 356/2012 Modificación. Bs. As., 18/7/2012.
- 606 European Community. (2006). Commission Regulation No. 1881/2006. Setting maximum
607 levels for certain contaminants in foodstuffs. DO L 364 20/12/2006, p. 5.
- 608 Hastie, T., Tibshirani, R., & Friedman, J. (2008). The elements of statistical learning: Data
609 mining, inference, and prediction, Second Edition, Stanford University, Springer.
- 610 Hernández-Pereira, E. M., Álvarez-Estévez, D., & Moret-Bonillo, V. (2015). Automatic
611 classification of respiratory patterns involving missing data imputation techniques.
612 *Biosystems Engineering*, 138, 65–76.
- 613 Inobeme, A., Nlemadim, A. B., Obigwa, P. A., Ikechukwu, G., & Ajai, A. I. (2014).
614 Determination of proximate and mineral compositions of white cowpea beans (*Vigna*
615 *Unguiculata*) collected from markets in Minna, Nigeria. *International Journal of*
616 *Scientific & Engineering Research*, 5, 502–504.
- 617 Karbowska, B. (2016). Presence of thallium in the environment: Sources of contaminations,
618 distribution and monitoring methods. *Environmental Monitoring and Assessment*, 188,
619 640–659.
- 620 Kato, L. S., Fernandes, E. A. D. N., Bacchi, M. A., Sarriés, G. A., & Reyes, A. E. L. (2015).

- 621 Elemental characterization of Brazilian beans using neutron activation analysis. *Journal*
622 *of Radioanalytical and Nuclear Chemistry*, 306, 701–706.
- 623 Lantz, B. (2015). *Machine Learning with R*. Second Edition, Packt Publishing Ltd.,
624 Birmingham.
- 625 Liu, X., Xue, C., Wang, Y., Li, Z., Xue, Y., & Xu, J. (2012). The classification of sea
626 cucumber (*Apostichopus japonicus*) according to region of origin using multi-element
627 analysis and pattern recognition techniques. *Food Control*, 23, 522–527.
- 628 McHugh, M. L. (2012). Interrater reliability : the kappa statistic. *Biochemia Medica*, 22, 276–
629 282.
- 630 Moncayo, S., Manzoor, S., & Caceres, J. O. (2015). Chemometrics and Intelligent Laboratory
631 Systems Evaluation of supervised chemometric methods for sample classification by
632 Laser Induced Breakdown Spectroscopy. *Chemometrics and Intelligent Laboratory*
633 *Systems*, 146, 354–364.
- 634 Moreira-Araújo, R. S. R., Sampaio, G. R., Soares, R. A. M., Silva, C. P., Araújo, M. A. M.,
635 & Arêas, J. A. G. (2018). Identification and quantification of phenolic compounds and
636 antioxidant activity in cowpeas of brs xiquexique cultivar. *Revista Caatinga*, 31, 209–
637 216.
- 638 Santos, W. P. C., Gramacho, D. R., & Teixeira, A. P. (2008). Use of Doehlert design for
639 optimizing the digestion of beans for multi-element determination by inductively
640 coupled plasma optical emission spectrometry. *Journal of the Brazilian Chemical*
641 *Society*, 19, 1–10.
- 642 Santos, W. P. C., Hatje, V., Lima, L. N., Trignano, S. V., Barros, F., Castro, J. T., & Korn,
643 M. G. A. (2008). Evaluation of sample preparation (grinding and sieving) of bivalves,
644 coffee and cowpea beans for multi-element analysis. *Microchemical Journal*, 89, 123–
645 130.

- 646 Santos, W. P. C., Santos, D. C. M. B., Fernandes, A. P., Castro, J. T., & Korn, M. G. A.
647 (2013). Geographical characterization of beans based on trace elements after
648 microwave-assisted digestion using diluted nitric acid. *Food Analytical Methods*, *6*,
649 1133–1143.
- 650 Santos, W. P. C., Teixeira, J., Almeida, M., Pires, A., Luis, S., Ferreira, C., Graças, M., Korn,
651 A. (2009). Application of multivariate optimization in the development of an ultrasound-
652 assisted extraction procedure for multielemental determination in bean seeds samples
653 using ICP OES. *Microchemical Journal*, *91*, 153–158.
- 654 Scolaro, B., Kim, H. S. J., & Castro, I. A. (2018). Bioactive compounds as an alternative for
655 drug co-therapy: Overcoming challenges in cardiovascular disease prevention. *Critical*
656 *Reviews in Food Science and Nutrition*, *58*, 958–971.
- 657 Souza, P. A. L., Marcadenti, A., & Portal, V. L. (2017). Effects of olive oil phenolic
658 compounds on inflammation in the prevention and treatment of coronary artery disease.
659 *Nutrients*, *9*, 1087–1109.
- 660 Stephens, D., & Diesing, M. (2014). A comparison of supervised classification methods for
661 the prediction of substrate type using multibeam acoustic and legacy grain-size data.
662 *PLoS One*, *9*, e93950.
- 663 Tang, W., Hu, J., Zhang, H., Wu, P., & He, H. (2015). Kappa coefficient: a popular measure
664 of rater agreement. *Shanghai Archives of Psychiatry*, *27*, 62–67.
- 665 Thompson, M., Ellison, S. L. R., & Wood, R. (2002). Harmonized guidelines for single-
666 laboratory validation of methods of analysis (IUPAC Technical Report). *Pure and*
667 *Applied Chemistry*, *74*, 835–855.
- 668 Varmuza, Kurt & Filzmoser, P., (2009). Introduction to multivariate statistical analysis in
669 chemometrics, CRC Press, Taylor & Francis Group.
- 670 Zhou, H., Yang, W. T., Zhou, X., Liu, L., Gu, J.F., Wang, W. L., Zou, J. L., Tian, T., Peng, P.

671 Q., & Liao, B. H. (2016). Accumulation of heavy metals in vegetable species planted in
672 contaminated soils and the health risk assessment. *International Journal of*
673 *Environmental Research and Public Health*, 13, 289–301.

674

675 **Figure captions:**

676

677 **Fig. 1.** Score plot of the first principal component (PC1) versus the second principal component
678 (PC2).

679

680 **Fig. 2.** Loading plot for the original variables in the first two principal components (PCs).

681

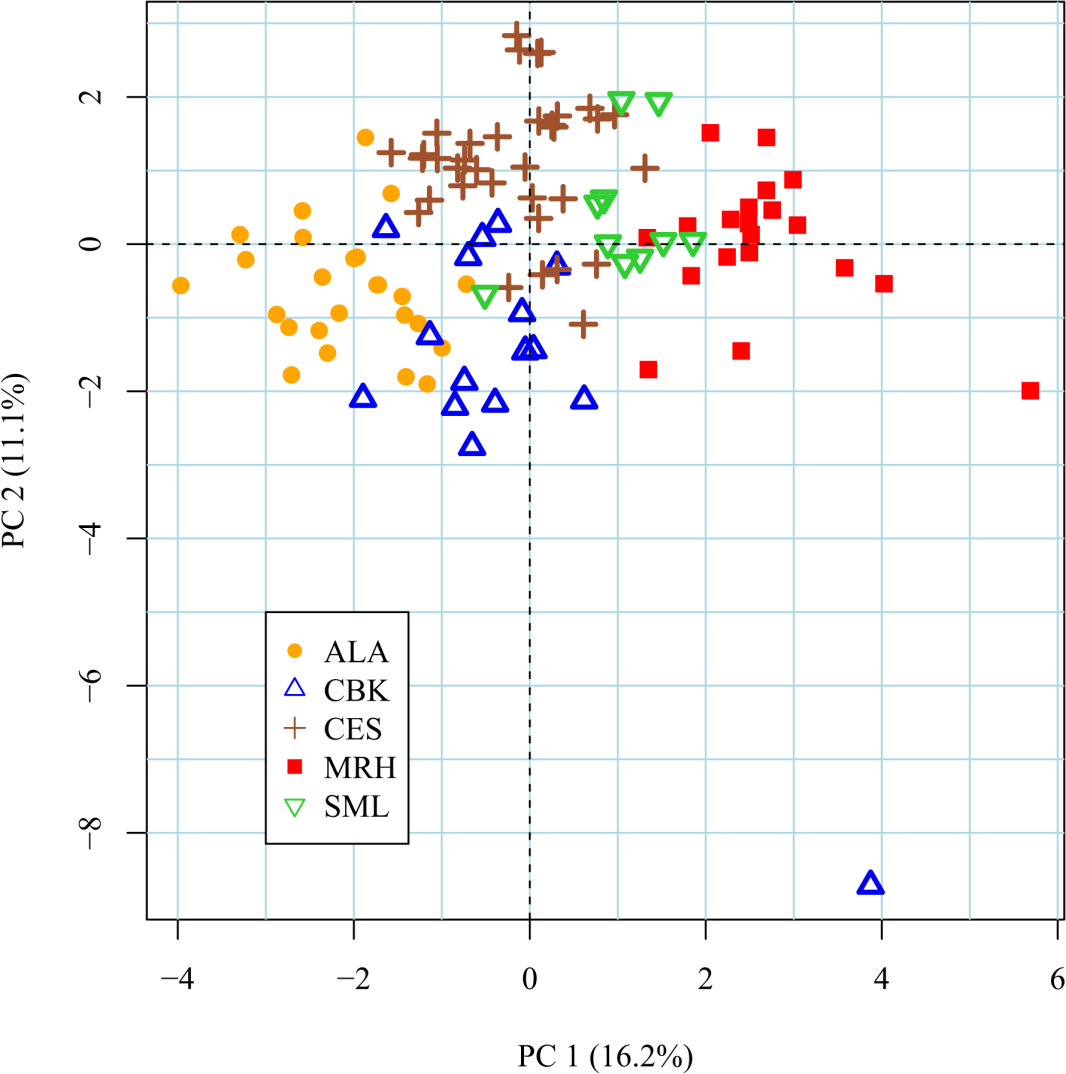
682 **Fig. 3.** Scatter plot of the first two discriminant functions obtained from linear discriminant
683 analysis of cowpea beans according to their botanical variety.

684

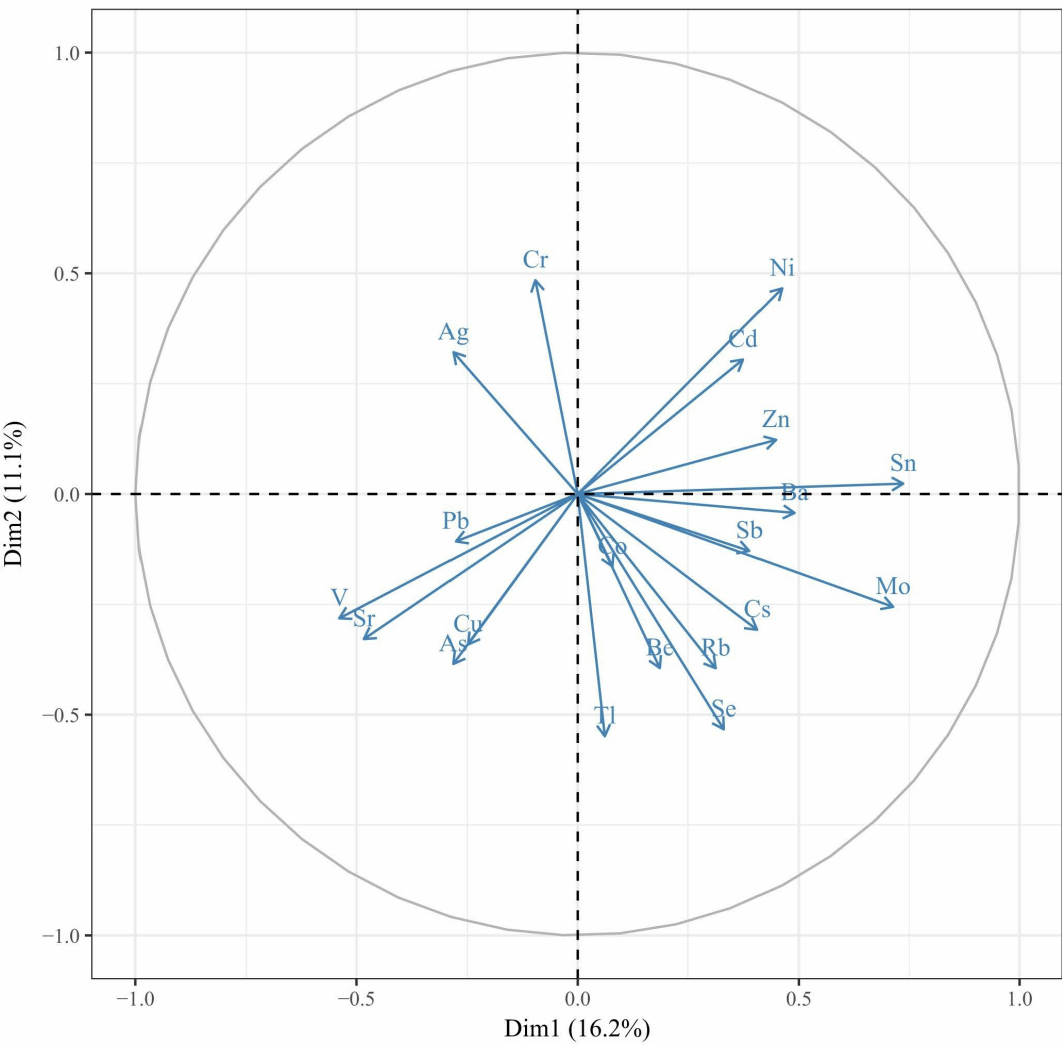
685 **Fig. 4.** Importance of variables for predicting cowpea genotype according to the PLS-DA model.

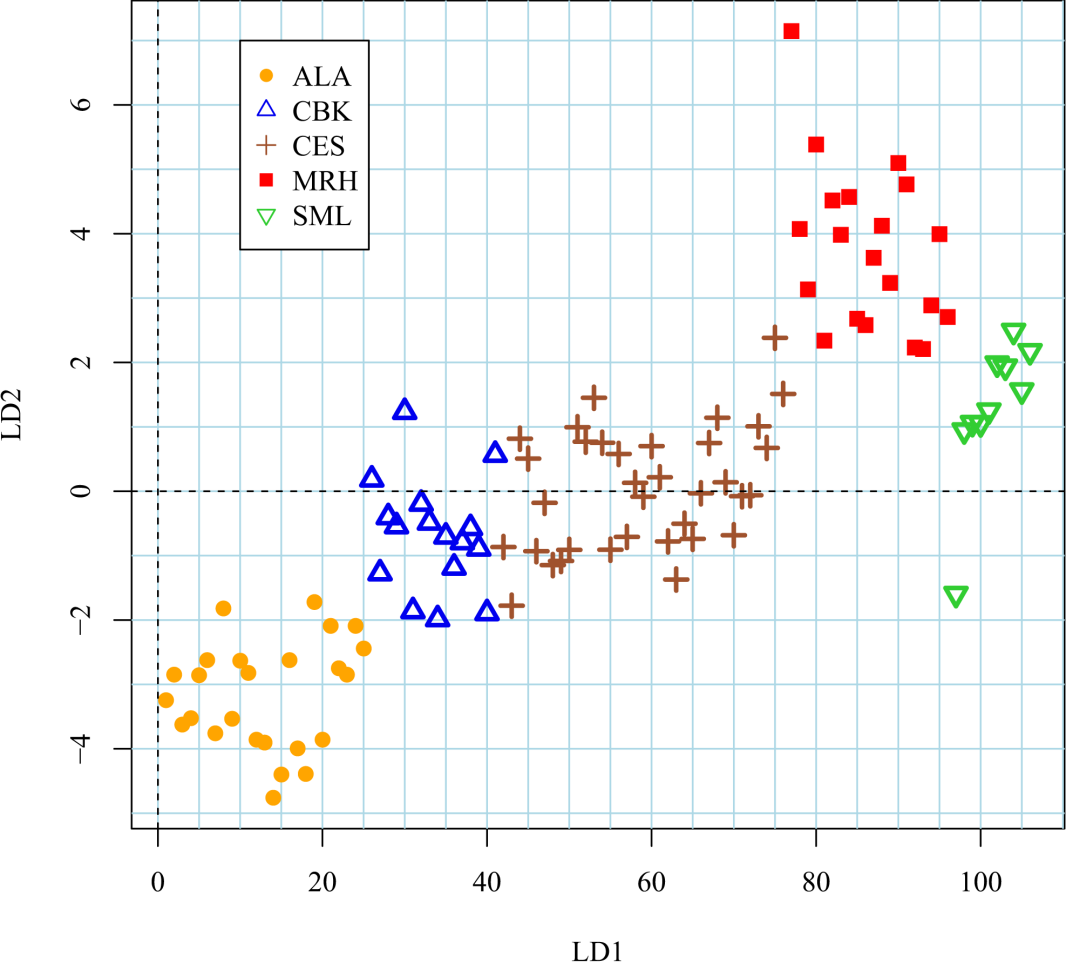
HIGHLIGHTS

- Multielemental fingerprinting was used to verify the authenticity of cowpea beans.
- LDA, PLS-DA, k-NN, RF and SVM were applied to differentiate cowpea genotypes.
- The best performance for predicting cowpea varieties was achieved by SVM classifier.



Variables - PCA





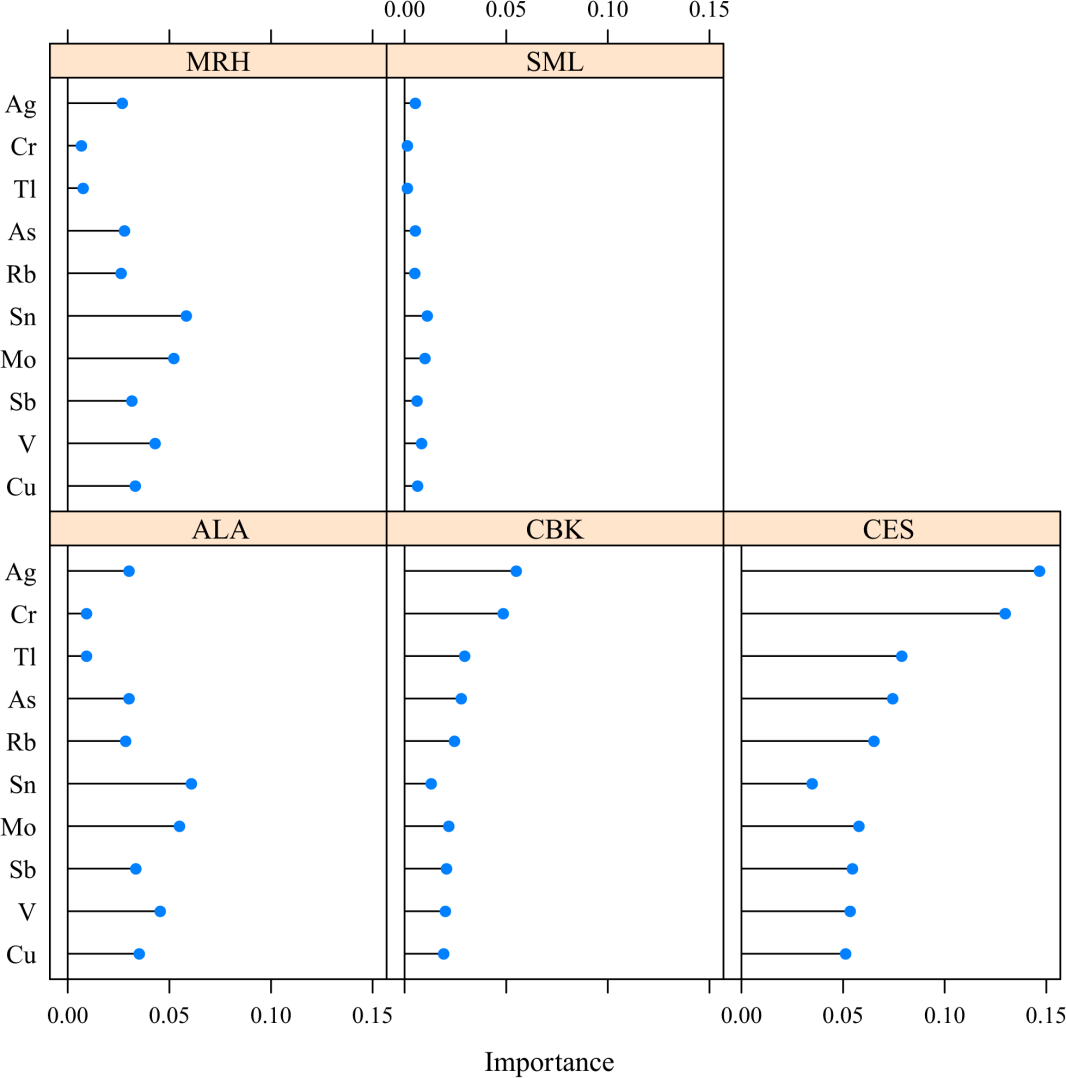


Table 1

Instrumental parameters for ICP-MS determinations.

Radiofrequency power (W)	1321
Sampling depth (mm)	10.0
Gas flow rate (L min ⁻¹)	
Plasma gas (Ar)	13.88
Nebulizer gas (Ar)	0.95
Carrier gas (Ar)	0.85
Auxiliary gas (Ar)	0.80
Sampling cone (Ni, mm)	1.0
Skimmer cone (Ni, mm)	0.4
Collection points/unit mass	3
Internal standard	⁷² Ge, ¹¹⁵ In

Table 2

Coefficients of determination (r^2), limits of quantification (LOQ), intra-day and inter-day precisions and overall recoveries obtained for the determination of multi-elements in bean samples by ICP-MS.

Element	r^2	LOQ (mg kg ⁻¹)	Intra-day (RSD%)	Inter-day (RSD%)	Certified values (mg kg ⁻¹)	Recovery (%) ($n = 5$)	
						CRM	Fortified samples
Ag	0.9989	0.002	4.1	5.5			101.5
As	0.9998	0.010	0.2	2.7	0.112 ± 0.004	97.1	109.2
Ba	0.9993	0.020	0.7	1.8			95.8
Be	0.9996	0.092	1.3	4.8			98.7
Cd	0.9987	0.011	0.5	3.9	1.52 ± 0.04	94.5	98.1
Cs	0.9986	0.003	0.6	5.7			95.3
Co	0.9995	0.011	1.7	1.9	0.57 ± 0.02	93.0	97.3
Cr	0.9997	0.079	1.4	4.9	1.99 ± 0.06	104.5	97.5
Cu	0.9992	0.545	2.6	6.5	4.70 ± 0.14	102.7	100.5
Mo	0.9988	0.205	4.8	6.9			96.0
Ni	0.9994	0.050	3.2	2.7	1.59 ± 0.07	94.2	89.8
Pb	0.9990	0.125	0.5	3.3			86.2
Rb	0.9988	0.015	2.0	3.8	14.89 ± 0.27	95.6	101.8
Sb	0.9998	0.018	0.7	2.8	0.063 ± 0.006	107.0	101.0
Se	0.9993	0.031	4.2	5.9	0.054 ± 0.003	98.4	103.7
Sn	0.9991	0.078	2.2	2.3			108.7
Sr	0.9989	0.007	1.2	7.2			92.7
Tl	0.9995	0.020	0.7	2.4			90.5
V	0.9997	0.010	2.8	1.5	0.835 ± 0.010	97.3	107.5
Zn	0.9985	0.063	1.0	3.7	30.9 ± 0.7	101.5	98.1

Table 3
Elemental composition of cowpea bean samples of different varieties analyzed by ICP-MS.

Element (mg kg ⁻¹)	Sampling varieties (average \pm SD)					Overall range
	ALA (n = 25)	CBK (n = 16)	CES (n = 35)	MRH (n = 20)	SML (n = 10)	
Ag	0.04 \pm 0.04	0.02 \pm 0.03	0.07 \pm 0.01	0.02 \pm 0.01	0.03 \pm 0.02	0.005 – 0.098
As	0.05 \pm 0.02	0.03 \pm 0.01	0.03 \pm 0.01	0.03 \pm 0.01	0.04 \pm 0.01	0.008 – 0.082
Ba	40.74 \pm 16.91	58.76 \pm 15.92	53.04 \pm 14.38	66.57 \pm 16.98	52.26 \pm 15.66	19.8 – 96.6
Be	0.22 \pm 0.02	0.22 \pm 0.03	0.22 \pm 0.02	0.23 \pm 0.02	0.23 \pm 0.02	0.17 – 0.30
Cd	0.015 \pm 0.002	0.013 \pm 0.005	0.018 \pm 0.004	0.020 \pm 0.002	0.020 \pm 0.003	0.010 – 0.025
Cs	0.010 \pm 0.002	0.010 \pm 0.003	0.010 \pm 0.001	0.011 \pm 0.004	0.011 \pm 0.003	0.005 – 0.025
Co	0.05 \pm 0.01	0.05 \pm 0.02	0.05 \pm 0.01	0.05 \pm 0.01	0.05 \pm 0.02	0.02 – 0.08
Cr	0.17 \pm 0.04	0.12 \pm 0.04	0.21 \pm 0.05	0.16 \pm 0.04	0.14 \pm 0.04	0.07 – 0.36
Cu	5.98 \pm 1.05	4.74 \pm 1.29	4.77 \pm 1.11	4.95 \pm 0.99	4.69 \pm 1.44	2.60 – 7.90
Mo	0.69 \pm 0.49	1.58 \pm 0.71	1.02 \pm 0.03	1.85 \pm 0.45	1.36 \pm 0.37	0.26 – 3.05
Ni	0.28 \pm 0.07	0.32 \pm 0.13	0.42 \pm 0.11	0.46 \pm 0.10	0.67 \pm 0.17	0.15 – 0.96
Pb	0.28 \pm 0.07	0.35 \pm 0.11	0.28 \pm 0.06	0.24 \pm 0.09	0.25 \pm 0.04	0.14 – 0.52
Sb	0.06 \pm 0.02	0.09 \pm 0.03	0.09 \pm 0.01	0.08 \pm 0.01	0.07 \pm 0.01	0.02 – 0.14
Se	0.08 \pm 0.01	0.09 \pm 0.04	0.08 \pm 0.01	0.08 \pm 0.01	0.08 \pm 0.01	0.06 – 0.24
Sn	0.16 \pm 0.04	0.20 \pm 0.04	0.23 \pm 0.04	0.48 \pm 0.16	0.23 \pm 0.11	0.10 – 0.99
Sr	4.35 \pm 1.45	5.94 \pm 1.26	3.91 \pm 0.99	2.55 \pm 0.88	3.32 \pm 1.43	1.50 – 8.12
Tl	0.11 \pm 0.02	0.10 \pm 0.02	0.09 \pm 0.03	0.10 \pm 0.02	0.11 \pm 0.02	0.05 – 0.17
Rb	0.17 \pm 0.04	0.26 \pm 0.09	0.17 \pm 0.05	0.23 \pm 0.04	0.19 \pm 0.09	0.05 – 0.45
V	0.15 \pm 0.03	0.14 \pm 0.02	0.13 \pm 0.02	0.11 \pm 0.02	0.11 \pm 0.02	0.08 – 0.20
Zn	32.56 \pm 8.46	32.45 \pm 4.74	35.82 \pm 8.13	42.12 \pm 8.40	33.19 \pm 5.39	17.34 – 60.27

Note: n indicates the number of samples per bean variety.

Table 4
Performance evaluation measurements for different classification models.

Varieties	LDA		PLS-DA (<i>ncomp</i> = 2) ^a		K-NN (<i>k</i> = 21) ^b		FR (<i>nt</i> = 500; <i>mtry</i> = 2) ^c		SVM (<i>C</i> = 0.5; γ = 0.096) ^d	
	Sens (%)	Spec (%)	Sens (%)	Spec (%)	Sens (%)	Spec (%)	Sens (%)	Spec (%)	Sens (%)	Spec (%)
ALA	100	91	100	83	57	70	100	94	100	96
CBK	75	100	–	100	–	100	67	100	100	100
CES	67	100	100	100	70	70	100	93	100	100
MRH	100	96	100	88	67	92	100	94	100	96
SML	100	96	–	100	–	100	–	100	67	100
Accuracy	87		77		50		86		93	
Kappa	0.82		0.68		0.32		0.81		0.91	

^a *ncomp*: number of significant components.

^b *k*: number of k neighbors.

^c *nt*: number of trees; *mtry*: number of variables tried at each split.

^d *C*: penalty factor; γ : γ -insensitive loss function.

Table 5
Interpretation of Kappa statistic.

Kappa value	Agreement level	% of data reliability ^a
0–0.20	None	0–4%
0.21–0.39	Minimal	4–15%
0.40–0.59	Weak	16–35%
0.60–0.79	Moderate	36–62%
0.80–0.90	Strong	64–81%
0.91–1.00	Almost perfect to perfect	82–100%

^a% agreement or squared kappa value