

1 TRACEABILITY OF SOYBEANS PRODUCED IN ARGENTINA BASED ON
2 THEIR TRACE ELEMENT PROFILES

3 Short Title: Chemometric Tools for Classification of Soybean Grains.

4 Melisa J. Hidalgo^{a*}, Diana C. Fechner^a, Davide Ballabio^b, Eduardo J. Marchevsky^c,
5 and Roberto G. Pellerano^a.

6
7
8 ^aInstituto de Química Básica y Aplicada del Nordeste Argentino (IQUIBA-NEA),
9 UNNE-CONICET, Facultad de Ciencias Exactas y Naturales y Agrimensura, Av.
10 Libertad 5400 (3400) Corrientes, Argentina.

11 ^b Milano Chemometrics y QSAR Research Group, Departamento de Ciencias
12 Ambientales, Universidad de Milano-Bicocca, Piazza della Scienza 1 (20126) Milano,
13 Italia.

14 ^cInstituto de Química San Luis (INQUISAL), UNSL-CONICET, Facultad de
15 Química Bioquímica y Farmacia, Av. Ejército de los Andes 950 (5700) San Luis,
16 Argentina.

17
18
19
20 * Corresponding author. Tel: +54-379-4457996; Fax: +54-379-4464793.

21 *E-mail address:* hidalgo.melisa@conicet.gov.ar (M.J. Hidalgo)

22

23

24 **Summary**

25 Soybean (*Glycine max* (L.) Merrill) is a popular foodstuff and crop plant, used in
26 human and animal food. In this work, multielement analysis of soybean grains samples
27 in combination with chemometric tools was used to classify the geographical origins.
28 For this purpose, 120 samples from three provinces of Argentina were analyzed for a
29 panel of 20 trace elements by inductively coupled plasma mass spectrometry (ICP-MS).
30 First, we used principal component analysis (PCA) for exploratory analysis. Then,
31 supervised classification techniques such as support vector machine discriminant
32 analysis (SVM-DA), random forest (RF), k- nearest neighbors (k-NN) and class-
33 modeling techniques such as soft independent modeling of class analogy (SIMCA),
34 potential functions (PF), and one class support vector machine (OC-SVM) were applied
35 as tools to establish a model of origin of samples. The performance of the techniques
36 was compare using global indexes. Among all the models tested, SVM and SIMCA
37 showed the highest percentages in terms of prediction ability in cross-validation with
38 average values of 99.3% for SVM-DA and a median value of balanced accuracy of
39 96.0%, 91.7%, 88.3% for the three origins using SIMCA. Results suggested that the
40 developed methodology by chemometric techniques is robust and reliable for the
41 geographical classification of soybean samples from Argentina.

42

43

44

45 **Keywords:** Soybean grains; Geographical origin; Class-modeling techniques.

46 **1. Introduction**

47 In recent years, the traceability of food products has become increasingly relevant
48 for citizens of many countries mainly interested in food safety and quality. The
49 traceability of a food means that movements can be traced back one step and one step
50 forward anywhere in the supply chain. Most traceability systems are registration
51 systems that document the path of a product from suppliers through intermediate steps
52 to consumers. However, traceability systems mainly depend on the quality of the
53 records and controls that are usually carried out by local food safety authorities¹. Having
54 a traceability system based on the chemical composition of food is a vital tool to
55 guarantee the origin of food, especially for producing countries, such as Argentina.

56 Currently, Argentina is one of the main soybeans (*Glycine max* (L.) Merrit)
57 exporters worldwide. However, this country does not currently offer chemical
58 traceability systems of origin for the soybean produced. This fact acquires great
59 relevance today, considering the soybean production often implies the indiscriminate
60 deforestation of native forests that still serve as a home for native communities, among
61 other environmental and health problems. Having an origin system of chemical
62 traceability of this food would allow its consumers to guarantee that this product was
63 produced respecting environmental and human issues. Today, major food producing
64 companies and even the European Union itself are committed to eliminating
65 deforestation from their global supply chains.

66 In this work, three important Argentine soybean producing areas were considered.
67 The province of Córdoba, which stands out as the main producing region in the heart of
68 Humid Pampas. The province of San Luis, in a border region of Humid Pampas with a
69 lower production than Córdoba but of great importance at the national level. Finally, the

70 province of Chaco, questioned in the last years since it produces soy on the agricultural
71 frontier occupying native forest regions.

72 The authenticity and traceability of the food product can be demonstrated by modern
73 analytical techniques². Reviews have been published about chemometric techniques
74 applied on the geographical origin of foods³⁻⁵. For instance, the geographical origin of
75 soybean seeds was characterized using X-ray fluorescence, showing differences
76 between the trace element contents in soybean from different geographical regions⁶.
77 These results showed that Mg, P, Cl, K, Mn, Cu, Br, and Ba were good parameters for
78 constructing a discriminant model for geographical origin characterization between
79 Japanese and imported samples (from Canada, China, and America). Beside, transgenic
80 and non-transgenic soybean seeds was differentiated according to mineral content
81 analyzed by ICP-MS⁷. The results show that transgenic and non-transgenic soybean
82 seeds show differences in concentrations of Cu, Fe, and Sr.

83 Supervised pattern recognition techniques goal is to establish a classification model
84 based on experimental data to assign unknown samples to a previously defined sample
85 class based on its pattern of measured features⁸. For instance, the discriminant
86 classification methods are appropriate when at least two classes are defined in the
87 problem under study and allow to properly address only multi-class situations. The
88 discriminant classification methods separate the hyperspace in as many regions as the
89 number of sample groups. So, if a sample is matched in the region of space
90 corresponding to a category, it is classified as that category. In reality, each sample is
91 always assigned to one group, even if this sample is not from the studied categories. The
92 most used discriminant classification methods in food chemistry include Linear
93 Discriminant Analysis (LDA), k- Nearest Neighbors (k-NN), Support Vector Machine

94 discriminant analysis (SVM-DA), Random Forest (RF), Partial Least-Squares
95 Discriminant Analysis (PLS-DA), among others⁹.

96 A different approach to supervised pattern recognition are the class-modeling
97 methods that are useful when the focus is on a single class¹⁰. The analogies among the
98 elements of a class in each category are modeled separately. For this reason, the class-
99 modeling methods can be used to study both one-class and multi-class problems. The
100 samples in agreement with the model are assumed as a member of the class, while
101 objects not in agreement are assumed as non-members.

102 When more than one class is modeled, three different situations are possible: each
103 sample can be assigned to a single category or assigned to several categories or not
104 assigned in any category. One of the main advantages offered by class modeling
105 methods is the possibility of recovering samples that are not represented in any of the
106 studied categories. As a consequence, these methods will be able to identify as
107 “foreign” samples those problem samples that correspond to external observations or
108 members of a new class not considered in the training stage. Another advantage of the
109 method is that any additional class can be added without recalculating the already
110 existing class models, as each category is modeled separately. In addition, all class-
111 modeling methods can be used as discriminant tools, while the reverse is not always the
112 case. The most used chemometric class-modeling techniques are Soft Independent
113 Modeling of Class Analogy (SIMCA), Potential Functions (PF), One-Class Support
114 Vector Machine (OC-SVM)¹¹.

115 Many studies in the literature compare the performance of different pattern
116 recognition techniques in food. Hence, the aims of this work were: (1) to characterize
117 the geographical origin of soybeans produced in three regions of Argentina using

118 chemometrics tools applied to trace element compositions, and (2) to compare the
119 classification performance of three discriminant classification methods (k-NN, SVM-
120 DA, RF) and three class-modeling methods (SIMCA, OC-SVM, PF).

121 **2. Experimental**

122 **2.1. Reagents**

123 Mono and multi-element standard solutions of trace analysis grade were purchased
124 from Sigma-Aldrich and Agilent. Ultrapure grade 65% (m/m) HNO₃ and 30% (m/m)
125 H₂O₂ was acquired from Sigma (St. Louis, MO, USA). Nitric acid was further purified
126 by sub-boiling distillation. Water with a resistivity of 18.1 MΩ cm⁻¹ was obtained from
127 a Milli-Q Pluswater purification system Millipore (Molsheim, France). Indium solution
128 100 µg/L obtained from Agilent (Santa Clara, CA) was used. All the chemicals used
129 were of the highest purity available and all the glass materials used were soaked in 10%
130 (v/v) nitric acid and washed with deionized water.

131 **2.2. Samples**

132 A total of 120 samples were collected from six test fields located in the north-central
133 region of Argentina, corresponding to the main soy producing region of this country.
134 The test fields were located at: Almirante Brown (Chaco province, 26°40' S 60°54' W),
135 San Justo (Córdoba province, 31°26' S 62°04' W) and General Pedernera (San Luis
136 province, 33°37' S, 65°19' W). Representative samples of soybeans were collected
137 during the 2018-2019 campaign at different times to create composite samples labeled
138 according to the sampling region. All the samples studied correspond to the botanical
139 species *G. max* and were grown by direct sowing.

140 In the laboratory, soybean seeds were Soybeans were manually separated from pods.
141 Then, allseeds were washed with tap water and rinsed with deionized water. After that,
142 seeds were homogenized with a domestic mixer and stored at $-20\text{ }^{\circ}\text{C}$ in a freezer until
143 analysis.

144 **2.3. Sample preparation**

145 Previous to multielemental determination in botanical samples organic matter
146 should be eliminated (digested). For the digestion of the soybean samples a microwave
147 digestion oven, Milestone® (Chicago, USA) model Ethos One was used. The
148 microwave digestion program was: (1) $25\text{-}200\text{ }^{\circ}\text{C}$ for 15 min, (2) $200\text{ }^{\circ}\text{C}$ for 15 min and
149 (3) $200\text{-}110\text{ }^{\circ}\text{C}$ for 15 min, followed by ventilation at room temperature for 20 min.
150 After cooling to room temperature, the volume was made up to 25 mL with deionized
151 water. Blank solutions and validation spiked samples were prepared in the same way.
152 Prior to use, all plastic containers were soaked in 10% v/v sub-boiling HNO_3 for at least
153 24 h and then rinsed extensively with deionized water. The samples were measured in
154 triplicate.

155 **2.4. ICP-MS analysis**

156 The measurements of trace elements concentrations have been carried out by using
157 an Agilent 7700 x ICP-MS spectrometer (Agilent Technologies, Santa Clara, CA). The
158 instrument is equipped with off-axis ion lens, a quadrupole mass analyzer and an
159 electron multiplier detector. MicroMist glass concentric nebulizer combined with a
160 cooled double-pass spray chamber made of quartz, an octopole collision/reaction system
161 (ORS). The operating parameters for the instrument were described as follow: RF
162 power (1350 W), plasma gas (14.0 L/min), the flow rate of auxiliary gas (0.9 L/min),
163 carrier gas (1.0 L/ min). Indium was used as internal standard. The selected isotopes for

164 measurement were ^{107}Ag , ^{11}B , ^{137}Ba , ^{59}Co , ^{53}Cr , ^{63}Cu , ^{56}Fe , ^7Li , ^{55}Mn , ^{95}Mo , ^{60}Ni , ^{208}Pb ,
165 ^{85}Rb , ^{121}Sb , ^{78}Se , ^{118}Sn , ^{88}Sr , ^{205}Tl , ^{51}V , ^{66}Zn .

166 The accuracy and precision of the ICP-MS method were verified with one Standard
167 Reference Materials from the National Institute of Standards and Technology (NIST),
168 namely SRM 1573a tomato leaves. The precision of the proposed procedure was
169 evaluated by measuring the repeatability and reproducibility. In the repeatability test
170 (within-day precision), the SRM was analyzed three times within one day; and in the
171 reproducibility test (day-to-day precision), sample digestion and ICP MS analysis were
172 studied by triplicate analyses of three aliquots of the SRM on three days for a period of
173 three weeks. All concentration values were found to be in good agreement with the
174 reference values (Table 1). Limit of detection (LOD) and limit of quantification (LOQ)
175 were obtained according to IUPAC guidelines. LOD was calculate as $LOD = \bar{x}_b +$
176 $k \cdot S_b$, where \bar{x}_b is the mean of the blank measures, S_b is the standard deviation of the
177 blank measures, and k is a numerical factor chosen according to a confidence level. The
178 LOQ was defined as 3.3-fold the LOD.

179 TABLE 1

180 2.5. Chemometric models

181 2.5.1. Exploratory data-analysis

182 The results obtained were organized in a matrix with dimension 120 rows (soybeans
183 samples) and 20 columns (element concentrations). Prior to the exploratory analysis of
184 the data matrix, the concentration values of each element corresponding to each sample
185 were autoscaled (each value is subtracted by the mean and divided by the standard
186 deviation). This pretreatment method allows to avoid dimensionality problems between
187 the levels of the different trace elements in the samples.

188 A basic exploratory analysis was performed using principal component analysis
189 (PCA). PCA is a strategic technique that allows knowing relationships between
190 variables, between samples as well as between variables and samples. Characterized by
191 orthogonal linear combinations called principal components. The first components
192 retain the highest percentage of variability present in the initial data set¹².

193 **2.5.2. Supervised classification models**

194 As the objective of this work is to provide a classification model capable to predict
195 the geographical origin of soybean seeds from three principal production regions of
196 Argentina, we firstly perform a comparative study on the performance of three learning
197 classification algorithms. In addition, in a subsequent stage we compare the perform of
198 three class-modeling methods, in order to exploit their comparative advantage in terms
199 of class prediction of future unknown samples.

200 Three supervised classification algorithms were used to classify provenance of
201 soybean seeds. The supervised model uses pre-defined classes to learn through a
202 training phase how data is organized, making possible to predict unlabeled samples
203 based on the classification model. k nearest neighbor (k-NN), support vector machines
204 (SVM-DA) and random forest (RF) are three techniques which have yielded good
205 results in small rectangular data arrays in the literature⁹.

206 k-NN is a distance based non-parametric procedure. The basic idea on which this
207 paradigm is based is that a new sample is going to be classified in the most frequent
208 class to which its k nearest neighbors belong. The value assumed by k is implicitly
209 related to the shape of the decision boundaries that separate the classes. In practice, the
210 optimal value of k is found by some validation procedure¹³.

211 SVM is a supervised technique that produces linear boundaries among the objects of
212 the groups in a transformed space. Three parameters affect the performance of this
213 technique: penalty factor (C value), regularization parameter (ϵ) and the type of kernel
214 function used. Radial basis function (RBF) kernel was used in this study. In the
215 optimization of the parameters (C and ϵ) a grid-search and cross-validation were used to
216 best fit the model and improve the accuracy results¹⁴.

217 RF algorithm is an ensemble learning method. The idea of ensemble learning is to
218 build and combine base learners to obtain a better classification capability. In this
219 technique, multiple trees are generated. Each tree gives a classification (vote for a
220 class). The result is the class with the highest number of votes in the whole forest. As
221 the base learner, random forest uses the CART (classification and regression tree)⁹.

222 Then, in order to propose robust predictive models, three class modeling methods
223 were performed to classify soybean samples. Soft independent modeling by class
224 analogy (SIMCA), one class SVM (OC-SVM) and potential functions (PF).

225 SIMCA was one of the first class-modeling technique introduced in the literature. A
226 principal component analysis is generated for each of the classes present in the data set.
227 The number of principal components that are retained by each class is generally
228 obtained by cross-validation and the number of principal components retained may be
229 different for each class¹⁵.

230 OC-SVM consists in estimating the function that encloses training samples in a
231 hypersphere with a reduced volume. This technique allows to classify only the objects
232 of a class and distinguish them from other objects. RBF was the selected kernel
233 function. This function allows to determine the radius of the hypersphere considering
234 the parameter γ ¹⁶.

235 PF try to estimate the shape of the probability density distribution of the class as a
236 sum of individual contributions of the samples of the class in the training phase. To
237 define the contributions of the samples, different functions can be used¹⁷.

238 **2.5.3. Selection of a test set for external validation of models**

239 In the classification and class modeling phase, the data matrix was random split in
240 training (n = 84) and test (n = 36) sets. The random sampling occurred within each class
241 and preserved the overall class distribution of the data. For discrimination models, the
242 training set was used to tuning the parameters of k-nearest neighbors (k-NN), support
243 vector machine (SVM-DA) and random forest (RF). Optimization of parameters was
244 made using k-fold cross validation (k = 10). Testing set was used to compare the
245 performance of each optimized method. To ensuring that the same resamples are used,
246 internal parameters in R software was used. Finally, to compare the performance of
247 optimized methods we resampled 50 iterations to avoid bias.

248 For class-modeling, internal cross validation, venetian blinds with 5 cross
249 validation groups has been used with training samples to select model parameters such
250 as number of PCs for Soft Independent Modeling of Class Analogy (SIMCA), kernel
251 for Potential Functions (PF) and One-Class Support Vector Machine (OC-SVM). The
252 results were achieved on the 100 iterations for each class with each classifier.

253 **2.5.4. Software**

254 Exploratory and supervised classification analysis were performed using R
255 software¹⁸, version 3.5. For class-modeling methods the Classification toolbox¹⁹ for
256 MATLAB[®] was used.

257 **3. Results and discussion**

258 **3.1. Trace elements in soybean samples**

259 Table 1 shows the median, minimum and maximum concentration of the elements
260 detected in soybean samples from the three geographical origins. Fe and Mn were the
261 most abundant elements in all samples, followed by Zn, Cu, Rb, V, Li, Ti, Sr, Ba, B,
262 Mo and Cr at levels above $1 \mu\text{g kg}^{-1}$, in decrescent order. The concentrations obtained
263 for the 120 samples are provided in the Supplementary Information (SI-1).

264 The non-parametric Kruskal-Wallis test was applied to evaluate differences between
265 the means population of three origins. Cr, Li, Mn and Zn concentrations were
266 significantly different among three pairs of origins ($p < 0.01$) (Table 1). Thirteen out of
267 20 elements exhibited significant differences between the mean ranks of at least one
268 pair of origins ($p < 0.05$), demonstrating that soybean samples from different regions
269 have a characteristic elemental profile. There was no evidence of variation in the
270 concentrations of Ag, Pb, or Se ($p > 0.05$) between any pair of origins.

271 **3.2. Exploratory data-analysis by PCA**

272 PCA was performed based on the concentration of 20 elements determined by ICP-
273 MS in samples of soybean (*G. max*) grains from three geographical origins of
274 Argentina. The first two principal components (PCs) accounted for 51.6% of the total
275 variance. The PC1 summarized 33.4% and the second 18.2% of the variance present in
276 the multielemental results of analyzed samples. As can be seen in the loading-plot
277 obtained from PCA (Figure 1a), PC1 presented a strong positive correlation with the
278 contents of Co, Rb and Sb, and in the direction of the negative values on the x-axis with
279 the contents of Fe and Sr. The representation of mathematical space defined by the first

280 two PCs is completed with the PC2, which shows a strong negative correlation with the
281 concentrations of B, Ni and Cu mainly.

282 FIGURE 1

283 Fig. 1b shows the distribution of samples in the space of the two first computed PCs
284 (score plot). No clear separation is achieved by the samples from Córdoba and San Luis,
285 although some trends can be observed. On the other hand, the Chaco samples appear in
286 negative values of PC1 clearly differentiated from the previous groups, indicating that
287 there are particular characteristics in the multielemental compositions of samples of this
288 group. These differentiation trends are in accordance with the results of the previous
289 Kruskal–Wallis multiple comparison test. As can be seen, the five variables showing
290 statistical differences, such as Fe, Sr, Co, Rb and Sb, also appear as most contributing to
291 the PC1, which is able to group samples in two principal groups.

292 These results of the exploratory analysis by PCA indicate that the contents of trace
293 elements in the samples studied would be useful for the proposal of supervised
294 classification models of geographical origin of soybeans produced in Argentina.

295 **3.3. Supervised classification methods**

296 In a first stage, we begin comparing the yield to correctly classify soybean samples
297 according to their geographical region of origin, applying supervised classification
298 methods. The algorithms selected were k-NN, SVM and RF. These methods were
299 selected due to their great ability to achieve high correct classification rates, especially
300 when only a small number of samples are available.

301 The dataset (120×20) was splitting up into training and testing sets, in a ratio
302 70/30. The partition of data matrix was carried out in a stratified form by random

303 sampling to create balanced splits of the data. The random sampling was performed
304 within each class to preserve the overall class distribution of the data. Training set was
305 used to tuning the hyperparameters of the algorithms k-NN, RF and SVM-DA.
306 Optimization of parameters was made using k-fold CV ($k = 10$). The split of data was
307 repeated 50 times. Finally, testing set was used to compare the performance of each
308 optimized method. Table 2 shows the results of the classification metrics achieved by
309 each technique in the data matrix.

310 TABLE 2

311 As can be seen, Table 2 shows the results achieved by each optimized algorithm in
312 terms of average accuracy, sensitivity and specificity. Sensitivity (also called the true
313 positive rate) describes the positive test samples of each group correctly classified as
314 such, and specificity (also called the true negative rate) measures the ratio of negative
315 test samples belonging to a different group which have been correctly predicted as such.
316 The results of optimized SVM-DA model were the highest, followed by RF and k-NN
317 in this order. This result is consistent with the findings from samples from Brazil²⁰,
318 since it is expected that non-linear techniques (SVM-DA) have greater flexibility in
319 solving non-linear systems. In addition, other commonly used supervised techniques
320 such as LDA or PLS-DA (results not shown) were simultaneously tested with worse
321 results.

322 To refine the classificatory results obtained by SVM-DA, we further optimize the
323 hyperparameters of the algorithm performing 10-fold-cross validation (repeated 50
324 times). Radial basis function (RBF) kernel SVM was selected because of its speed and
325 great capacity to obtain good results in complex systems. The hyperparameters $C = 8$
326 and $\gamma = 0.039$ were the best to obtain the minimal performance error in training setting.

327 Using these optimized hyperparameters, the SVM-DA algorithm achieved a
328 classification rate with a range of 97.3% - 100% of global accuracy. A perfect
329 classification rate was obtained for samples from the provinces of Chaco and San Luis,
330 while only one sample from Cordoba could not be classified correctly. These results
331 indicate that the SVM-DA method is suitable for the geographical classification of
332 soybean samples, being able to differentiate even the samples coming from neighboring
333 provinces (SLS and COR), which as could be observed in the exploratory analysis
334 showed a high degree of similarity.

335 **3.4. Class-modeling methods**

336 Three different class-modeling algorithms (one class classifiers) were used to model
337 trace element compositions of soybean produced in Argentina: SIMCA, PF (Gaussian
338 Kernels) and OC-SVM. Being class modeling methods, they model one class at a time,
339 we have considered the three classes (CHC, COR, SLS) separately.

340 For each class and for each type of classifiers, we used the following validation
341 protocol. We made a double cross-validation with the following protocol, which have
342 been repeated 100 times (iterations): I) Random split samples in training (70%) and test
343 (30%) sets; II) Use the training samples to calibrate the model; internal cross validation
344 has been used with training samples to select model parameters (such as number of PCs
345 for SIMCA, kernel for Potential Functions and one-class SVM); III) Predict the test
346 samples and calculate sensitivity, specificity and their average (balanced accuracy).
347 Thus, test samples do not participate in the model optimization along each iteration.

348 Figure 2 shows the distribution of balanced accuracy obtained for each modelled
349 class obtained on the 100 validation iterations summarized in box plots. Boxplots are a
350 standardized way of displaying the distribution of data based on a five parameters

351 summary (“minimum”, first quartile (Q1), median, third quartile (Q3), and
352 “maximum”).

353 **FIGURE 2**

354 As is shown in Fig. 2, CHC was best modelled class with a median of 96.0% for the
355 three methods studied. However, the OC-SVM method shows a higher density of results
356 around 100% success, indicating a better performance for modelling the samples of this
357 group. It is also observed that class modelling methods have greater difficulties in
358 differentiating samples from COR and SLS provinces, being these classes associated to
359 lower average balanced accuracies (91.7% and 88.3% respectively) and greater
360 dispersion of results over the validation iterations with respect to CHC. The best method
361 for the COR group was SIMCA with a higher average accuracy and a lower dispersion
362 of results. The samples of the SLS group presented greater difficulties to be modelled
363 correctly, showing similar performances between OC-SVM and PF methods. These
364 results, while compatible with the performances achieved by the supervised
365 classification methods, indicate that the recommendation of class modelling techniques
366 for CHC samples is appropriate, with their respective advantages. As an example, Fig. 3
367 shows the Hotelling T2 versus residual Q plot based on 1 PC SIMCA model for CHC
368 class, where the separation of samples from this group is clearly observed.

369 **FIGURE 3**

370 **4. Conclusion**

371 In this study, ICP-MS in combination with chemometric tools was successfully used
372 to classify soybean grains samples of three different geographical origins from
373 Argentina. Supervised classification techniques (RF, SVM-DA, k-NN) and class-

374 modeling techniques (SIMCA, OC-SVM and PF) were performed on trace element
375 compositions of samples. Among all the models tested, SVM-DA and SIMCA showed
376 the highest percentages in terms of prediction ability in cross-validation with average
377 values of 99.3% for SVM-DA and a median value of balanced accuracy of 96.0% for
378 CHC, 92.0% for COR, 88.0% for SLS using SIMCA. It is important to highlight that
379 although the average accuracy for SIMCA reached lower values, the main advantage of
380 this method is that it has a high capacity to identify unmodelled samples, which in the
381 case of supervised classification techniques are unable to detect. For future studies, we
382 expect that some limitations found in our present research can be addressed, such as the
383 expansion of soybean data from other regions.

384 **5. References**

- 385 1. Kemsley EK, Defernez M, Marini F. Multivariate statistics: Considerations and
386 confidences in food authenticity problems. *Food Control*. 2019;**105**:102-112.
387 doi:10.1016/J.FOODCONT.2019.05.021.
- 388 2. Wadood SA, Boli G, Xiaowen Z, Hussain I, Yimin W. Recent development in
389 the application of analytical techniques for the traceability and authenticity of
390 food of plant origin. *Microchem J*. 2020;**152**:104295.
391 doi:10.1016/j.microc.2019.104295.
- 392 3. Badia-Melis R, Mishra P, Ruiz-García L. Food traceability: New trends and
393 recent advances. A review. *Food Control*. 2015;**57**:393-401.
394 doi:10.1016/j.foodcont.2015.05.005.
- 395 4. Borràs E, Ferré J, Boqué R, Mestres M, Aceña L, Busto O. Data fusion
396 methodologies for food and beverage authentication and quality assessment – A

- 397 review. *Anal Chim Acta*. 2015;**891**:1-14. doi:10.1016/j.aca.2015.04.042.
- 398 5. Granato D, Putnik P, Kovačević DB, et al. Trends in Chemometrics: Food
399 Authentication, Microbiology, and Effects of Processing. *Compr Rev Food Sci*
400 *Food Saf*. 2018;**17**(3):663-677. doi:10.1111/1541-4337.12341.
- 401 6. Otaka A, Hokura A, Nakai I. Determination of trace elements in soybean by X-
402 ray fluorescence analysis and its application to identification of their production
403 areas. *Food Chem*. 2014;**147**:318-326. doi:10.1016/j.foodchem.2013.09.142.
- 404 7. Mataveli LRV, Pohl P, Mounicou S, Arruda MAZ, Szpunar J. A comparative
405 study of element concentrations and binding in transgenic and non-transgenic
406 soybean seeds. *Metallomics*. 2010;**2**(12):800-805. doi:10.1039/c0mt00040j.
- 407 8. Berrueta LA, Alonso-Salces RM, Héberger K. Supervised pattern recognition in
408 food analysis. *J Chromatogr A*. 2007;**1158**(1-2):196-214.
409 doi:10.1016/J.CHROMA.2007.05.024.
- 410 9. James G, Witten D, Hastie T, Tibshirani R. *An Introduction to Statistical*
411 *Learning*. Springer; 2013.
- 412 10. Oliveri P, Downey G. Multivariate class modeling for the verification of food-
413 authenticity claims. *TrAC Trends Anal Chem*. 2012;**35**:74-86.
414 doi:10.1016/j.trac.2012.02.005.
- 415 11. Oliveri P. Class-modelling in food analytical chemistry: Development, sampling,
416 optimisation and validation issues – A tutorial. *Anal Chim Acta*. 2017;**982**:9-19.
417 doi:10.1016/J.ACA.2017.05.013.
- 418 12. Bro R, Smilde AK. Principal component analysis. *Anal Methods*.
419 2014;**6**(9):2812-2831. doi:10.1039/C3AY41907J.

- 420 13. Gemperline P. *Practical Guide to Chemometrics*. CRC press; 2006.
- 421 14. Varmuza K, Filzmoser P. *Introduction to Multivariate Statistical Analysis in*
422 *Chemometrics*. CRC press; 2016.
- 423 15. Mees C, Souard F, Delporte C, et al. Identification of coffee leaves using FT-NIR
424 spectroscopy and SIMCA. *Talanta*. 2018;**177**:4-11.
425 doi:10.1016/J.TALANTA.2017.09.056.
- 426 16. Guerbai Y, Chibani Y, Hadjadji B. The effective use of the one-class SVM
427 classifier for handwritten signature verification based on writer-independent
428 parameters. *Pattern Recognit*. 2015;**48**(1):103-113.
429 doi:10.1016/J.PATCOG.2014.07.016.
- 430 17. Marini F. *Chemometrics in Food Chemistry*. Vol 28. Newnes; 2013.
- 431 18. Team RC. R: A language and environment for statistical computing. 2015.
- 432 19. Ballabio D, Consonni V. Classification tools in chemistry. Part 1: linear models.
433 PLS-DA. *Anal Methods*. 2013;**5**(16):3790-3798. doi:10.1039/C3AY40582F.
- 434 20. Barbosa RM, de Paula ES, Paulelli AC, et al. Recognition of organic rice samples
435 based on trace elements and support vector machines. *J Food Compos Anal*.
436 2016;**45**:95-100. doi:10.1016/j.jfca.2015.09.010.

437 **Figure captions**

438 **Fig. 1.** PCA results: (a) Loading plot of PC1 vs PC2 of PCA performed using all the
439 determined trace element concentrations; (b) the corresponding score plot of PC1 vs
440 PC2 with the scores identified according to their geographical origin: CHC Chaco, COR
441 Córdoba and SLS San Luis.

442

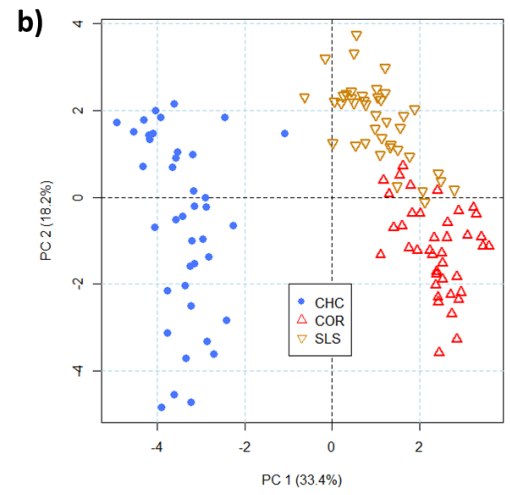
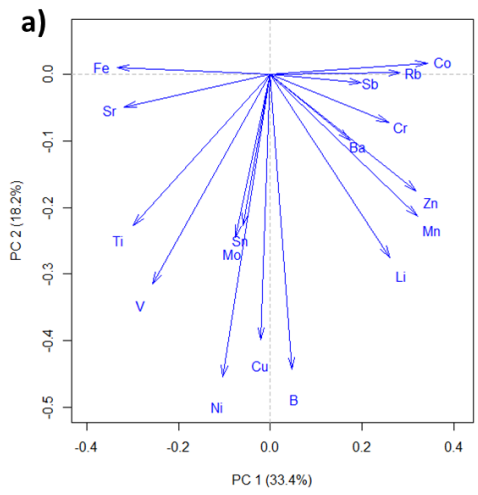
443 **Fig. 2.** Box plot comparing the supervised classification chemometrics models
444 applied for soybean seeds classification.

445

446 **Fig. 3.** Hotelling T2 versus Q residuals for the samples from Chaco (CHC) SIMCA
447 model.

448

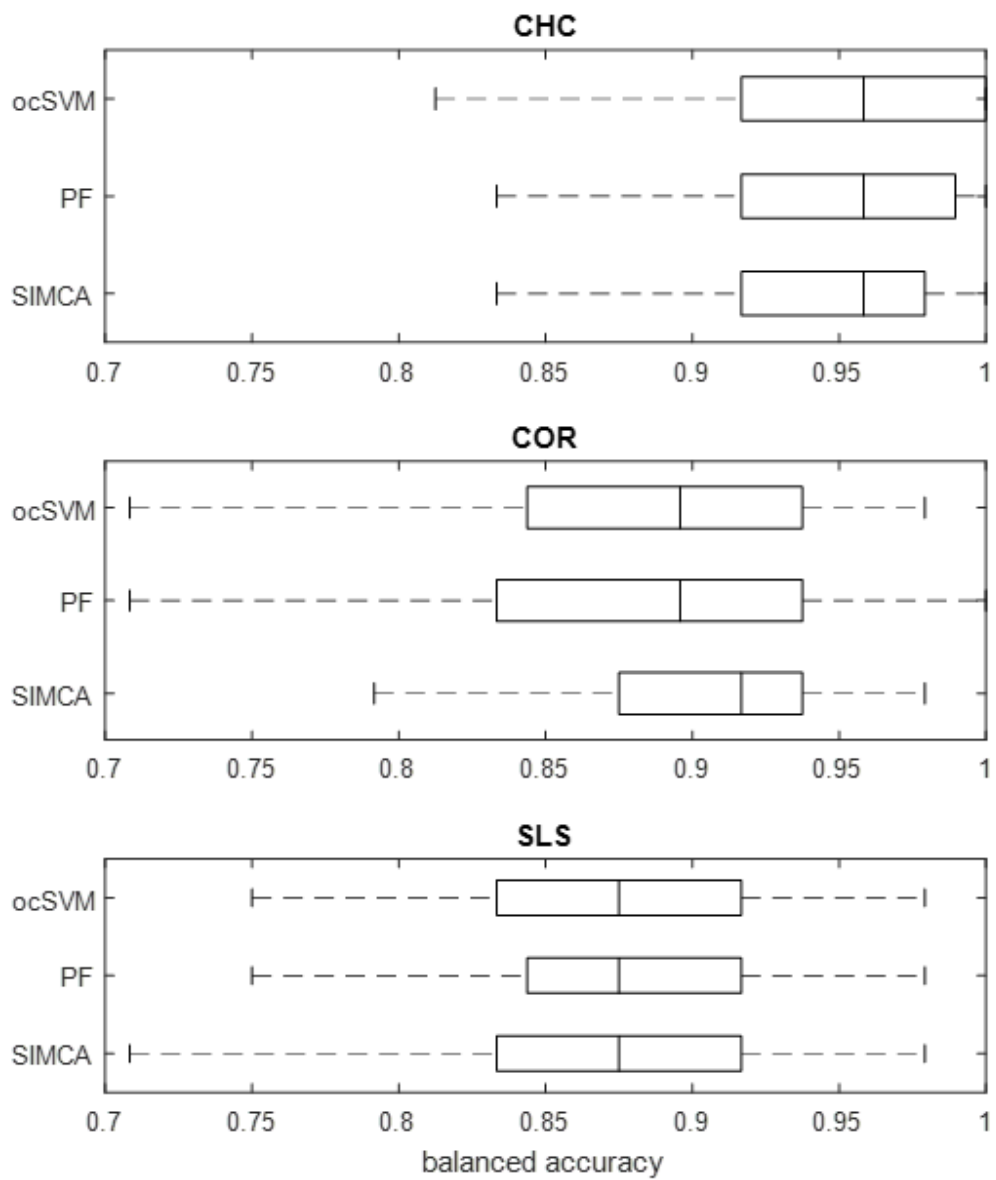
449



450

451 Figure 1

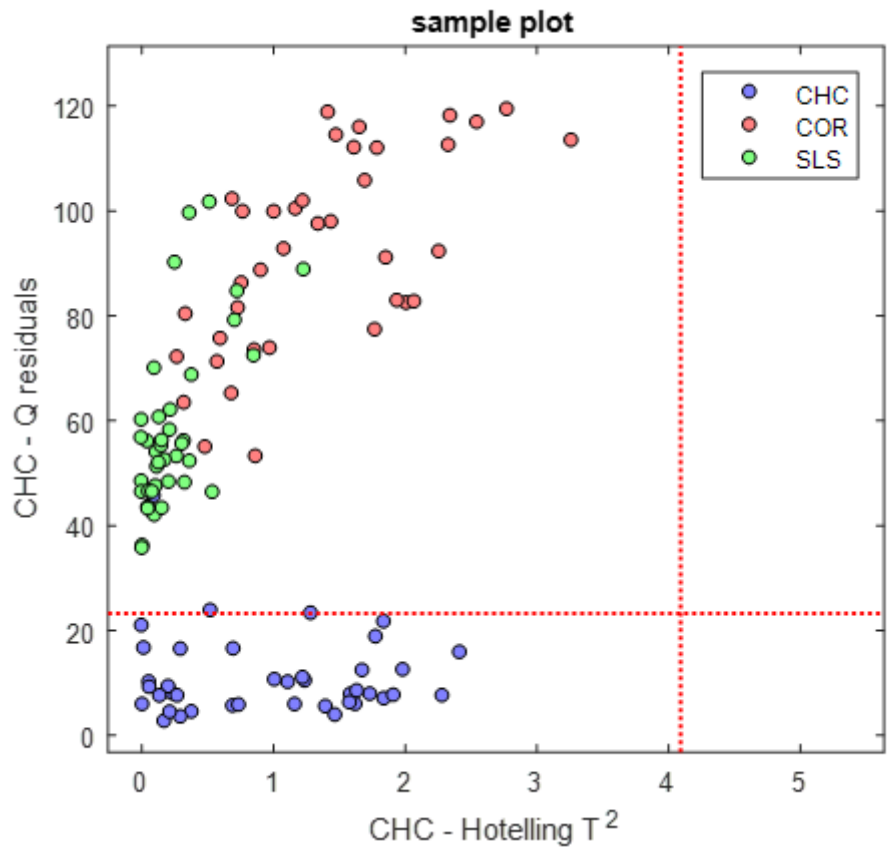
452



453

454 Figure 2

455



456

457 Figure 3

458

459 Table I. Trace element composition of soybean seeds according to their geographical origin.

Element	Certified Values ($\mu\text{g g}^{-1}$)	Recovery Percentage (%)	Geographical origin ($\mu\text{g kg}^{-1}$)			P value
			Chaco <i>n</i> = 40	Córdoba <i>n</i> = 40	San Luis <i>n</i> = 40	
Ag	-	-	0.05 (nd – 0.15)	0.05 (nd – 0.15)	0.05 (nd – 0.10)	ns
B	33.1	98.0	3.3 a (0.3 – 4.6)	3.4 a (1.7 – 4.8)	1.5 b (0.5 – 2.8)	***
Ba	-	-	2.0 a (0.6 – 4.9)	4.8 b (0.4 – 7.2)	3.4 b (2.0 – 6.4)	**
Co	0.58	97.7	1.1 a (0.8 – 1.4)	1.4 b (1.3 – 1.5)	1.4 b (1.3 – 1.4)	**
Cr	1.99	100.1	1.2 a (0.8 – 3.4)	3.2 b (1.3 – 3.9)	1.8 c (0.9 – 3.8)	***
Cu	4.70	98.2	20.2 a (2.8 – 44.1)	18.3 b (3.9 – 27.8)	5.1 b (1.9 – 18.4)	***
Fe	367.5	97.9	162 a (140 – 190)	135 b (120 – 148)	138 b (120 – 150)	**
Li	-	-	3.3 a (0.9 – 6.8)	10.6 b (3.0 – 16.7)	4.9 c (0.9 – 14.8)	***
Mn	246.3	104.0	24 a (10 – 69)	98 b (56 – 132)	64 c (33 – 109)	***
Mo	-	-	2.4 a (1.1 – 5.3)	2.6 b (0.5 – 4.6)	0.9 b (0.2 – 3.7)	***
Ni	1.58	100.8	4.2 a (1.3 – 10.8)	4.1 a (1.9 – 5.2)	1.5 b (0.6 – 3.7)	***
Pb	-	-	0.05 (nd – 0.10)	0.05 (nd – 0.08)	0.05 (nd – 0.10)	ns
Rb	14.8	99.3	8.8 a (5.2 – 16.5)	15.4 b (9.2 – 16.3)	14.3 b (6.4 – 16.3)	***
Sb	-	-	0.06 a (nd – 0.06)	0.07 b (nd – 0.07)	0.06 a (nd – 0.07)	**
Se	0.05	98.5	0.12 (nd – 0.12)	0.12 (nd – 0.12)	0.12 (nd – 0.15)	ns
Sn	-	-	0.8 a (nd – 1.4)	2.5 b (nd – 3.2)	1.8 b (nd – 2.7)	***
Sr	-	-	5.1 a (2.1 – 6.2)	3.1 a (2.1 – 3.9)	3.3 b (2.1 – 5.4)	***
Ti	-	-	6.4 a (2.0 – 7.9)	3.3 b (2.2 – 6.1)	2.8 b (2.1 – 3.8)	**
V	0.83	97.6	11.2 a (3.9 – 17.4)	2.4 b (1.7 – 2.5)	2.2 b (2.0 – 3.8)	**
Zn	30.9	99.9	11.3 a (10.5 – 20.4)	30.6 b (20.8 – 31.4)	20.5 c (12.2 – 31.2)	***

460 Nonparametric Kruskal-Wallis test was applied: ns. not significant at $p > 0.05$; *. $p < 0.05$; **. p
461 < 0.01 ; ***. $p < 0.001$. Pairwise comparison, different letters a, b or c, in the same row indicate
462 significant differences ($p < 0.05$).
463
464
465

466 Table II. Classification results achieved with the different chemometrics models.

Method	Number of samples		Classification metrics		
	Training set	Testing set	Balanced accuracy (%)	Sensitivity (%)	Precision (%)
<i>k</i> -NN ^a	28	12	83.4	83.4	83.4
SVM-DA^b	28	12	91.7	91.7	91.7
RF ^c	28	12	83.4	83.4	91.7

467 ^a *k*: number of neighbors = 5

468 ^b C: penalty factor = 16; Gamma: intensive loss function: 0.039

469 ^c nt: number of trees = 500; mtry: number of variables tried in each split = 7

470