# Accepted Manuscript

Use of principal component analysis (PCA) and hierarchical cluster analysis (HCA) for multivariate association between bioactive compounds and functional properties in foods: A critical perspective

Daniel Granato, Jânio S. Santos, Graziela B. Escher, Bruno L. Ferreira, Rubén M. Maggio
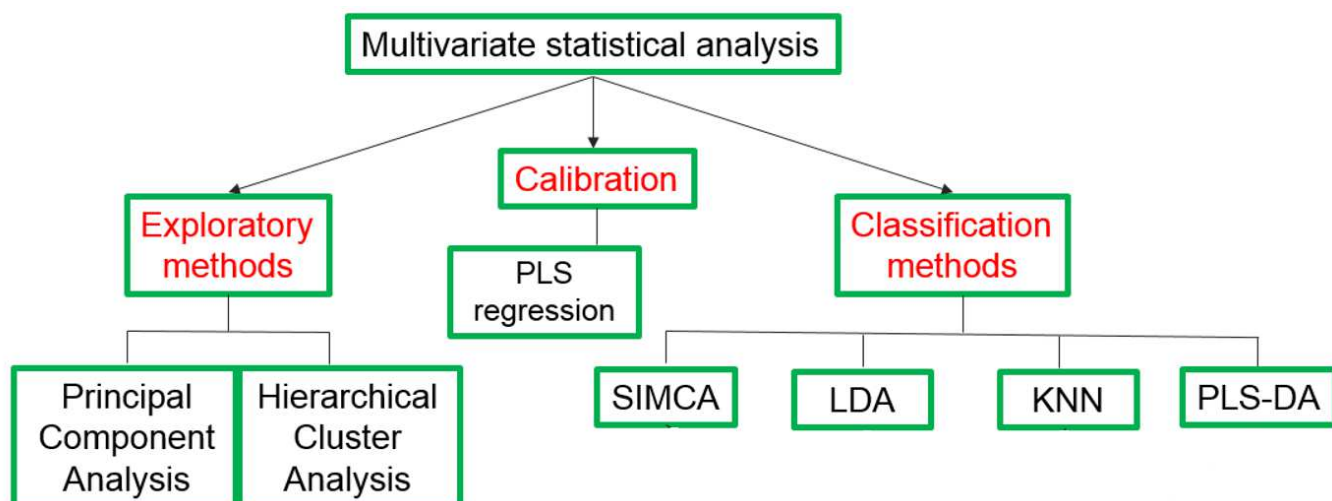
Please cite this article as: Granato, D., Santos, Jâ.S., Escher, G.B., Ferreira, B.L., Maggio, Rubé.M., Use of principal component analysis (PCA) and hierarchical cluster analysis (HCA) for multivariate association between bioactive compounds and functional properties in foods: A critical perspective, *Trends in Food Science & Technology* (2018), doi: 10.1016/j.tifs.2017.12.006.

**GRAPHICAL ABSTRACT**

1   **Use of principal component analysis (PCA) and hierarchical cluster**

2   **analysis (HCA) for multivariate association between bioactive compounds**

3   **and functional properties in foods: a critical perspective**

4

5   Daniel Granato[1]*, Jânio S. Santos[1], Graziela B. Escher[1], Bruno L. Ferreira[2],

6   Rubén M. Maggio[3*]

7

8   [1]Graduation Program in Food Science and Technology, State University of

9   Ponta Grossa. Av. Carlos Cavalcanti, 4748, 84030-900, Ponta Grossa, Brazil.

10  E-mail: dgranato@uepg.br

11  [2]Graduation Program in Food Science, Federal University of Santa Catarina.

12  Rodovia Admar Gonzaga, 1346, Itacorubi, Florianópolis, Brazil.

13  [3]Área Análisis de Medicamentos, Facultad de Ciencias Bioquímicas y

14  Farmacéuticas, Universidad Nacional de Rosario, Instituto de Química de

15  Rosario (IQUIR-CONICET), Suipacha 531, S2002LRK Rosario, Argentina. E-

16  mail: maggio@iquir-conicet.gov.ar

17

**Abstract**

*Background:* The development of statistical software has enabled food scientists to perform a wide variety of mathematical/statistical analyses and solve problems. Therefore, not only sophisticated analytical methods but also the application of multivariate statistical methods have increased considerably. Herein, principal component analysis (PCA) and hierarchical cluster analysis (HCA) are the most widely used tools to explore similarities and hidden patterns among samples where relationship on data and grouping are until unclear. Usually, larger chemical data sets, bioactive compounds and functional properties are the target of these methodologies. *Scope and approach:* In this article, we criticize these methods when correlation analysis should be performed and results analyzed. *Key findings and conclusions:* The use of PCA and HCA in food chemistry studies has increased because the results are easy to interpret and discuss. However, their indiscriminate use to assess the association between bioactive compounds and *in vitro* functional properties is criticized as they provide a qualitative view of the data. When appropriate, one should bear in mind that the correlation between the content of chemical compounds and bioactivity could be duly discussed using correlation coefficients.

**Keywords**: chemometrics; principal component analysis; cluster analysis; correlation analysis; bioactive compounds; functional properties.

**Abbreviations**

ABTS - 2,2'-azino-bis(3-ethylbenzothiazoline-6-sulphonic acid)

2

43 ANN - artificial neural networks

44 CAIMAN - classification and influence matrix analysis

45 DD-SIMCA - data-driven soft independent modeling of class analogy

46 DPPH - 2,2-diphenyl-1-picrylhydrazyl

47 FRAP - ferric reducing antioxidant power

48 FuRES - fuzzy rule-building expert system

49 HCA – hierarchical cluster analysis

50 HPLC – high performance liquid chromatography

51 IMS - ion mobility spectrometry

52 $k$-NN - $k$-nearest neighbors

53 LDA – linear discriminant analysis

54 NMR – nuclear magnetic resonance

55 OPLS-DA - orthogonal partial least squared discriminant analysis

56 ORAC – oxygen radical absorbance capacity

57 PCA – principal component analysis

58 PLS-DA - partial least squared discriminant analysis

59 PRIMA - pattern recognition by independent multi-category analysis

60 QDA - quadratic discriminant analysis

61 RF - random forests

62 SIMCA - soft independent modeling of class analogy

63 sPLS-DA - super partial least squared discriminant analysis

64 SVM - support vector machine

65 UHPLC-MS – ultra-high performance liquid chromatography – mass

66 spectrometry

67

**Introduction**

As well stressed by Ropodi, Panagou, and Nychas (2016), in the 21$^{st}$ century, governmental, industrial, and academic problems need to be addressed by using sophisticated analytical tools with proper data collection, analysis and interpretation. In this sense, data mining and data analysis are two interrelated approaches developed rapidly to address problems related to engineering and technology, as well as medicine, economics, biology, and food science (Brown, 2017).

Chemometrics is an interfacial discipline that extracts useful information from large chemical and biochemical data sets using different mathematical and statistical methods (Nunes et al., 2015, Brown, 2017). In applied chemistry, the use of chemometrics has been spread and well recognized since 1960 (Brereton, 2014), but in food sciences and technology the applications of chemometrics and sensometrics (multivariate methods applied to sensory data and studies consumers) are somewhat new (Munck, Nørgaard, Engelsen, Bro, & Andersson, 1998; Aquino et al., 2014; Qannari, 2017). Conversely, the application of chemometrics for assessing the adulteration and geographical origin of foods based on chemical markers is well established in food science (Granato, Koot, Schnitzler, & van Ruth, 2015; Granato, Margraf, Brotzakis, Capuano, & van Ruth, 2015; Paneque, Morales, Burgos, Ponce, & Callejón, 2017; Giannetti, Mariani, Mannino, & Marini, 2017; Opatić et al., 2018). For example, Garrido-Delgado, Muñoz-Pérez, and Arce (2018) used ion mobility spectrometry (IMS) to determine the origin of the olive oil, quality and adulteration with low-cost vegetable oils. Using different statistical tools, authors were able to predict the level of contaminating oil in olive oil. Therefore, there is

4

93  no doubt that chemometric tools is of fundamental importance to solve real life

94  problems.

95      Granato, Nunes, and Barba (2017) stated that the use of design of

96  experiments together with appropriate statistical data analysis is of pivotal

97  importance to assess the association between nutrition, biology, pharmacology,

98  functional properties and the chemical components of foods and their extracts.

99  In this sense, chemometric tools and other statistical methodologies may be of

100 interest when different food extracts and bioactivities need to be evaluated

101 (Granato, de Araújo Calado, & Jarvis, 2014).

102     In real life applications, chemometrics may be employed in food science

103 and technology studies either to assess similarities/differences between multiple

104 objects (samples) or to project the objects in a two/three-dimensional factor-

105 plane based on various characteristics. Therefore, clusterings can be observed

106 and the reasons for the grouping can be pinpointed (Jandrić, & Cannavan,

107 2017; Lund, Brown, & Shipley, 2017; Erasmus, Muller, Butler, & Hoffman,

108 2018). Additionally, multivariate techniques have been widely used to

109 authenticate/trace the geographical origin of foods, to verify the farming system

110 employed by a company and check whether it complies to the information

111 declared on the label, and to check for adulterations (intentional or not) of foods

112 and raw materials (Granato, Koot, & van Ruth, 2015; Chiesa et al., 2016;

113 Müller-Maatsch, Schweiggert, & Carle, 2016; Tavares et al., 2016; Zhu, Wang,

114 & Chen, 2017; Karabagias et al., 2017; Chung et al., 2017; Giannetti, Mariani,

115 Mannino, & Marini, 2017; Acierno et al., 2018).

116     For example, Luo, Shi, and Feng (2017) aimed to characterize the

117 metabolites of Zhi-Zi-Hou-Po decoction, a traditional Chinese medicine, in rat

bile, urine and feces after oral administration, using untargeted liquid chromatography time of flight mass spectrometry combined with orthogonal partial least squared discriminant (OPLS-DA). After analyzing the experimental data, authors were able to identify 83 compounds, in which 39 were metabolites, in the biological samples. In addition, the metabolic pathway (glucoronidation) by which these metabolites formed after oral administration of the decoction was identified by using OPLS-DA. This research is an example on how chemometric tools are important aids in not only in the food chemistry field but also in the experimental nutrition studies.

According to Brereton (2015), chemometrics users tend to 'follow the crowd' and use indiscriminately the available software without knowing the principles and fundamentals of each method applied in their research data analysis. In food chemistry studies, Principal Components Analysis (PCA) and Hierarchical Cluster Analysis (HCA) are widely (and, sometimes, improperly) applied as "*unsupervised classification*" methods to assess the association between bioactive compounds and *in vitro* functional properties (*i.e.*, antioxidant and inhibition of enzymes). Herein, a critical perspective on these display techniques (PCA and HCA) is made together with some comments on their use in the field of bioactive compounds.

**Study of bioactive compounds and *in vitro* potential functional properties with the use of chemometrics**

Chemometrics may be used for both *qualitative* and *quantitative* analysis of experimental data (Szymanska et al., 2015; Martínez et al., 2017). Determining whether a rice sample comes from European countries or

143 elsewhere based on the NMR spectra or the presence or absence of a chemical

144 compound in a HPLC chromatogram are two typical examples of *qualitative*

145 data. On the other hand, assessing the correlation between the content of

146 chlorogenic acid derivatives and antioxidant activity of coffee brews represents

147 a *quantitative* approach. A summary of selected multivariate statistical methods

148 is shown in Figure 1.

149      Overall, chemometrics may be divided into *calibration*, *classification* and

150 *exploratory* methods. According to Oliveri and Simonetti (2016), chemometrics

151 may be divided into *supervised* and *unsupervised* methods. The first class

152 encompasses a varied number of methods/algorithms, including both qualitative

153 and quantitative approaches. Among qualitative methods, *k*-nearest neighbors

154 (*k*-NN), partial least squares discriminant analysis (PLS-DA), super PLS-DA

155 (sPLS-DA), fuzzy rule-building expert system (FuRES), soft independent

156 modeling of class analogy (SIMCA) and linear or quadratic discriminant analysis

157 (LDA or QDA) are the most used techniques. However, some methods, such as

158 classification and influence matrix analysis (CAIMAN), pattern recognition by

159 independent multi-category analysis (PRIMA), support vector machine (SVM),

160 random forests (RF), and artificial neural networks (ANN), show several

161 applications in food science and technology, especially in the classification and

162 authentication problems (Tian et al., 2017; Torkashvand, Ahmadi, & Nikravesh,

163 2017; Aloglu et al., 2017; Mehretie, Al Riza, Yoshito, & Kondo, 2018).

164      Unsupervised methods, also named *clustering* or *displays methods*, are

165 used to study the data structure, look for similarities between multiple objects,

166 and check for outliers in the data set (Liu, Koot, Hettinga, de Jong, & van Ruth,

167 2018). Mixture models, self-organizing maps, *k*-means, HCA and PCA are

168 representatives of unsupervised methods. However, PCA and HCA are the

169 most used in food and chemistry field, representing both sub-classes

170 visualization and agglomerative algorithms, respectively (Wang, Zeng,

171 Contreras, & Wang, 2017).

172 The goal of multivariate unsupervised methods is to evaluate whether

173 clustering exists in a dataset without using class membership information in the

174 calculations (Beebe, Pell, & Seasholtz, 1998). Natural clustering of

175 samples/objects is the result of understanding the measurement system used to

176 characterize the samples and this union between statistical analysis and

177 analytical methods aids in elucidating the physical reasons for the

178 presence/absence of clustering in the data. For further information on these

179 methods, the reader is referred to existing literature (Oliveri & Downey, 2012; de

180 Oliveira et al., 2015).

181 Here we show some recent applications of unsupervised multivariate

182 techniques in the field of bioactivity of food components. When it comes to

183 studies relating bioactive compounds, almost all reports aim to associate the

184 level of certain chemical compounds, *i.e.*, phenolic compounds and carotenoids,

185 with antioxidant activity and other functionalities. Additionally, a critical

186 perspective on the use of display techniques (PCA and HCA) is made together

187 with some comments on their use in the field of bioactive compounds.

188

189 *Principal component analysis*

190 The term PCA is statistical test that belongs to a group of factor analysis.

191 PCA is a mathematical tool that aims to represent the variation present in the

192 dataset (*i.e.*, responses used to characterize the samples) using a small

193    number of factors. For visual analysis, usually two-dimensional or three-

194    dimensional projection of samples is constructed having the axes (principal

195    components, PC) as the factors. Each PC is a linear combination of the original

196    responses (that retain some correlation among) and PCs are orthogonal to each

197    other. PCs iteratively calculated hold as much variation from original data set as

198    possible, in a way that PC1 explains more the data variation than PC2, and PC2

199    explain more data variation than PC3 and so on. That is why a few PCs explain

200    the variation of a large number of original responses. One possible way to

201    determine the number of PC is based on the Kaiser criterion (Kaiser, 1960):

202    eigenvalues higher than 1 are considered as "significant" in the PCA analysis. In

203    addition, the use of Bartlett's test of sphericity is of interest to check for

204    correlation between responses. This test indicates that the responses are

205    (un)related and therefore (un)suitable for structure detection.

206    Figure 2 contains an example of PCA of fruit juices (*i.e.,* orange, lemon

207    and grape) based on chemical composition and antioxidant activity: the

208    responses used to generate the 2D-scatter plot are based on correlation

209    analysis of each response with the first three PCs. As first step an exploration of

210    cumulative variance explained should be carried out and the Kaiser criterion

211    (eigenvalues higher than 1) may be used to define the number of significant PC.

212    Usually this decision is taken according to pre-established level of variance (90,

213    95, 99, or 99.9%) or based on experimental error.

214    Using a factor loadings analysis (Table 1), PC1 retained about 50% of

215    data variation and differentiate the juice samples according to the contents of

216    caffeic acid, (-)-epicatechin, (+)-catechin, quercetin, luteolin and antioxidant

217    activity   (2,2-diphenyl-1-picrylhydrazyl   –   DPPH,   2,2'-azino-bis   (3-

218  ethylbenzothiazoline-6-sulfonic acid – ABTS, and ferric reducing antioxidant

219  power - FRAP). Similarly, PC2 explained another 30% of variability in the

220  original responses and separates the juices based on FRAP, gallic acid, and 5-

221  *O*-caffeoylquinic acid. PC3 and PC4 explain only 11% of data variance and

222  barely does not differentiate the juice samples. The factor loadings from PC3

223  and PC4 were very low (except for quercetrin/luteolin and ellagic acid,

224  respectively). Factor loadings lower than 0.60 indicate that those variables that

225  do not fit well with the factor solution should possibly be dropped from the

226  analysis, especially if the projection of samples on a factor-plane is based on a

227  2-dimensional graph. As a final comment, the first two PCs explain about 81%

228  of data variance but there remains room for about 19% unexplained variation.

229      Once the representative PCs were found, on the basis of samples

230  differentiation/grouping and variance explained, loading analysis is started in

231  order to find the underlying relationships in the original data structure. In this

232  step loading could be visualized as a regression vector (a vector of correlation

233  coefficients between the original variables with each PC-score). The positive

234  factor loadings indicate that the factor will be higher in the positive axis of that

235  PC. For example, for DPPH, a factor loading of 0.69 was obtained with PC1,

236  which means that the samples located in the right-hand side (*i.e.*, violet stars) of

237  the graph have higher mean DPPH values than the samples located in the left-

238  hand side (*i.e.*, red stars). Similarly, the negative factor loadings indicate that

239  the factor will be higher in the positive axis of that PC. For example, for (-)-

240  epicatechin a factor loading of -0.75 was obtained for PC1, meaning that the

241  samples located in the right-hand side (*i.e.*, violet stars) of the graph have lower

10

242 mean concentrations than the samples located in the left-hand side (*i.e.*, red

243 stars).

244       As a complementary analysis, as an illustrative example, PCA data may

245 be compared to correlation coefficients (Table 2). As shown, the antioxidant

246 activity measured by three different assays (*i.e.*, ABTS, FRAP, and DPPH) is

247 mainly correlated ($p < 0.05$) to caffeic acid, (-)-epicatechin, (+)-catechin,

248 quercetin, and luteolin. FRAP also correlated significantly with gallic acid and 5-

249 O-caffeoylquinic acid. However, if the main objective is to check for association

250 between bioactive compounds and functional properties, correlation analysis

251 should be carried out.

252       For instance, Pearson's correlation coefficients or Spearman's rank

253 correlation coefficients are the choices for normally distributed data and for data

254 do not conform to the normal distribution, respectively (de Oliveira et al., 2015).

255       As a final comment on this topic, there is no scientific need to perform

256 PCA or HCA for data sets that have a similar conclusion as the one shown in

257 the above-mentioned example. However, if the number of responses and

258 samples is quite large and data are quite complex (*i.e.,* NMR spectra), PCA is

259 highly indicated.

260       Dos Santos et al. (2017) quantified 13 phenolic compounds in 96 guava

261 fruit pulps (*Psidium guajava* L.) by HPLC, including (+)-catechin, gallic, ferulic,

262 *trans*-cinnamic, chlorogenic, caffeic, *p*-coumaric, syringic, vanillic, and ellagic

263 acid, rutin, quercetin, and kaempferol. The extraction procedure was optimized

264 using different concentrations of ethyl alcohol and methyl alcohol for 15 to 90

265 min using a sample to solvent ratio between 1:30 and 1:100 w/v. The extracts

266 were also analyzed for total phenolic content, ascorbic acid, and flavonoids,

267 together with the antioxidant activity toward DPPH and ABTS radicals. PCA was

268 able to explain only 60% of data variability with 2 PC, but a clear separation

269 between ripe and green guava fruits was observed from the scatter plot. The

270 main responses that separated the groups were syringic acid, (+)-catechin, *p*-

271 coumaric acid, caffeic acid, ellagic acid, *trans*-cinnamic acid and rutin for the

272 green guava, while for ripe and white guava, the better markers were gallic acid

273 and chlorogenic acid.  As rational subsequent step, authors applied ANN (a

274 supervised algorithm) on same data set to obtain a reliable methodology to

275 classify their samples. ANN showed a suitable separation between not only

276 green and white variety but also ripe and unripe guava fruits. It should be

277 stressed that as data were successfully analyzed by PCA, a linear algorithm,

278 LDA or PLS-DA was the logical way to try.

279     However, in some cases, the differentiation between classes is not so

280 clear (Figure 3A) and outliers (one or more observation point(s) that is/are

281 unusually distant from the other observations) can be detected in the dataset. In

282 this case, the researcher cannot expect a straightforward separation between

283 classes. Almost perfect segregation was obtained when all samples are

284 analyzed after outliers removal (in synthetic data) using only two principal

285 components (PCs), as shown in Figure 3B.

286     Fidelis et al. (2017) evaluated multiple juices from different botanical

287 origins (fruits and other vegetables) in relation to some classes of

288 phenolics/bioactive compounds (tannins, total phenols, flavonoids, *ortho*-

289 diphenols, flavonols, total anthocyanins, and betalains), physicochemical

290 properties (pH, soluble solids, and acidity), and antioxidant effects ($Fe^{2+}$

291 chelating properties, antiradical effect (DPPH, ABTS, and FRAP), Folin-

292  Ciocalteu's reducing capacity, and total reducing capacity. A total of 570 data

293  points (38 juices and 15 responses) were analyzed for patterns using PCA,

294  which explained 72% of data variability with 2 PC and it was possible to pinpoint

295  the juices with higher bioactive compounds and antioxidant activity. PLS-DA

296  was used to discriminate juice groups and authors were able to separate *Citrus*

297  juices from *Super juices* (made with berries) with correct classification rates

298  above 73%, while data-driven SIMCA, which is a one-class classification

299  method, was able to discriminate the juices samples with accuracy higher than

300  86%. In this research, authors concluded that the use of DD-SIMCA may be of

301  interest when the authentication of juices based on phenolic compounds and

302  antioxidant activity need to be performed, especially in quality control programs

303  in the juice industry.

304      Kalaycıoğlu, Kaygusuz, Döker, Kolaylı, & Erim (2017) used PCA to

305  explore only n = 10 Turkish honeybee pollens from distinct origins based on

306  organic acids, carbohydrates, 14 minerals, total phenolic content, and

307  antioxidant activity measured by the DPPH assay. Not surprisingly, the first

308  three principal components explained 71% of data variability and authors claim

309  they "classified" the pollen samples according to the geographical origin of the

310  samples (less than five samples per class, in which, n = 2 chestnuts, n = 1 oak,

311  n = 1 Abana, n = 1 Bayburt, n = 1 Balikesir, n = 1 buckwheat, and n = 3 Anzer).

312  However, results are untrustworthy when such low number of samples are

313  available, so conclusions based on the PC plots should be pondered. According

314  to de Oliveira et al. (2015), for PCA, at least five responses and five objects

315  (samples) need to be part of the dataset.

316    Santos et al. (2016) used PCA to reveal the effects of time (5 – 10 min)

317    and extraction temperature (65 – 85 $^o$C) on phenolic composition and functional

318    properties of aqueous extracts of fermented rooibos (*Aspalathus linearis*). For

319    this purpose, a $2^2$ factorial design with three central points was used to

320    manufacture beverages in which some phenolic acids and flavonoids were

321    quantified using LC-MS/MS, the antioxidant activity (ABTS, FRAP, and total

322    reducing capacity), and the inhibition of $\alpha$-amylase and $\alpha$-glucosidase were

323    determined. As a large amount of data were generated (210 data points),

324    authors performed a PCA to reduce dimensionality of the data. Authors verified

325    that rooibos extracted at 85 $^o$C, regardless of the extraction time, presented the

326    highest levels of phenolic compounds, *in vitro* antioxidant activity, and highest

327    inhibition of the digestive enzymes. Although correlation coefficients were

328    calculated to know which compounds exerted the *in vitro* antioxidant effect,

329    PCA was effective in showing the best technological conditions to produce the

330    infusions with higher bioactive compounds.

331    Farag, Ezzat, Salama, and Tadros (2016) studied the anti-

332    acetylcholinesterase activity and bioactive compounds of four sweet basil

333    species (*Ocimum basilicum, Ocimum africanum, Ocimum americanum* and

334    *Ocimum minimum*) by ultra-performance liquid chromatography quadrupole

335    time of flight mass spectrometry (UPLC/qTOF/MS), PCA was used as

336    exploratory tool and OPLS-DA was used for its further analysis. Twenty one

337    hydroxycinnamic acids, 4 benzoic acid conjugates, 14 flavonoid conjugates, 2

338    alcohols, 5 acyl sugars, 4 triterpenes and 12 fatty acids were identified in the

339    extracts. Using these responses, authors applied PCA and HCA to pinpoint the

340    sweet basis species with higher anti-acetylcholinesterase activity: *O.*

341 *americanum, O. africanum,* and *O. basilicum.* Additionally, OPLS-DA was used

342 to distinguish between *O. basilicum* (official drug) from *O. americanum*, with

343 more than 96% of data variability explained by the classification model.

344

345 *Hierarchical cluster analysis*

346　　　　HCA is a clustering method which explore the organization of samples in

347 groups and among groups depicting a hierarchy (Lee & Yang, 2009). The result

348 of HCA is usually presented in a dendrogram, a plot which shows the

349 organization of samples and its relationships in tree form. There are two main

350 approaches to resolve the grouping problem in HCA, agglomerative or divisive

351 (Figure 4).

352　　　　In the first one, each sample is initially considered a cluster, and

353 subsequently pairs of clusters are merged. In divisive approach algorithm start

354 with one cluster including al samples, recursive splits are performed. Clustering

355 is achieved by use of an appropriate metric of samples distance (usually,

356 Euclidean, Mahalanobis or Manhattan distance) and linkage criterion among

357 groups. Complete, single and average and Ward's linkage are the more

358 common variants of linkage criterions.  Ward's method, based in optimal value

359 of a target function, is a possible choice (Granato, Karnopp, & van Ruth, 2015).

360　　　　HCA has also been extensively used to evaluate the multivariate

361 association between bioactive compounds and bioactivity of foods, beverages

362 and their extracts. For instance, Viapiana et al. (2016) used HCA aiming to

363 associate the relationship between phenolic composition measured by HPLC

364 with the *in vitro* antioxidant activity (FRAP and DPPH assays) of 19 chamomile

365 commercial samples (*Matricaria chamomilla* L.). Overall, caffeic, ferulic, and

15

366  syringic acids were the most effective phenolics in exerting antioxidant activity

367  in the herbal extracts (MeOH:H$_2$O, 80:20 v/v). Linear correlation coefficients

368  were also calculated to display a mathematical proof of such findings (r>0.70,

369  p<0.05). HCA and PCA were used with the aim to tentatively "classify" the

370  commercial samples based on the HPLC fingerprint but no differentiation

371  between samples was achieved. This study shows that PCA/HCA methods not

372  always provide sufficient means to group samples according to the

373  concentrations of bioactive compounds and antioxidant activity indices.

374  Sánchez-Salcedo et al. (2016) used HCA as a tool to propose a

375  polyphenolic fingerprint of white (*Morus alba* L., n=4) and black (*Morus nigra* L.,

376  n=4) mulberry leaves clones. UHPLC-MS identified 31 phenolic compounds,

377  mostly important 20 flavonoids, in more than 120 spectrums analyzed, a very

378  high number of variables for such low number of samples. Ward's method

379  based on Euclidean distance generated three major groups, first characterized

380  by 4 clones of both species presenting high amount of caffeic acid-hexoside,

381  caffeoylquinic acid and kaempferol-malonyl-rutinoside and low content of O-

382  hexoside flavonols. Only one clone of *Morus nigra* formed the second group,

383  representing caffeic acid and cryptochlorogenic acid as characteristics. The last

384  group was formed of 3 clones of two mulberry species, presenting high

385  flavonols containing O-hexoside and a low content of caffeic acid.

386  To study the geographical influence on phenolic content and antioxidant

387  activity in Napirira bean (*Phaseolus vulgaris* L.), Fan and Beta (2017) applied

388  an unsupervised pattern recognition method (HCA) based on the Euclidean

389  distance and Ward's method. Total phenolic content, antioxidant activity, and

390  phenolic compounds (protocatechuic acid, *p*-hydroxybenzoic acid, catechin, *p*-

391 coumaric acid, ferulic acid and sinapic acid) were analyzed in eighteen bean

392 samples from four locations in Malawi. HCA was able to differentiate 3 major

393 groups: group 1 clustered samples with high contents of phenolic compounds

394 and antioxidant activity which were from the high-altitude region; group 2

395 clustered samples that presented low contents of phenolic compounds and

396 antioxidant activity which were from a lower-altitude region; and group 3

397 contained samples with intermediate values of phenolic compounds and

398 antioxidant activity and included samples from both intermediate regions of

399 Malawi. As a conclusion, HCA was a useful tool to associate the phenolic

400 compounds/antioxidant activity with the cultivation region.

401    A good example where algorithm configuration could be decisive to

402 obtain a valid conclusion is illustrated by Kaškonienė et al. (2015). Authors

403 analyzed the total phenolic and flavonoids contents, antioxidant activity and

404 individual phenolic compounds (gallic acid, caffeic acid, ferulic acid, 2-

405 hydroxycinnamic acid, rutin naringenin and quercetin) in 14 pollen samples

406 collected in the Baltic region (Latvia and Lithuania) and two others from Spain

407 and China. Data were treated by HCA using both Spearman's distance and

408 Euclidean distance. Samples were clustered in two groups according to the

409 antioxidant activity. Similarly, Euclidean distance clustered the samples into

410 three major groups according to the geographical regions with clear differences

411 in the phenolic composition. As a conclusion the choice of distance function is

412 not a trivial matter and should be tested when HCA is applied. The use of the

413 only one clustering technique (*i.e.,* k-means or tree-clustering), amalgamation

414 rule (*i.e.,* single linkage, complete linkage, or Ward's method), and distance

415 measure (*i.e*, Euclidean, Manhattan, 1- Pearson r) is not recommended.

17

416    Nayik and Nanda (2016) analyzed the minerals, phenolic composition

417    and antioxidant activity of n = 37 unifloral honeys from Kashmir, India. PCA and

418    HCA were used to assess the effects of the botanical origins of those samples

419    based on the quality parameters and verified that PCA was able to group the

420    samples according to the origin (apple, cherry, saffron and wild bush). The

421    authors claimed that "minerals presented the highest discriminating power" in

422    PCA while samples were "classified" using HCA. The terms "discriminating

423    power" and "classification" are related to supervised chemometric tools, such as

424    LDA/QDA or SIMCA, among other techniques (Popek, Halagarda, & Kursa,

425    2017; Mapelli-Brahm, Hernanz-Vila, Stinco, Heredia, & Meléndez-Martínez,

426    2018; Kasprzyk, Depciuch, Grabek-Lejko, & Parlinska-Wojtan, 2018).

427    Therefore, such terms should be avoided when PCA or HCA are employed.

428

429    *Overall comments on PCA and HCA*

430    Both PCA and HCA are usually used concomitantly in studies covering

431    bioactive compounds and functional properties. To illustrate what is widely seen

432    in published articles, consider the following: n = 20 samples coming from two

433    fruits (A and B) are analyzed for the concentrations of total phenolics,

434    carotenoids, antioxidant activity measured by the oxygen radical absorbance

435    capacity (ORAC) assay, and inhibition of amylase and lipase. Results were

436    analyzed using PCA and the 2D projection is given in Figure 5A: it is possible to

437    see a defined cluster containing fruit "B" and another group containing most "A"

438    fruits. However, there are n = 3 "A" samples that are far from the main "A"

439    group. One could say they are outliers simply by looking at the projection, but

440    this cannot be done as PCA does not "classify" objects. In Figure 5B, HCA was

441 applied using the Ward's method as the amalgamation rule and Euclidean

442 distances were calculated between fruits. Using a linkage distance of 15, only

443 two groups are formed, one containing the "A" fruits and the other containing all

444 "B" fruits. Similarly, if a distance of 5 is considered, there are 1 group containing

445 the "B" fruits and two other groups containing the "A" fruits, which is similar to

446 the PCA results. However, if a linkage distance of 1.5 is considered, a total of 6

447 small groups can be visualized. Using this simple example, it is possible to

448 conclude that HCA is an arbitrary method and should be used for exploratory

449 purposes only. Additionally, neither PCA nor HCA creates a "mathematical

450 model" for classification and authentication purposes. Rather, they only project

451 or display the objects under investigation based on selected responses and

452 grouping of samples may be identified by the user. Moreover, neither PCA nor

453 HCA provides a statistical significance of such similarities (Andrić, Bajusz,

454 Rácz, Šegan, & Héberger, 2016).

455      If the aim is to find an association between bioactive compounds and

456 functional properties using HCA, the method may be applied (Figure 5C). It is

457 an easy and straightforward result: total phenolics and carotenoids are

458 associated with ORAC and inhibition of α-amylase. Conversely, the inhibition of

459 lipase does not seem to be associated with any of the responses. Although

460 HCA shows the existence of association between responses,  however it does

461 not provide a measure of the association (qualitative approach). One alternative

462 to overcome this limitation is to calculate the correlation coefficient and provide

463 a quantitative measure of the correlation between responses. As a matter of

464 fact, the inhibition of lipase is not correlated to the concentrations of total

465 phenolics ($r = -0.022$, $p = 0.927$), carotenoids ($r = 0.213$, $p = 0.367$), and ORAC

466 (r = 0.304, p = 0.193). In this case, the use of HCA is near meaningless as

467 correlation coefficients are robust enough to draw the association between the

468 chemical composition and the functional properties of the fruits.

469       Although PCA and HCA are very useful to study the data structure and

470 find similarities among samples, in most cases, linear correlation coefficients

471 would render very similar interpretations of the results. Indeed, it is widely

472 known and recognized that higher levels of phenolic compounds will render a

473 higher antioxidant activity measured by chemical reactions *in vitro* (Guo, Sun,

474 Yu, & Qi, 2017; Lv, Zhang, Shi, & Lin, 2017). Another main disadvantage of

475 using PCA/HCA in those studies is the real applicability of the observations: it

476 seems that most researchers only use PCA and HCA to increment their data

477 analysis rather than to explain the mechanisms of action and have a strong and

478 in-depth discussion based on a solid hypothesis. In fact, in the field of bioactive

479 compounds, when *in vitro* assays are used, it is somewhat obvious that almost

480 all carotenoids and phenolic compounds will exert antioxidant activity. In this

481 case, correlation coefficients should be calculated and results analyzed.

482

483 **Final comments and recommendations**

484       The use of PCA and HCA in food chemistry studies has increased in the

485 past years because the results are easy to interpret and discuss, especially of a

486 large data set is analyzed. However, the indiscriminate use of multivariate

487 exploratory statistical techniques (PCA and HCA) to assess the association

488 between bioactive compounds and *in vitro* functional properties is criticized as

489 the results will be, in most cases, a *sine qua non* observation. When

490 appropriate, the researcher should bear in mind that the correlation between the

491 content of chemical compounds and bioactivity could be duly discussed using

492 simple linear correlation coefficients.

493

494 **Acknowledgements**

498

499 **References**

500 Acierno, V., Alewijn, M., Zomer, P., & van Ruth, S. M. (2018). Making cocoa

501 origin traceable: fingerprints of chocolates using Flow Infusion - Electro Spray

502 Ionization - Mass Spectrometry. *Food Control*, 85, 245-252.

503 Aloglu, A. K., Harrington, P. B., Sahin, S., Demir, C., & Gunes, M. E. (2017).

504 Chemical profiling of floral and chestnut *honey* using high-performance liquid

505 chromatography-ultraviolet detection. *Journal of Food Composition and*

506 *Analysis*, 62, 205-210.

507 Andrić, F., Bajusz, D., Rácz, A., Šegan, S., & Héberger, K. (2016). Multivariate

508 assessment of lipophilicity scales—computational and reversed phase thin-layer

509 chromatographic indices. *Journal of Pharmaceutical and Biomedical Analysis*,

510 127, 81-93.

511 Aquino, L. F. M. C., Silva, A. C. O., Freitas, M. Q., Felicio, T. L., Cruz, A. G., &

512 Conte-Junior, C. A. (2014). Identifying cheese whey an adulterant in milk:

513 Limited contribution of a sensometric approach. *Food Research International*,

514 62, 233-237.

515    Beebe, K. R., Pell, R. J., & Seasholtz, M. B. Chemometrics: a practical guide.

516    1st ed. New York: Wiley & Sons, 1998, 348 p.

517    Brereton, R. G. (2014). A short history of chemometrics: a personal view.

518    *Journal of Chemometrics*, 28(10), 749-760.

519    Brereton, R. G. (2015). Pattern recognition in chemometrics. *Chemometrics and*

520    *Intelligent Laboratory Systems*, 149, 90–96

521    Brescia, M. A., Alviti, G., Liuzzi, V., & Sacco, A. (2003). Chemometric

522    classification of olive cultivars based on compositional data of oils. *Journal of*

523    *the American Oil Chemists' Society*, 80(10), 945-950.

524    Brown, S. D. (2017). The chemometrics revolution re-examined. *Journal of*

525    *Chemometrics*, 31(1), e2856.

526    Chiesa, L., Panseri, S., Bonacci, S., Procopio, A., Zecconi, A., Arioli, F.,

527    Cuevas, F. J., & Moreno-Rojas, J. M. (2016). Authentication of Italian PDO lard

528    using NIR spectroscopy, volatile profile and fatty acid composition combined

529    with chemometrics. *Food Chemistry*, 212, 296-304.

530    Chung, I. M., Kim, J. K., Yang, J. H.,  Lee, J. H., Park, S. K., Son, N. Y., & Kim,

531    S. H. (2017). Effects of soil type and organic fertilizers on fatty acids and vitamin

532    E in Korean ginseng (*Panax ginseng* Meyer). *Food Research International*, 102,

533    265-273.

534    de Oliveira, C. C., de Araújo Calado, V. M., Ares, G. & Granato, D. (2015).

535    Statistical approaches to assess the association between phenolic compounds

536    and the *in vitro* antioxidant activity of *Camellia sinensis* and *Ilex paraguariensis*

537    teas. *Critical Reviews in Food Science and Nutrition*, 55, 1456-1473.

538    dos Santos, W. N. L., Sauthier, M. C. S., dos Santos, A. M. P., Santana, D. A.,

539    Azevedo, R. S. A., & Caldas, J. C. (2017). Simultaneous determination of 13

540 phenolic bioactive compounds in guava (*Psidium guajava* L.) by HPLC-PAD

541 with evaluation using PCA and Neural Network Analysis (NNA). *Microchemical*

542 *Journal*, 133, 583-592.

543 Erasmus, S. W., Muller, M., Butler, M., & Hoffman, L. C., (2018). The truth is in

544 the isotopes: Authenticating regionally unique South African lamb. *Food*

545 *Chemistry*, 239, 926-934.

546 Farag, M. A., Ezzat, S. M., Salama, M. M., & Tadros, M. G. (2016). Anti-

547 acetylcholinesterase potential and metabolome classification of 4 *Ocimum*

548 species as determined via UPLC/qTOF/MS and chemometric tools. *Journal of*

549 *Pharmaceutical and Biomedical Analysis*, 125, 292-302.

550 Fidelis, M., Santos, J. S., Coelho, A. L. K., Rodionova, O. Y., Pomerantsev, A.,

551 & Granato, D. (2017). Authentication of juices from antioxidant and chemical

552 perspectives: A feasibility quality control study using chemometrics. *Food*

553 *Control*, 73, 796-805.

554 Garrido-Delgado, R., Muñoz-Pérez, M. A., & Arce, L. (2018). Detection of

555 adulteration in extra virgin olive oils by using UV-IMS and chemometric

556 analysis. *Food Control*, 85, 292-299.

557 Giannetti, V., Mariani, M. B., Mannino, P., & Marini, F. (2017). Volatile fraction

558 analysis by HS-SPME/GC-MS and chemometric modeling for traceability of

559 apples cultivated in the Northeast Italy. *Food Control*, 78, 215-221.

560 Granato, D., de Araújo Calado, V. M., & Jarvis, B. (2014). Observations on the

561 use of statistical methods in food science and technology. *Food Research*

562 *International*, *55*, 137-149.

563 Granato, D., Karnopp, A. R., & van Ruth, S. M. (2015). Characterization and

564 comparison of phenolic composition, antioxidant capacity and instrumental taste

565  profile of juices from different botanical origins. *Journal of the Science of Food*

566  *and Agriculture*, 95 (10), 1997-2006.

567  Granato, D., Koot, A., Schnitzler, E., & van Ruth, S. M. (2015). Authentication of

568  geographical origin and crop system of grape juices by phenolic compounds

569  and antioxidant activity using chemometrics. *Journal of Food Science*, 80(3),

570  C584-C593.

571  Granato, D., Margraf, T., Brotzakis, I., Capuano, E., & Ruth, S. M. (2015).

572  Characterization of conventional, biodynamic, and organic purple grape juices

573  by  chemical  markers,  antioxidant  capacity,  and  instrumental  taste

574  profile. *Journal of Food Science*, 80(1), C55-C65.

575  Granato, D., Koot, A., & van Ruth, S. M. (2015). Geographical provenancing of

576  purple grape juices from different farming systems by proton transfer reaction

577  mass  spectrometry  using  supervised  statistical  techniques. *Journal  of  the*

578  *Science of Food and Agriculture*, 95(13), 2668-2677.

579  Granato, D., Nunes, D. S., & Barba, F. J. (2017). An integrated strategy

580  between food chemistry, biology, nutrition, pharmacology, and statistics in the

581  development of functional foods: A proposal. *Trends in Food Science and*

582  *Technology*, *62*, 13-22.

583  Guo, Y., Sun, L., Yu, B., & Qi, J. (2017). An integrated antioxidant activity

584  fingerprint for commercial teas based on their capacities to scavenge reactive

585  oxygen species. *Food Chemistry*, 237, 645-653.

586  Jandrić, Z., & Cannavan, A. (2017). An investigative study on differentiation of

587  citrus fruit/fruit juices by UPLC-QToF MS and chemometrics. *Food Control*, 72,

588  173-180.

589 Kaiser, H. F. (1960). The application of electronic computers to factor analysis.

590 *Educational and Psychological Measurement*, 20,141–151.

591 Kalaycıoğlu, Z., Kaygusuz, H., Döker, S., Kolaylı, S., & Erim, F. B. (2017).

592 Characterization of Turkish honeybee pollens by principal component analysis

593 based on their individual organic acids, sugars, minerals, and antioxidant

594 activities. *LWT – Food Science and Technology*, 84, 402-408.

595 Karabagias, I. K., Louppis, A. P., Karabournioti, S., Kontakos, S.,

596 Papastephanou, C., & Kontominas, M. G. (2017). Characterization and

597 geographical discrimination of commercial *Citrus* spp. honeys produced in

598 different Mediterranean countries based on minerals, volatile compounds and

599 physicochemical parameters, using chemometrics. *Food Chemistry*, 217, 445-

600 455.

601 Kaškonienė, V., Ruočkuvienė, G., Kaškonas, P., Akuneca, I., & Maruška, A.

602 (2014). Chemometric analysis of bee pollen based on volatile and phenolic

603 compound compositions and antioxidant properties. *Food Analytical*

604 *Methods*, 8(5), 1150-1163.

605 Kasprzyk, I., Depciuch, J., Grabek-Lejko, D., & Parlinska-Wojtan, M. (2018).

606 FTIR-ATR spectroscopy of pollen and honey as a tool for unifloral honey

607 authentication. The case study of rape honey. *Food Control*, 84, 33-40.

608 Lee, I., Yang, J. (2009). Common clustering algorithms, in: S.D. Brown, R.

609 Tauler, B. Walczak (Eds.), Comprehensive Chemometrics*,* Elsevier, Oxford,

610 England, 2009, pp. 577-618.

611 Liu, N., Koot, A., Hettinga, K., de Jong, J., & van Ruth, S. M. (2018). Portraying

612 and tracing the impact of different production systems on the volatile organic

613 compound composition of milk by PTR-(Quad)MS and PTR-(ToF)MS. *Food*

614 *Chemistry*, 239, 201-207.

615 Lv, H., Zhang, Y., Shi, J., & Lin, Z. (2017). Phytochemical profiles and

616 antioxidant activities of Chinese dark teas obtained by different processing

617 technologies. *Food Research International*, 100, 486-493.

618 Lund, J. A., Brown, P. N., Shipley, P. R. (2017). Differentiation of *Crataegus*

619 spp. guided by nuclear magnetic resonance spectrometry with chemometric

620 analyses. *Phytochemistry*, 141, 11-19.

621 Luo, K., Shi, Q., & Feng, F. (2017). Characterization of global metabolic profile

622 of Zhi-Zi-Hou-Po decoction in rat bile, urine and feces after oral administration

623 based on a strategy combining LC–MS and chemometrics. *Journal of*

624 *Chromatography B*, 1040, 260-272.

625 Mapelli-Brahm, P., Hernanz-Vila, D., Stinco, C. M., Heredia, F. J., & Meléndez-

626 Martínez, A. J. (2018). Isoprenoids composition and colour to differentiate virgin

627 olive oils from a specific mill. *LWT - Food Science and Technology*, 89, 18-23.

628 Martínez, E. B., Ramos, E. F., Hernández, N. P., Vallejo, L. G. Z., Ruano, N. V.,

629 Ponce, M. V., Mendoza, F. G., & Hernández, A. E. B. (2017). [1]H NMR-based

630 metabolomic fingerprinting to determine metabolite levels in serrano peppers

631 (*Capsicum annum* L.) grown in two different regions. *Food Research*

632 *International*, In Press.

633 Mehretie, S., Al Riza, D. F., Yoshito, S., & Kondo, N. (2018). Classification of

634 raw Ethiopian honeys using front face fluorescence spectra with multivariate

635 analysis. *Food Control*, 84, 83-88.

636 Müller-Maatsch, J., Schweiggert, R. M., & Carle, R. (2016). Adulteration of

637 anthocyanin- and betalain-based coloring foodstuffs with the textile dye

638 'Reactive Red 195' and its detection by spectrophotometric, chromatic and

639 HPLC-PDA-MS/MS analyses. *Food Control*, 70, 333-338.

640 Munck, L., Nørgaard, L., Engelsen, S. B., Bro, R., & Andersson, C. A. (1998).

641 Chemometrics in food science – a demonstration of the feasibility of a highly

642 exploratory, inductive evaluation strategy of fundamental scientific significance.

643 *Chemometrics and Intelligent Laboratory Systems*, 44, 31-60.

644 Nayik, G. A., & Nanda, V. (2016). A chemometric approach to evaluate the

645 phenolic compounds, antioxidant activity and mineral content of different

646 unifloral honey types from Kashmir, India. *LWT - Food Science and*

647 *Technology*, 74, 504-513.

648 Nunes, C. A., Alvarenga, V. O., Sant'Ana, A. S., Santos, J. S., & Granato, D.

649 (2015). The use of statistical software in food science and technology:

650 Advantages, limitations and misuses. *Food Research International,* 75, 270–

651 280.

652 Oliveri, P., & Downey, G. (2012). Multivariate class modeling for the verification

653 of food authenticity claims. *Trends in Analytical Chemistry*, 35, 74-86.

654 Oliveri, P., & Simonetti, R. (2016). Chemometrics for Food Authenticity

655 Applications. In: Downey, G. Advances in Food Authenticity Testing.

656 Amsterdam: Elsevier. 1st ed., p. 701-728.

657 Opatić, A. M., Nečemer, M., Lojen, S., Masten, J., Zlatić, E., Šircelj, H., Stopar,

658 D., Vidrih, R. (2018). Determination of geographical origin of commercial tomato

659 through analysis of stable isotopes, elemental composition and chemical

660 markers. *Food Control*, doi.org/10.1016/j.foodcont.2017.11.013

661 Paneque, P., Morales, M. L., Burgos, P., Ponce, L., & Callejón, R. M. (2017).

662 Elemental characterisation of Andalusian wine vinegars with protected

663  designation of origin by ICP-OES and chemometric approach. *Food Control*, 75,

664  203-210.

665  Popek, S., Halagarda, M., & Kursa, K. (2017). A new model to identify botanical

666  origin of Polish honeys based on the physicochemical parameters and

667  chemometric analysis. *LWT - Food Science and Technology*, 77, 482-487.

668  Qannari, E. M. (2017). Sensometrics approaches in sensory and consumer

669  research. *Current Opinion in Food Science*, 15, 8-13.

670  Ropodi, A. I., Panagou, E. Z., & Nychas, G. J. E. (2016). Data mining derived

671  from food analyses using non-invasive/nondestructive analytical techniques;

672  determination of food authenticity, quality & safety in tandem with computer

673  science disciplines. *Trends in Food Science and Technology*, 50, 11-25.

674  Santos, J. S., Deolindo, C. T. P., Fujita, A., Genovese, M. I., Daguer, H.,

675  Valese, A., Marques, M. B., Rosso, N. D., & Granato, D. (2016). Effects of time

676  and extraction temperature on phenolic composition and functional properties of

677  red rooibos (*Aspalathus linearis*). *Food Research International*, 89, 476-487.

678  Szymanska, E., Gerretzen, J., Engel, J., Geurts, B., Blanchet, L., & Buydens, L.

679  M. C. (2015). Chemometrics and qualitative analysis have a vibrant relationship.

680  *Trends in Analytical Chemistry*, 69, 34–51.

681  Tavares, K. M., Lima, A. R., Nunes, C. A., Silva, V. A., Mendes, E., Casal, S., &

682  Pereira, R. G. F. A. (2016). Free tocopherols as chemical markers for Arabica

683  coffee adulteration with maize and coffee by-products. *Food Control*, 70, 318-

684  324.

685  Tian, Y., Yan, C., Zhang, T., Tang, H., Li, H., Yu, J., Bernard, J., Chen, L.,

686  Martin, S., Delepine-Gilon, N., Bocková, J., Veis, P., Chen, Y., & Yu, J. (2017).

687  Classification of wines according to their production regions with the contained

688     trace elements using laser-induced breakdown spectroscopy. *Spectrochimica*

689     *Acta Part B: Atomic Spectroscopy*, 135, 91-101.

690     Torkashvand, A. M., Ahmadi, A., & Nikravesh, N. L. (2017). Prediction of

691     kiwifruit firmness using fruit mineral nutrient concentration by artificial neural

692     network (ANN) and multiple linear regressions (MLR). *Journal of Integrative*

693     *Agriculture*, 16(7), 1634-1644.

694     Viapiana, A., Struck-Lewicka, W., Konieczynski, P., Wesolowski, M., &

695     Kaliszan, R. (2016). An approach based on HPLC-fingerprint and chemometrics

696     to quality consistency evaluation of *Matricaria chamomilla* L. commercial

697     samples. *Frontiers in Plant Science*, 7, 1-11.

698     Wang, X., Zeng, Q., Contreras, M. M., & Wang, L. (2017). Profiling and

699     quantification of phenolic compounds in *Camellia* seed oils: Natural tea

700     polyphenols in vegetable oil. *Food Research International*, 102, 184-194.

701     Zhu, W., Wang, X., & Chen, L. (2017). Rapid detection of peanut oil adulteration

702     using low-field nuclear magnetic resonance and chemometrics. *Food*

703     *Chemistry*, 216, 268-274.

704

705     **FIGURE HEADINGS**

706

707     **Figure 1**: Summary of selected multivariate statistical methods applied in food

708     research.

709     **Figure 2**: PCA of juice samples based on chemical composition and antioxidant

710     activity: A – represents the number of PCs e the explained variance. B-

711     represents the projection of samples on the factor-plane. For illustration

712 purposes, red starts represent orange juice, green stars represent lemon juices,

713 and violet stars represent grape juices.

714 **Figure 3**: Principal components analysis, PCA, to project different samples (i.e.,

715 fruits from different varieties) based on some selected responses: outlier

716 detection with no separation between varieties (A), no outliers with a clear

717 separation between fruit varieties (B).

718 **Figure 4**: HCA dendrogram for agglomerative algorithm (A) and divisive

719 algorithm grouping flow (B).

720 **Figure 5**: Example of PCA (A) and HCA (B, C) applied to a data set composed

721 of n = 20 fruit samples (A and B) according to the concentrations of total

722 phenolics, carotenoids, antioxidant activity measured by the ORAC assay, and

723 inhibition of lipase and α-amylase.

724

**Table 1**: Factor loadings for illustrating the interpretation of Figure 2.

| Factor | PC1 | PC2 | PC3 | PC4 |
|---|---|---|---|---|
| DPPH | **0.69** | -0.47 | 0.16 | -0.42 |
| ABTS | **0.68** | 0.06 | -0.40 | 0.44 |
| FRAP | **0.63** | **-0.65** | -0.12 | -0.02 |
| Gallic acid | 0.50 | **-0.66** | 0.09 | -0.30 |
| Caffeic acid | **0.81** | -0.23 | 0.22 | 0.15 |
| 5-*O*-caffeoylquinic acid | 0.04 | **-0.70** | -0.50 | 0.27 |
| (+)-Epicatechin | **-0.75** | -0.54 | -0.30 | -0.09 |
| (+)-Catechin | **-0.90** | -0.07 | 0.08 | -0.17 |
| Quercetin | **-0.90** | -0.19 | 0.03 | -0.08 |
| Quercetrin | -0.52 | -0.26 | **0.70** | -0.36 |
| Luteolin | -0.78 | -0.05 | **-0.65** | 0.09 |
| Ellagic acid | 0.13 | 0.48 | -0.46 | **-0.73** |
| Eigenvalue | 5.39 | 3.78 | 0.56 | 0.23 |
| Explained variance (%) | 50.35 | 30.56 | 8.05 | 3.18 |

**Table 2**: Illustrative correlation coefficients to help in the interpretation of the

example shown in Figure 2.

| Responses | DPPH | ABTS | FRAP |
|---|---|---|---|
| DPPH | 1 | | |
| ABTS | 0.899 | 1 | |
| FRAP | 0.946 | 0.947 | 1 |
| Gallic acid | 0.564* | 0.529* | 0.608 |
| Caffeic acid | 0.895 | 0.911 | 0.935 |
| 5-O-caffeoylquinic acid | 0.523* | 0.518* | 0.622 |
| (+)-Epicatechin | 0.875 | 0.812 | 0.804 |
| (+)-Catechin | 0.926 | 0.874 | 0.935 |
| Quercetin | 0.873 | 0.924 | 0.901 |
| Quercetrin | 0.425* | 0.378* | 0.333* |
| Luteolin | 0.788 | 0.829 | 0.845 |
| Ellagic acid | 0.238* | 0.356* | 0.458* |

Note: * denotes $p > 0.05$ while the other correlation coefficients present $p <$

$0.05$.

**A**



**B**

**A**

Principal component 2: 20.39%

Principal component 1: 75.00%



**B**

Linkage Distance

Group 1    Group 2

Group 1    Group 2    Group 3

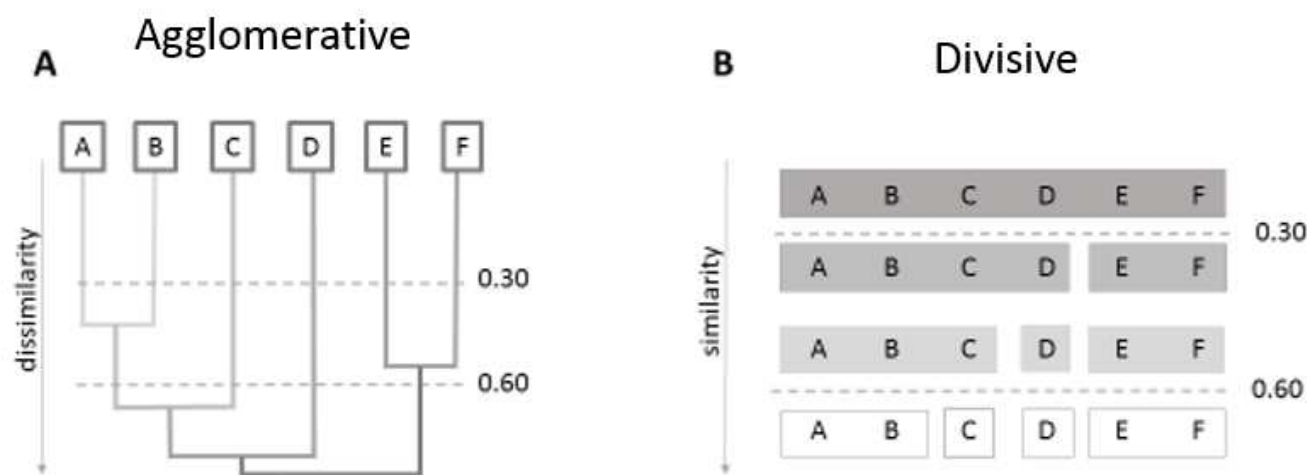Group 1    Group 2    Group 3    Group 4    Group 5    Group 6

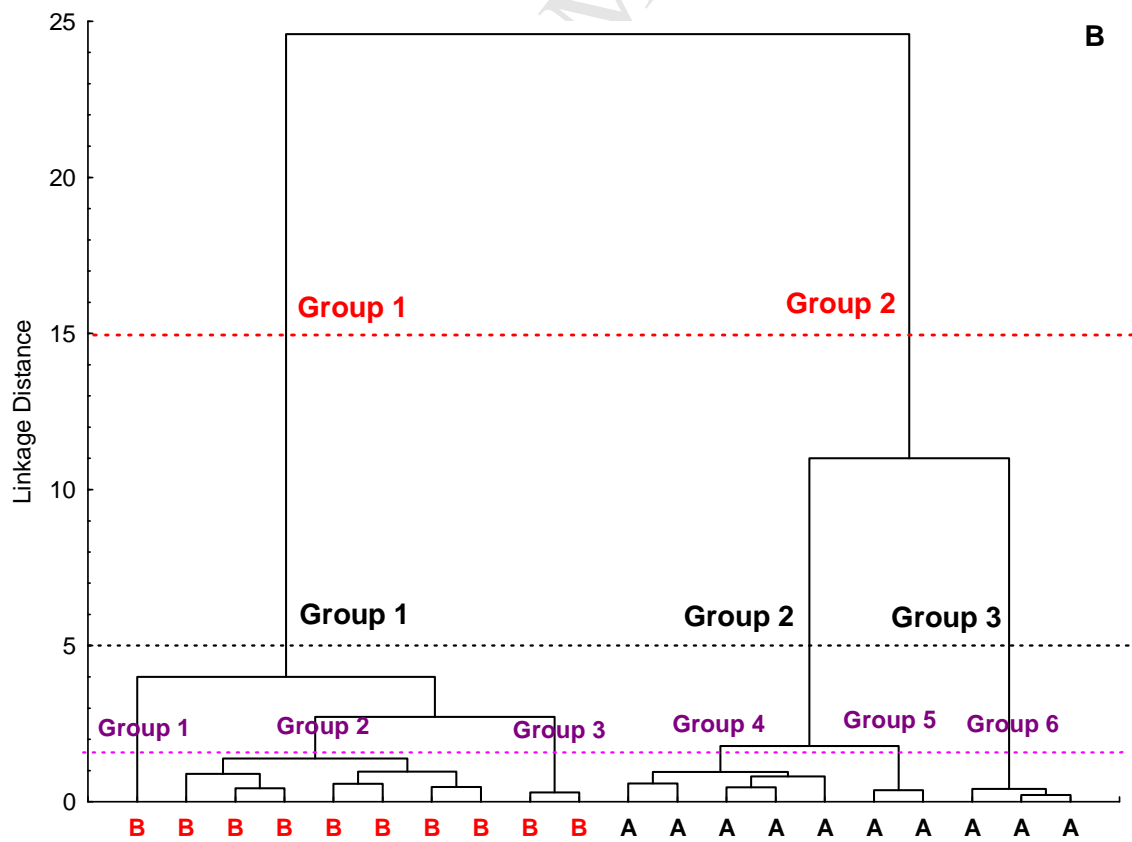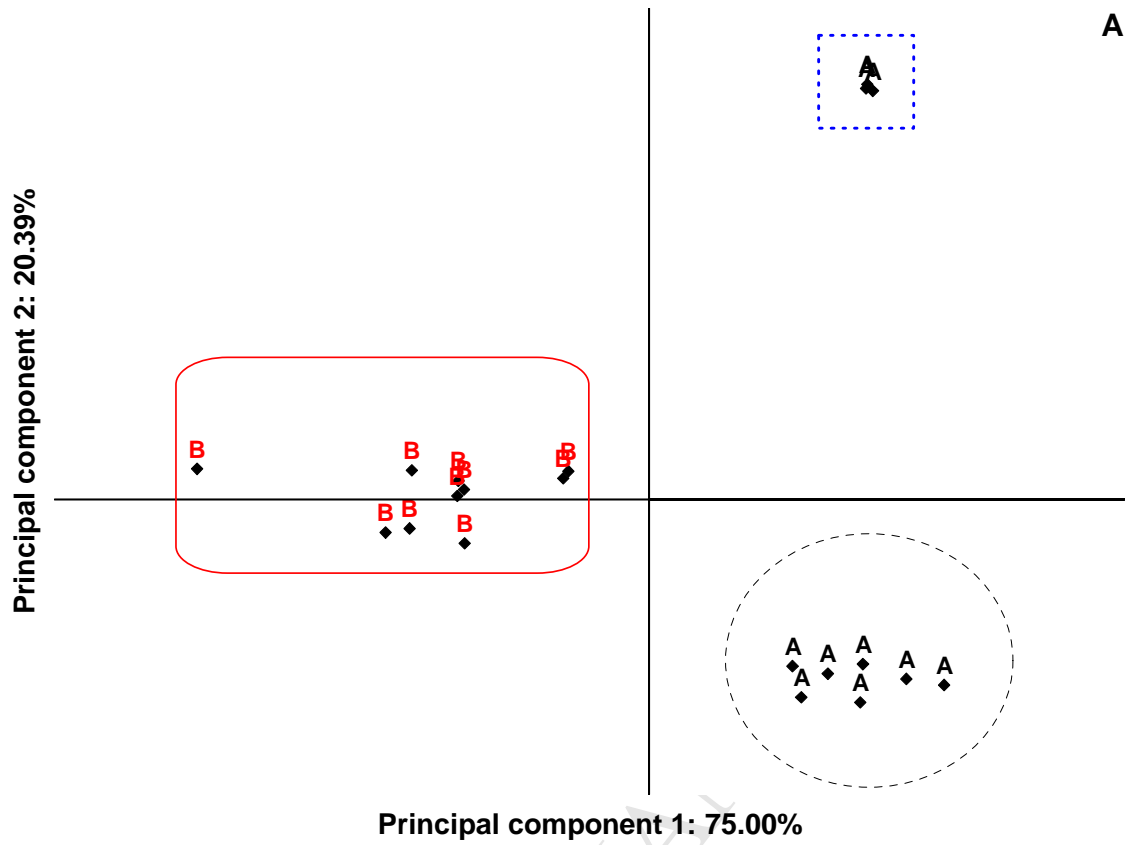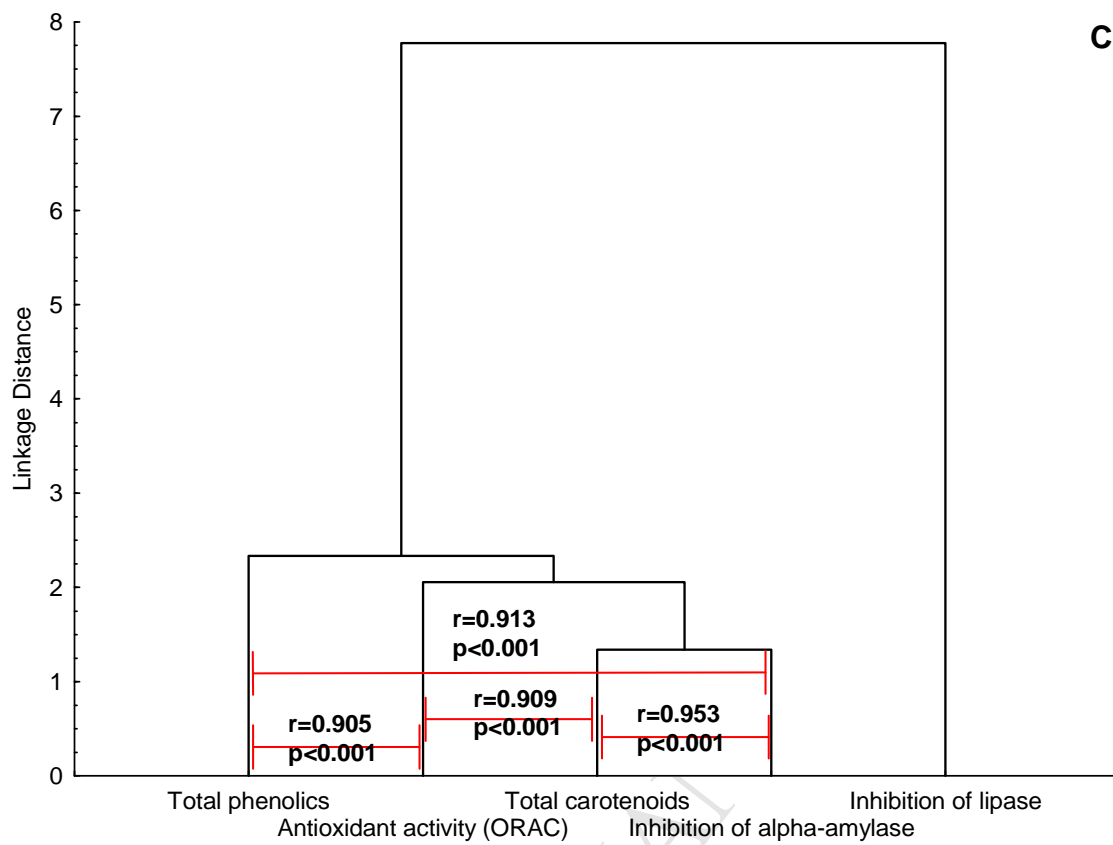B B B B B B B B B B A A A A A A A A A A

**HIGHLIGHTS**

- Chemometric tools are widely used for classification, calibration and exploratory issues

- Unsupervised statistical methods are used to study data structure and look for clusters of samples

- PCA and CA are the most widely used methods

- PCA and CA can be useful in studies regarding bioactive compounds in foods

- We criticize the indiscriminate use of PCA and CA