

Accepted Manuscript

Development of a methodology to analyze leaves from *Prunus dulcis* varieties using near infrared spectroscopy

Sergio Borraz-Martínez, Ricard Boqué, Joan Simó, Mariàngela Mestre, Anna Gras



PII: S0039-9140(19)30599-5

DOI: <https://doi.org/10.1016/j.talanta.2019.05.105>

Reference: TAL 19990

To appear in: *Talanta*

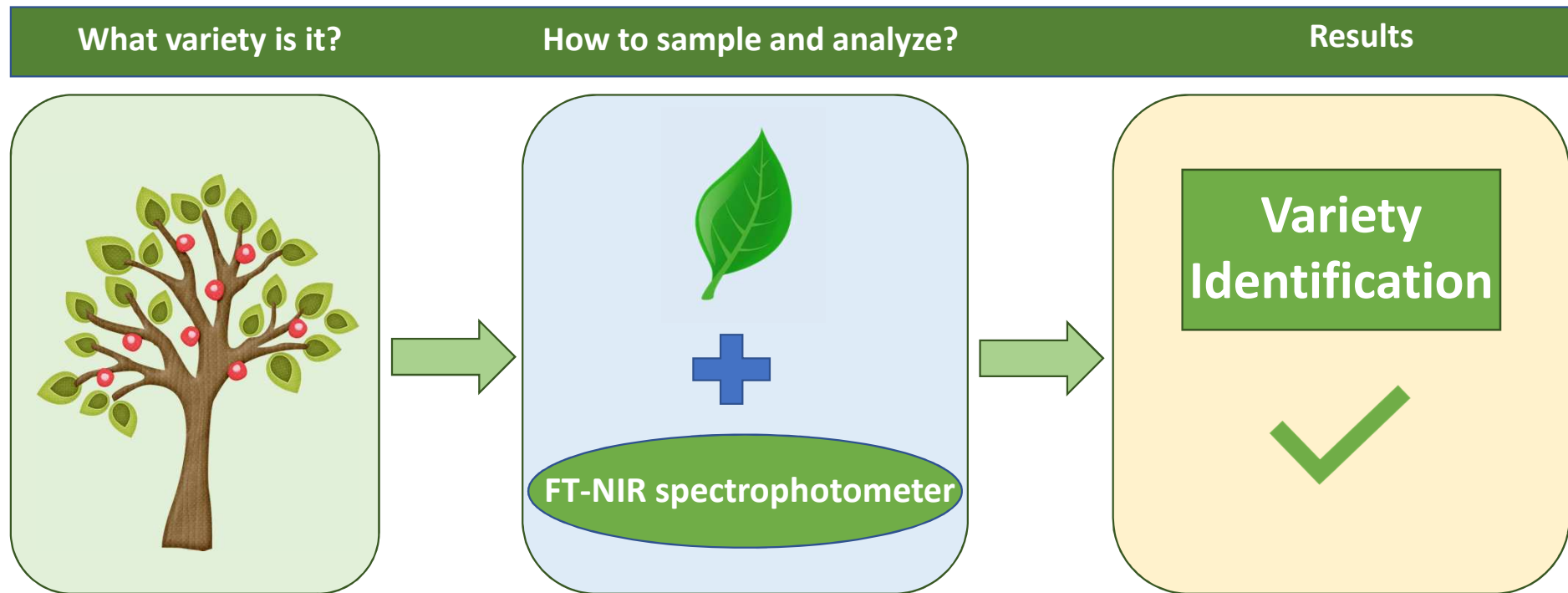
Received Date: 22 March 2019

Revised Date: 21 May 2019

Accepted Date: 27 May 2019

Please cite this article as: S. Borraz-Martínez, R. Boqué, J. Simó, Marià. Mestre, A. Gras, Development of a methodology to analyze leaves from *Prunus dulcis* varieties using near infrared spectroscopy, *Talanta* (2019), doi: <https://doi.org/10.1016/j.talanta.2019.05.105>.

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



1 **DEVELOPMENT OF A METHODOLOGY TO ANALYZE LEAVES FROM *PRUNUS***
2 ***DULCIS* VARIETIES USING NEAR INFRARED SPECTROSCOPY**

3
4 **Sergio Borraz-Martínez^{a,c}, Ricard Boqué^b, Joan Simó^{a,d}, Mariàngela Mestre^c, Anna Gras^a**

5
6 ^a *Universitat Politècnica de Catalunya, Department of Agri-Food Engineering and Biotechnology,*
7 *Esteve Terrades 8, 08860 Castelldefels, Spain*

8 ^b *Universitat Rovira i Virgili, Department of Analytical Chemistry and Organic Chemistry, Campus*
9 *Sescelades, 43007 Tarragona, Spain*

10 ^c *Agromillora Iberia S.L.U, Center of Initial Materials, Ctra. BV-2247 km. 3, 08770-Sant Sadurní*
11 *d'Anoia*

12 ^d *Fundació Miquel Agustí, Esteve Terrades 8, 08860 Castelldefels, Spain*

13
14 Corresponding author: S. Borraz-Martínez: sergio.borraz@upc.edu

15
16 **ORCID numbers of the authors:** Sergio Borraz-Martínez: 0000-0002-5607-9462; Ricard Boqué:
17 0000-0001-7311-4824; Joan Simó: 0000-0002-2853-3381; Anna Gras: 0000-0003-0111-7246

18 **Abstract**

19 Near-infrared spectroscopy (NIRS) can be a faster and more economical alternative to traditional
20 methods for screening varietal mixtures of nursery plants during the propagation process to ensure
21 varietal purity and to avoid errors in the dispatch batches. The global objective of this work was to
22 develop and optimize a NIR spectral collection method for construction of robust multivariate
23 discrimination models. Three different varieties of *Prunus dulcis* (*Avijor*, *Guara*, and *Pentacebas*)
24 of agricultural interest were used for this study. Sources of variation were investigated, including
25 the position of the leaves on the trees, differences among trees of the same variety, and differences
26 at the varietal level. Three types of processed samples were investigated. Fresh leaves, dried leaves,
27 and dried leaves in powder form were included in each analysis. A study of spectral pre-treatment
28 methods was also performed, and multivariate methods were applied to analyze the influence of
29 different factors on classification. These included principal component analysis (PCA), partial least
30 squares discriminant analysis (PLS-DA), and ANOVA simultaneous component analysis (ASCA).
31 The results indicated that variety was the most important factor for classification. The spectral pre-
32 treatment that provided the best results was a combination of standard normal variate (SNV),
33 Savitzky-Golay first derivative, and mean-centering methods. With regard to the type of processed
34 sample, the highest percentages of correct classifications were obtained with fresh and dried
35 powdered leaves at both the training set and test set validation levels. This study represents the first
36 step towards the consolidation of NIRS as a method to identify *Prunus dulcis* varieties.

37 **Keywords:** Optimization; Almond trees; Leaf analysis; Varietal purity; NIR; PLS-DA.

38 1. INTRODUCTION

39 Rapid discrimination between vegetal varieties is a key requirement for all nursery plant production.
40 The huge diversity of vegetal materials necessitates the incorporation of new control systems along
41 the nursery plant production chain to avoid mixing varieties and to ensure varietal purity in the
42 dispatch batches.

43 Nowadays, the most extensively used methods for varietal identification are based on DNA
44 analysis. These techniques include DNA amplification by the polymerase chain reaction (PCR) [1]
45 followed by analysis of genetic variations, such as single nucleotide polymorphisms (SNPs) [2].
46 However, these biomolecular techniques are very expensive for routine analysis of a large number
47 of samples. In this context, the use of spectroscopic analysis combined with chemometrics has
48 recently increased. This combination comprises a rapid, accurate, and nondestructive methodology
49 for the classification and authentication of agricultural products [3].

50 Near-infrared (NIR) spectroscopy has proved to be a powerful analytical tool and has been widely
51 used in various sectors, including the petrochemical [4] and pharmaceutical industries [5]. It has
52 also become a well-established technique for the quantitative and qualitative analysis of agricultural
53 products [6]. Several recent studies have employed spectroscopic techniques for species
54 discrimination [7,8], or differentiation of varieties within a species, such as tomato [9], rice [10] and
55 lettuce [11]. For these reasons, NIRS can be considered a potential candidate for the differentiation
56 of *Prunus dulcis* varieties.

57 Despite recent studies, there is a lack of knowledge regarding the best methodology for accurate
58 sampling of leaves. Most of the published works on species discrimination do not consider factors
59 derived from the nature of the samples, which are potential sources of variance. For example,
60 mature trees have a heterogeneous canopy composed of leaves in different phenological stages. It is
61 thus important to take the sampling procedure into account, especially when an analysis is
62 performed with whole leaves. Improper sampling may generate invalid data, the use of which could
63 lead to incorrect conclusions [12]. To perform correct sampling, it is important to recognize sources
64 of variation and to control for factors from which variation originates. Therefore, it is necessary to
65 first develop a sampling protocol and to select the best material for use.

66 Another analytically relevant aspect is the study of sample processing methods, which may
67 considerably alter the vibrational spectrum of a sample compared to that collected with the sample
68 in its native state. Due to economic and time constraints, it is generally best to avoid any type of
69 sample processing. Moreover, modifying the native architecture of biological tissues can result in

70 the loss of information. Thus, performing analyses in vivo is preferred whenever possible [13].
71 Occasionally, however, sample processing is an indispensable step. In any case, the option that best
72 accomplishes the objective of the study must be selected.

73 The aim of this work is to determine how sampling of vegetal material affects the collection of NIR
74 spectra for the construction of a multivariate discriminant model for *Prunus dulcis* varietal
75 classification. The specific objectives are to 1) determine whether there are differences among the
76 analyzed regions of the leaves or between their upper and lower surfaces; 2) to determine whether
77 differences exist due to the age of the leaves; 3) identify the best sampling procedure for varietal
78 discrimination of almond trees; and 4) study pre-treatment of spectral data and to identify the pre-
79 treatment that leads to the best classification model.

80 2. MATERIAL AND METHODS

81 2.1 Experimental design

82 2.1.1 Assay one

83 The first assay was performed to obtain information about the analyzed regions of fresh and dried
84 leaves. Specifically, the NIR spectra differences resulting from including or excluding the primary
85 veins of the leaves were examined (**Error! Reference source not found.**) together with analysis of
86 the upper (adaxial) and lower (abaxial) leaf surfaces. Twenty samples from two varieties of almond
87 trees, *Guara* and *Pentacebas*, were used for each experiment (Table 1). Results were evaluated by
88 using PCA and PLS-DA models. This assay focused on aspects that affected only fresh and dried
89 leaves. The information obtained in this assay was used for the development of the next two assays.

90 **[Insert Fig.1]**

91 2.1.2 Assay two

92 The second assay was designed to study the NIR spectra differences between young leaves and
93 adult leaves and among samples from different trees of the same variety. The assay was performed
94 on the *Guara* and *Pentacebas* varieties and on a third variety, *Avijor*. Four trees per variety were
95 sampled, twelve in total. Twenty leaves were collected from each tree. Ten of the leaves were
96 collected from the upper part of the branch (apex), which corresponded to young leaves, while the
97 other ten were adult leaves that were collected from the lower part of the branch. Two hundred forty
98 leaves were sampled in total (Table 1). Results were evaluated by using PCA and ASCA-ANOVA
99 models.

100 2.1.3 Assay three

101 Three different leaf processing methodologies were studied in the third assay; one for fresh leaves,
102 one for dried leaves, and the other for dried powdered leaves. The aim was to determine whether the
103 water content and macrostructures of the leaves had any influence on discrimination results. It is
104 important to note that every sample was processed with each of the three methods in order to
105 increase the comparative robustness. We also identified the most suitable pre-treatment method for
106 NIR spectral analysis. The applicability of NIRS for discriminating between *Prunus dulcis* varieties
107 was evaluated by mean partial least squares discriminant analysis (PLS-DA). The available data
108 were randomly divided into calibration (70%) and validation (30%) sets, but both sets contained the
109 same proportion of each variety to prevent unbalanced representation of the almond tree classes. To
110 improve the robustness of comparing results from the three sample processing methods, the same
111 samples included in the three sample processing datasets were used for both cross validation and
112 test set validation. All of the samples used for assay two were also used for this assay (Table 1).

113 **[Insert Table 1]**

114 2.2 Description of the sampling field

115 Vegetal material used in this study came from almond trees located at the mother plant field from
116 the Center of Initial Materials of Agromillora Iberia, S.L.U. in Sant Sadurní d'Anoia (Catalonia,
117 Spain). These trees are under a strict control in order to prevent the appearance of diseases and to
118 ensure the sanitary quality of nursery plants. The use of molecular biology techniques to assess the
119 traceability of the varieties was not necessary in this case because the almond trees were previously
120 certificated by the company.

121 The samples were stored in a plastic bag after collection, assigned identifiers, and stored at 4 °C
122 until analysis.

123 2.3 Sample pre-processing

124 Samples were analyzed either as fresh leaves without processing, as dried leaves, or as dried
125 powdered leaves. To obtain dried leaves, fresh leaves were heated in an oven at 65 °C for 48 hours.
126 A weight was placed on the leaves to keep them flat and to facilitate their posterior analysis. Once
127 dried, the leaves were pulverized to a homogeneous powder with a grinder. Once samples were
128 dried, they were stored in a desiccator with silica gel to prevent any influence from moisture. Only
129 one leaf was used per experiment. Each sample was analyzed in the three ways. First, they were

130 analyzed in fresh, second in dried and finally in powdered. In all the experiments each sample was
131 composed of one leaf only.

132 *2.4 Acquisition of NIR spectra*

133 Samples were scanned in reflectance mode using an Antaris II FT-NIR analyzer (Thermo Scientific,
134 USA) equipped with an integrating sphere module. Samples were measured in the spectral range of
135 12000–3800 cm^{-1} (833–2630 nm). For each spectrum, 32 scans were averaged with a resolution of 4
136 cm^{-1} . Each sample was analyzed in triplicate. Fresh leaves and dried leaves were placed directly
137 over the sphere and covered to prevent interference from environmental light. The powdered leaf
138 samples were measured in a standard sample cup that came with the instrument. A background
139 spectrum was collected every 20 minutes. All spectra were recorded as $\log(1/R)$, where R was the
140 reflectance. Room temperature was maintained at ~ 25 °C, and the humidity remained constant
141 throughout the spectral acquisition process.

142 *2.5 Spectral data pre-treatment*

143 This was an important step, because although different pre-treatments have been reported on
144 extensively [14–16], there is still no clear consensus regarding the best pre-treatment or a guideline
145 to follow. As can be seen in **Error! Reference source not found.**, the spectra contained very little
146 noise. The raw spectra had to be corrected for additive and multiplicative effects that were probably
147 due to light scattering.

148 **[Insert Fig. 2]**

149 A basic pre-treatment was performed in assays one and two, which consisted of the standard normal
150 variate (SNV) method with mean centering. In the assay three, four different pre-treatments were
151 applied and compared to identify the combination that provided the best results in the PLS-DA
152 model. The combinations used were: SNV method with mean centering; SNV method with
153 Savitzky-Golay (SG) first derivative and mean centering; and finally, SNV method followed by de-
154 trending and mean centering. Spectral pre-treatments were performed using PLS_Toolbox
155 (Eigenvector Research Incorporated, Manson, WA) with MATLAB R2017b (MathWorks, Natick,
156 MA).

157 SNV is a normalization procedure for spectral light scattering correction. It is used to correct
158 additive and multiplicative effects in the spectra due to particle size variation. SNV calculates the
159 standard deviation of all the variables in a given sample spectrum. The entire data set is then
160 normalized by this value, which yields a unit standard deviation ($s = 1$) for the sample spectrum

161 [17]. De-trending is sometimes used to remove constant, linear, or curved offsets and is often used
162 in conjunction with SNV. With this method, the mean value or linear trend is subtracted from a
163 vector or matrix. To achieve this, a polynomial of a given order is fitted to the entire data set, and
164 the polynomial is simply subtracted. This algorithm fits all points in the baseline and the signal.
165 [17]. SG first derivative was applied to remove baseline drift and to enhance small spectral
166 differences. The SG derivative method includes a smoothing step, the Savitzky-Golay algorithm,
167 which corrects for the increased noise due to application of the derivative. The SG derivatization
168 algorithm requires selection of the filter width, which is the size of the window, the order of the
169 polynomial, and the order of the derivative [18]. In this work, we selected a 15-point window and
170 applied a second order polynomial. Mean centering is one of the most common pre-processing
171 methods, in which the mean value of each column is calculated and subtracted from each individual
172 value in the column. After mean centering, the mean of each column equals zero, and each row of
173 mean-centered data reflects only how it differs from the average sample in the original data matrix
174 [16].

175 *2.6 Principal component analysis (PCA)*

176 PCA captures the largest amount of variance in the data and reduces the dimensionality of the
177 original dataset through calculation of a new set of variables called principal components (PCs).
178 The PCs are linear combinations of the original variables. Samples and variables are projected onto
179 the new PCs in the calculated PCA space. Samples are defined by their scores, and variables are
180 defined by their loadings. Inspection of the scores and loading plots can lead to a better
181 understanding of the different sources of variation in the data. As a data reduction technique, PCA
182 is frequently the first step in the analysis of a high-dimensional data set. It can then be followed by
183 classification, clustering, or other multivariate techniques [19].

184 *2.7 Partial Least Squares Discriminant Analysis (PLS-DA)*

185 PLS-DA is a classification technique widely used in research studies concerning both varietal
186 classification and authentication of geographical origin [10,20]. PLS-DA is based on the PLS
187 regression algorithm, which searches for linear combinations of the original variables (latent
188 variables) that display maximum covariance with the Y-variables (classes). A discriminator, or
189 threshold, is created that separates the different classes [21]. This technique allows determination of
190 whether or not a given sample belongs in a specific predefined class [22]. The optimal number of
191 factors or latent variables (LVs) for the PLS-DA models was estimated with a cross-validation
192 procedure, and the number yielding the minimum classification error was selected. Venetian blinds

193 cross validation was used for the calibration with a data split of 10 and one sample per blind
194 (thickness).

195 2.8 ASCA-ANOVA

196 Designed experiments with a single dependent variable are typically analyzed with ANOVA [23].
197 Problems occur when hundreds or thousands of variables are measured simultaneously, which is the
198 case in spectroscopic analysis. ANOVA is thus not useful for analyzing multivariate data.
199 Multivariate ANOVA (MANOVA) [24], the natural multivariate extension of ANOVA, breaks
200 down when the number of measurements is smaller than the number of variables [25].

201 ANOVA-simultaneous component analysis (ASCA) [26] is a method used to determine which
202 factors in a fixed-effect experimental design are significant relative to the residual error. ASCA
203 allows an ANOVA-like analysis, even when there are more variables than samples. Two matrices
204 are used to perform the procedure. The X-matrix contains the experimental data, while the F-matrix
205 represents the experimental design. PCA of each factor in the effect (X) matrix reduces the number
206 of variables to a smaller number of principal components. In this way, the parameter estimation
207 functionality of ANOVA is merged with PCA, and the presence of more variables than samples is
208 no longer problematic [27]. Due to the hierarchy of factors analyzed in the present study, a nested
209 design referred to as multi-level simultaneous component analysis (MLSCA) [28] was applied.
210 Hence, the leaf age factor was nested within the tree factor, which in turn was nested within the
211 variety factor.

212 3. RESULTS AND DISCUSSION

213 3.1 Assay one

214 3.1.1 Comparison of leaf midvein and lamina

215 Whether differences exist within the same leaf is a question that frequently arises. For this reason,
216 spectra were collected in different areas of healthy leaves. The two regions of the leaves used for
217 comparison are shown in Fig. 1. PCA was performed with two of the almond tree varieties, *Guara*
218 and *Pentacebas*, to identify possible differences between the measurement areas on fresh and dried
219 samples. These results are shown in Fig. 2.

220

221 **[Insert Fig. 3]**

222

223 Differences when including or not the primary vein were detected. The data clouds with and
224 without midvein form separate clusters in both kinds of pre-processed samples. This cluster

225 separation can be observed in both varieties, although the separation is clearer for the *Pentacebas*
226 variety. Differences were detected whether or not the primary vein was included. The data clouds
227 with and without the midvein formed separate clusters for both processed sample types. This cluster
228 separation was observed with both varieties, although the separation was more pronounced in the
229 results from the *Pentacebas* variety. Considering the macrostructures and compositions of the
230 analysis regions were not equivalent, which was reflected in their spectral signatures, these
231 differences were justifiable. When the primary vein was scanned, the reflectance spectra of both the
232 primary vein and the laminar regions located on either side of the primary vein were collected.
233 Taking into account that secondary veins were present in the laminar regions, identifying
234 differences between these regions indicated the primary vein had a profound influence on the
235 spectra.

236 The apical region and a region adjacent to the leaf margin showed more damage and decay than the
237 central region of the leaves. Consequently, the central region was usually more stable. The leaf size
238 could make it difficult to completely exclude the primary vein during measurement of the laminar
239 region. Collecting spectra in the central region, including the primary vein, could therefore provide
240 a standardized measure.

241

242 *3.1.2 Comparison of adaxial and abaxial surfaces*

243 Differences between the upper and lower surfaces of the leaves were also investigated. These
244 results are shown in **Error! Reference source not found.** In both fresh and dried samples, results
245 of PCA revealed differences between the spectra obtained from the upper and lower leaf surfaces.
246 However, this difference was not as clear in fresh leaves of the *Pentacebas* variety. The upper and
247 lower surfaces of leaves in all plants are different. In addition, the stomas are usually present on the
248 abaxial surface together with trichomes and others surface features. The differences between these
249 two surfaces could be the cause for separation of their spectra in the PCA plots.

250

251 **[Insert Fig. 4]**

252

253 A PLS-DA model was built to determine which surfaces were most suitable for discriminating
254 between two almond tree varieties using fresh or dried samples. The classification results are shown
255 in Table 2. The PLS-DA model had a classification score of 100% for both types of processed
256 samples when the upper leaf surface was analyzed. Perfect discrimination was obtained using the
257 lower leaf surface as well. Based on these results, the differences identified by PCA did not affect
258 the discrimination results with either surface.

259

260 **[Insert Table 2]**

261

262 *3.2 Assay two*263 *3.2.1 Variability between trees of the same variety*

264

265 Differences among trees of the same variety are important to consider when building a classification
266 model. This source of variation determines the number of trees of each variety that must be sampled
267 for development of the final model. If the variance is very large, it could affect the model's
268 discrimination capability. The PCA results from assay 2 are shown in Fig 4.

269

270 **[Insert Fig. 5]**

271

272 No differences were identified among the four trees studied within each variety. This was the case
273 for fresh, dried, and dried powdered leaves. This was remarkable, because if significant differences
274 were found, it would have been more difficult to build a good classification model. Also noteworthy
275 was that the same results were obtained with samples processed with the three different methods,
276 and with samples of different varieties. Such similar behavior in all cases is a positive indicator
277 when creating a classification model. A more exhaustive study of the variability between trees was
278 performed using the ASCA-ANOVA method, which is discussed in section 3.2.3.

279

280 *3.2.2 Variability between leaves of the same variety*

281 Since differences among almond trees of the same variety were not detected at the PCA level, we
282 decided to include all samples of the same variety in a single PCA model. This made it easier to
283 study the variability among samples within each variety while increasing the robustness of the
284 model with more samples. The results of PCA modelling are shown in **Error! Reference source**
285 **not found..**

286 **[Insert Fig. 6]**

287 Two clusters could be distinguished using only the first two principal components. This separation
288 was very clear in some cases, such as the dried processed samples of the *Pentacebas* variety, for
289 which the two clusters were completely separated (Fig. 6f). The results of all of the PCA models
290 were similar, regardless of the sample processing method or the variety studied. However, overlap
291 between the two data clusters was observed in some cases, such as dried samples of the *Avijor*

292 variety (Fig. 6d). The overlap could be explained by the presence of leaves in a phenological
293 stadium intermediate between young and adult. It was possible to observe the progressive growth of
294 the leaves, although this was not the goal of the assay. In any case, the results indicated there were
295 differences between young and adult leaves at the spectral level. This difference should be
296 considered at the time of sampling.

297

298 3.2.3 ASCA-ANOVA analysis

299 To study variability between *Prunus dulcis* varieties more deeply, an ASCA-ANOVA model was
300 constructed for young and adult leaves from trees of the same variety. The modelling results are
301 shown in **Error! Not a valid bookmark self-reference.** The raw spectra pre-treatment used to
302 develop the model, SNV with mean centering, was the same as that used for the PCA models.

303

304 [Insert Table 3]

305

306 Tree variety was the most influential factor for variance among fresh and dried powdered leaves
307 and accounted for 30.26% and 24.99%, respectively, of the total effect in these samples. Despite
308 explaining 19.25% of the effect for dried leaves, tree variety was not the factor that accounted for
309 the majority of variance. For two of the three processing methods, the variety factor had the greatest
310 effect, which indicated that differences between varieties were important. The tree factor explained
311 little of the variance for the three processing methods, which was in agreement with the PCA results
312 shown in Fig 4. This indicated strong homogeneity between trees of the same variety, an aspect that
313 could be key for effective discrimination between varieties. For fresh and dried powdered leaves,
314 the age (young/adult) factor explained a higher percentage of variance than the tree factor, but it
315 accounted for less of the variance than tree variety. In dried powdered leaves, the difference
316 between the age and tree factors was not large. The age factor accounted for 6.68% of the variance,
317 while the tree factor explained 1.87%. The difference was more notable for fresh leaves, as the age
318 factor accounted for 19.11% of the explained variance. The age factor was most significant for
319 dried leaves, accounting for 24.18% of the explained variance. Therefore, the age factor had a
320 greater effect in non-powder samples. These results also correlated with the results of the PCA
321 (**Error! Reference source not found.**), in which differences due to leaf age were observed, but
322 overlap of the cluster regions was detected.

323

324 All of the variance not explained by the studied factors accumulated in the residual term. In the
325 three types of processed samples, the residual accounted for a high percentage of the variance.

326 Fresh leaves had a lower residual than either the dried or dried powdered leaves. It was thought that
327 the main source of uncontrolled variance was the physiological state of the leaves, which included
328 damage to the leaves and climatologic agents. The combination of these abiotic factors with biotic
329 factors influences plant physiology [29,30]. It is important to note that the leaves used in this study
330 came from trees located in an outdoor field.

331

332 *3.3 Assay three*

333 *3.3.1 Spectral pre-treatment study*

334 **[Insert Table 4]**

335

336 **Error! Reference source not found.** shows the results of the PLS-DA modelling using different
337 spectral pre-treatments. The best classification results for the three types of samples were obtained
338 with the SNV pre-treatment and application of the SG first derivative and mean centering. This was
339 curious, because although modelling was performed for one material (almond tree leaves), the
340 samples analyzed were completely different in terms of their macrostructures and dry compositions.
341 With this spectral pre-treatment, 100% classification accuracy was achieved for at least one variety
342 with each sample processing method. Results were even more remarkable with dried powdered
343 leaves, for which 100% accuracy was attained in the test set validation for all three varieties. The
344 lowest accuracy obtained with this spectral pre-treatment was 97.5% at both the cross-validation
345 and test set validation levels. No relevant differences between the other two spectral pre-treatments
346 were observed, so de-trending did not appear to have a significant effect. It is important to note that
347 in the case of fresh leaves, similar results were obtained with the three different spectral pre-
348 treatments.

349

350 *3.3.2 Sample processing study*

351 Each sample processing method had its advantages and disadvantages. Fresh leaves did not require
352 any processing, so measurement was faster and easier than it was with the other types of samples.
353 However, the water content of the leaves was a disadvantage, because it generated wide bands in
354 the NIR spectra. This could make discrimination between varieties more difficult. Samples can be
355 dehydrated to circumvent the effects of water, but this process is time-consuming (48 h), so it is not
356 the best option if rapid identification is required.

357 To evaluate which of the processed samples was the most suitable for varietal classification, the
358 advantages and disadvantages of each were considered together with the PLS-DA classification

359 results obtained with SNV spectral pre-treatment and application of the SG first derivative and
360 mean centering (Table 4).

361 The results obtained with the three types of sample processing at the calibration level could be
362 considered quite good, although those obtained with fresh leaves were less stellar. The dried
363 powdered leaves provided a higher percentage of correct classifications. For the test set validation,
364 high percentages of correct classifications were obtained with all varieties and processed sample
365 types. The results provided by the dried leaves were not as good as those obtained with the other
366 two processed sample types, although the *Pentacebas* variety was correctly classified in 100% of
367 the test set validations. Fresh leaves provided almost perfect classification, and nearly 100% correct
368 classification was attained with dried powdered leaves. Taking only the PLS-DA results into
369 account, the best sample processing method was drying and powdering the leaves. Considering the
370 methodological aspects, using fresh leaves was the fastest and easiest option. The biggest drawback
371 of fresh leaves was their water content, but this did not seem to hinder discrimination between the
372 varieties studied.

373 In the ASCA-ANOVA model performed in assay two (Table 3), the strongest effect on dried leaves
374 was contributed by the leaf age factor. The age factor accounted for more variability than even the
375 tree variety factor, which could be problematic. Fresh leaves exhibited more favorable behavior in
376 the ASCA-ANOVA model. Results of the ASCA-ANOVA model with dried powdered leaves were
377 similar to those obtained with fresh leaves, but the residual was higher.

378

379 **4. CONCLUSIONS AND PERSPECTIVES**

380

381 In this study, we defined a methodology for construction of a classification model that could
382 discriminate between *Prunus dulcis* varieties using NIRS. We also identified the most important
383 sampling and analysis aspects. In assay one, differences were seen in the PCA whether or not the
384 midvein was included. The central leaf region provided more useful information for discriminating
385 between almond tree varieties, because it contained both the primary vein and the laminar tissues.
386 We also attempted to determine which surface of the leaves, adaxial or abaxial, was the most
387 suitable for analysis. Despite the spectral differences observed, the comparison made using the
388 PLS-DA model indicated this was not an important aspect.

389 In assay two, no notable differences were detected between trees of the same variety, which
390 indicated that trees within each variety were quite homogeneous. Differences were observed at the
391 PCA level between young and adult leaves, which indicated age was important to consider during
392 the sampling process.

393 The best results from the PLS-DA models in assay three were obtained with dried powdered leaves
394 when SNV was used for spectral pre-treatment with application of the SG first derivative (15-point
395 window, second order) and mean centering. However, fresh leaves appeared to be the easiest and
396 most suitable samples for laboratory or industrial analysis. These results indicated that both fresh
397 leaves and dried powdered leaves could be useful for discriminating between *Prunus dulcis*
398 varieties using NIR spectroscopy.

399 All the information gathered in the present study will be used to build a classification model that
400 includes more *Prunus dulcis* varieties. The potential of NIR spectroscopy for the classification of
401 almond tree varieties and its implementation as a quality control tool in the nursery plant industry
402 will be studied.

403 **Acknowledgment**

404 The authors thank Thermo Scientific for temporary assignment of the NIR equipment.

405 **Funding**

406 This work was supported by Generalitat de Catalunya through a grant from Program of Industrial
407 Doctorates (DI-COF 2017) and by The Spanish Ministry of Economy and Competitiveness (project
408 AGL 2015-70106-R).

409 **REFERENCES**

- 410 [1] K. Mullis, F. Faloona, S. Scharf, R. Saiki, G. Horn, H. Erlich, Specific enzymatic
411 amplification of DNA in vitro: the polymerase chain reaction., Cold Spring Harb. Symp.
412 Quant. Biol. 51 Pt 1 (1986) 263–73. doi:10.1101/SQB.1986.051.01.032.
- 413 [2] P.K. Gupta, J.K. Roy, M. Prasad, Single nucleotide polymorphism a new paradigm for
414 molecular marker technology and DNA polymorphism detection with emphasis on their use
415 in plant.pdf, Curr. Sci. 80 (2001) 524–535. doi:10.2307/24104242.
- 416 [3] M. Makky, P. Soni, In situ quality assessment of intact oil palm fresh fruit bunches using
417 rapid portable non-contact and non-destructive approach, J. Food Eng. 120 (2014) 248–259.
418 doi:10.1016/J.JFOODENG.2013.08.011.
- 419 [4] M. V. Reboucas, J.B. dos Santos, D. Domingos, A.R.C.G. Massa, Near-infrared
420 spectroscopic prediction of chemical composition of a series of petrochemical process
421 streams for aromatics production, Vib. Spectrosc. 52 (2010) 97–102.
422 doi:10.1016/J.VIBSPEC.2009.09.006.
- 423 [5] M. Verstraeten, D. Van Hauwermeiren, M. Hellings, E. Hermans, J. Geens, C. Vervae, I.
424 Nopens, T. De Beer, Model-based NIR spectroscopy implementation for in-line assay
425 monitoring during a pharmaceutical suspension manufacturing process, Int. J. Pharm. 546
426 (2018) 247–254. doi:10.1016/j.ijpharm.2018.05.043.
- 427 [6] Z. Seregély, T. Deák, G.D. Bisztray, Distinguishing melon genotypes using NIR

- 428 spectroscopy, *Chemom. Intell. Lab. Syst.* 72 (2004) 195–203.
429 doi:10.1016/j.chemolab.2004.01.013.
- 430 [7] Q. Fan, Y. Wang, P. Sun, S. Liu, Y. Li, Discrimination of Ephedra plants with diffuse
431 reflectance FT-NIRS and multivariate analysis, *Talanta*. (2010).
432 doi:10.1016/j.talanta.2009.09.018.
- 433 [8] F.M. Durgante, N. Higuchi, A. Almeida, A. Vicentini, Species Spectral Signature:
434 Discriminating closely related plant species in the Amazon with Near-Infrared Leaf-
435 Spectroscopy, *For. Ecol. Manage.* 291 (2013) 240–248. doi:10.1016/j.foreco.2012.10.045.
- 436 [9] Y.N. Shao, C.Q. Xie, L.J. Jiang, J.H. Shi, J.J. Zhu, Y. He, Discrimination of tomatoes bred
437 by spaceflight mutagenesis using visible/near infrared spectroscopy and chemometrics,
438 *Spectrochim. Acta - Part A Mol. Biomol. Spectrosc.* (2015). doi:10.1016/j.saa.2015.01.018.
- 439 [10] L. Zhang, S.S. Wang, Y.F. Wang, J.R. Pan, C. Zhu, Discrimination of Transgenic Rice
440 Based on Near Infrared Reflectance Spectroscopy and Partial Least Squares Regression
441 Discriminant Analysis, *Rice Sci.* (2015). doi:10.1016/j.rsci.2015.09.004.
- 442 [11] L. de Oliveira Moura, D. de Carvalho Lopes, A.J. Steidle Neto, L. de Castro Louback
443 Ferraz, L. de Almeida Carlos, L.M. Martins, Evaluation of Techniques for Automatic
444 Classification of Lettuce Based on Spectral Reflectance, *Food Anal. Methods*. 9 (2016)
445 1799–1806. doi:10.1007/s12161-015-0366-5.
- 446 [12] G.D. Batten, Plant analysis using near infrared reflectance spectroscopy: the potential and
447 the limitations, *Aust. J. Exp. Agric.* 38 (1998) 697–706. doi:10.1071/EA97146.
- 448 [13] P. Skolik, M.R. McAinsh, F.L. Martin, Biospectroscopy for Plant and Crop Science, in: C.S.
449 João Lopes (Ed.), *Compr. Anal. Chem.*, Elsevier, 2018: pp. 15–49.
450 doi:10.1016/BS.COAC.2018.03.001.
- 451 [14] A. Candolfi, R. De Maesschalck, D. Jouan-Rimbaud, P.A. Hailey, D.L. Massart, The
452 influence of data pre-processing in the pattern recognition of excipients near-infrared
453 spectra, *J. Pharm. Biomed. Anal.* 21 (1999) 115–132. doi:10.1080/13554790490495140.
- 454 [15] L. Xu, Y.P. Zhou, L.J. Tang, H.L. Wu, J.H. Jiang, G.L. Shen, R.Q. Yu, Ensemble
455 preprocessing of near-infrared (NIR) spectra for multivariate calibration, *Anal. Chim. Acta.*
456 616 (2008) 138–143. doi:10.1016/j.aca.2008.04.031.
- 457 [16] Å. Rinnan, F. van den Berg, S.B. Engelsen, Review of the most common pre-processing
458 techniques for near-infrared spectra, *TrAC - Trends Anal. Chem.* 28 (2009) 1201–1222.
459 doi:10.1016/j.trac.2009.07.007.
- 460 [17] R.J. Barnes, M.S. Dhanoa, S.J. Lister, Standard normal variate transformation and de-
461 trending of near-infrared diffuse reflectance spectra, *Appl. Spectrosc.* 43 (1989) 772–777.
462 doi:10.1366/0003702894202201.
- 463 [18] A. Savitzky, M.J.E. Golay, Smoothing and Differentiation of Data by Simplified Least
464 Squares Procedures., *Anal. Chem.* 36 (1964) 1627–1639. doi:10.1021/ac60214a047.
- 465 [19] R.S.K. Barry M. Wise, Neal B. Gallagher, Rasmus Bro, Jeremy M. Shaver, Willem Windig,
466 *Chemometrics Tutorial for PLS _ Toolbox and Solo*, 2006. doi:10.1016/j.cplett.2004.08.130.
- 467 [20] P. Wang, Z. Yu, Species authentication and geographical origin discrimination of herbal
468 medicines by near infrared spectroscopy: A review, *J. Pharm. Anal.* 5 (2015) 277–284.
469 doi:10.1016/j.jpha.2015.04.001.

- 470 [21] R.G. Brereton, G.R. Lloyd, Partial least squares discriminant analysis: Taking the magic
471 away, *J. Chemom.* 28 (2014) 213–225. doi:10.1002/cem.2609.
- 472 [22] D. Ballabio, V. Consonni, Classification tools in chemistry. Part 1: linear models. PLS-DA,
473 *Anal. Methods.* 5 (2013) 3790. doi:10.1039/c3ay40582f.
- 474 [23] L. Støhle, S. Wold, Analysis of variance (ANOVA), *Chemom. Intell. Lab. Syst.* 6 (1989)
475 259–272. doi:10.1016/0169-7439(89)80095-4.
- 476 [24] L. Støhle, S. Wold, Multivariate analysis of variance (MANOVA), *Chemom. Intell. Lab.*
477 *Syst.* 9 (1990) 127–141. doi:10.1016/0169-7439(90)80094-M.
- 478 [25] G. Zwanenburg, H.C.J. Hoefsloot, J.A. Westerhuis, J.J. Jansen, A.K. Smilde, ANOVA–
479 principal component analysis and ANOVA–simultaneous component analysis: A
480 comparison, *J. Chemom.* 25 (2011) 561–567. doi:10.1002/cem.1400.
- 481 [26] A.K. Smilde, J.J. Jansen, H.C.J. Hoefsloot, R.J.A.N. Lamers, J. van der Greef, M.E.
482 Timmerman, ANOVA-simultaneous component analysis (ASCA): A new tool for analyzing
483 designed metabolomics data, *Bioinformatics.* 21 (2005) 3043–3048.
484 doi:10.1093/bioinformatics/bti476.
- 485 [27] J.J. Jansen, H.C.J. Hoefsloot, J. Van Der Greef, M.E. Timmerman, J.A. Westerhuis, A.K.
486 Smilde, ASCA: Analysis of multivariate data obtained from an experimental design, *J.*
487 *Chemom.* 19 (2005) 469–481. doi:10.1002/cem.952.
- 488 [28] O.E. de Noord, E.H. Theobald, Multilevel component analysis and multilevel PLS of
489 chemical process data, *J. Chemom.* 19 (2005) 301–307. doi:10.1002/cem.933.
- 490 [29] N. Suzuki, R.M. Rivero, V. Shulaev, E. Blumwald, R. Mittler, Abiotic and biotic stress
491 combinations, *New Phytol.* 203 (2014) 32–43. doi:10.1111/nph.12797.
- 492 [30] J.C.M.S. Moura, C.A.V. Bonine, J. de Oliveira Fernandes Viana, M.C. Dornelas, P.
493 Mazzafera, Abiotic and Biotic Stresses and Changes in the Lignin Content and Composition
494 in Plants, *J. Integr. Plant Biol.* 52 (2010) 360–376. doi:10.1111/j.1744-7909.2010.00892.x.

495

496 **FIGURE CAPTIONS**

497 Fig. 1. Image of an almond leaf showing the two studied regions.

498 Fig. 2. Mean raw spectra from the three processed sample types. Fresh leaf (green dashed line);
499 dried powdered leaf (blue solid line); and dried leaf (red dotted line).

500 Fig. 2. PCA results from the Guara and Pentacebas varieties with and without inclusion of the
501 midvein. The presence of the midvein is indicated by red diamonds, and absence of the midvein is
502 indicated by green squares. a) Dried leaf of the Guara variety; b) dried leaf of the Pentacebas
503 variety; c) fresh leaf of the Guara variety; d) fresh leaf of the Pentacebas variety.

504 Fig. 3. PCA results showing the differences between the adaxial (red diamonds) and abaxial (green
505 squares) leaf surfaces. a) Dried leaf of the Guara variety; b) dried leaf of the Pentacebas variety; c)
506 fresh leaf of the Guara variety; d) fresh leaf of the Pentacebas variety.

507 Fig 4. PCA results from the study of differences between trees of the same variety. Each tree is
508 represented by a different symbol (triangle, circle, diamond, and square). a) Fresh leaf of the Avijor
509 variety; b) fresh leaf of the Guara variety; c) fresh leaf of the Pentacebas variety; d) dried leaf of the
510 Avijor variety; e) dried leaf of the Guara variety; f) dried leaf of the Pentacebas variety; g) dried
511 powdered leaf of the Avijor variety; h) dried powdered leaf of the Guara variety; i) dried powdered
512 leaf of the Pentacebas variety.

513 Fig. 6. PCA results from the study of differences between young (yellow circles) and adult (pink
514 stars) leaves. a) Fresh leaf of *Avijor* variety; b) fresh leaf of *Guara* variety; c) fresh leaf of
515 *Pentacebas* variety; d) dried leaf of *Avijor* variety; e) dried leaf of *Guara* variety; f) dried leaf of
516 *Pentacebas* variety; g) dried powdered leaf of *Avijor* variety; h) dried powdered leaf of *Guara*
517 variety; i) dried-powdered leaf of *Pentacebas* variety.

518

519

520

521

522

523

524

525

526

527

528

529

530

531

532

533

534 TABLES

535 **Table 1.** Summary of the samples used in the study.

	Varieties	Fresh samples	Dried samples	Dried-powdered samples	
		Number of samples			
Assay one	<i>Guara</i>	10	10	not used	536
	<i>Pentacebas</i>	10	10	not used	538
Assay two	<i>Avijor</i>	80	80	80	539
	<i>Guara</i>	80	80	80	
	<i>Pentacebas</i>	80	80	80	540
Assay three	<i>Avijor</i>	80	80	80	541
	<i>Guara</i>	80	80	80	
	<i>Pentacebas</i>	80	80	80	542

543

	Real class	Data set	Fresh samples		Dried samples	
			Assigned class		Assigned class	
			<i>Guara</i>	<i>Pentacebas</i>	<i>Guara</i>	<i>Pentacebas</i>
Adaxial	<i>Guara</i>	Cross-validation	100 %	100 %	100 %	100 %
	<i>Pentacebas</i>		100 %	100 %	100 %	100 %
Abaxial	<i>Guara</i>	Cross-validation	100 %	100 %	100 %	100 %
	<i>Pentacebas</i>		100 %	100 %	100 %	100 %

544 **Table 2.** PLS-DA results from the comparison of adaxial and abaxial leaf surfaces.

545

546 **Table 3.** Results of ASCA-ANOVA modelling to study variance of the factors.

Factor	Fresh leaves		Dried leaves		Dried-powdered leaves	
	Principal components	Effect %	Principal components	Effect %	Principal components	Effect %
Variety	2	30.26	2	19.25	2	24.99
Tree	3	1.83	3	4.69	3	1.87
Young / adult	1	19.11	1	24.18	1	6.68
Residual	6	48.80	3	51.88	3	66.46

547

548

549

550

551

552

553 **Table 4.** PLS-DA model results of the spectra pre-treatment and study of the types of pre-processed
 554 samples.

Dried-powdered leaves				
Real class	Data set	Assigned class		
		<i>SNV + Mean center</i>	<i>SNV + 1st derivative + Mean center</i>	<i>SNV + De-trending + Mean center</i>
<i>Avijor</i>	Cross-validation	87.4 %	99.2 %	86.6 %
	Test set validation	97.5 %	100 %	97.5 %
<i>Guara</i>	Cross-validation	89.9 %	99.2 %	89.1 %
	Test set validation	96.6 %	100 %	96.6 %
<i>Pentacebas</i>	Cross-validation	97.5 %	100 %	97.5 %
	Test set validation	99.2 %	100 %	99.2 %
Dried leaves				
Real class	Data set	Assigned class		
		<i>SNV + Mean center</i>	<i>SNV + 1st derivative + Mean center</i>	<i>SNV + De-trending + Mean center</i>
<i>Avijor</i>	Cross-validation	97.5 %	99.2 %	95.0 %
	Test set validation	95.0 %	98.3 %	93.3 %
<i>Guara</i>	Cross-validation	95.0 %	100 %	93.3 %
	Test set validation	92.5 %	97.5 %	91.7 %
<i>Pentacebas</i>	Cross-validation	97.5 %	99.2 %	93.3 %
	Test set validation	97.5 %	100 %	98.3 %
Fresh leaves				
Real class	Data set	Assigned class		
		<i>SNV + Mean center</i>	<i>SNV + 1st derivative + Mean center</i>	<i>SNV + De-trending + Mean center</i>
<i>Avijor</i>	Cross-validation	100 %	97.5 %	100 %
	Test set	99.2 %	100 %	99.2 %

	validation			
<i>Guara</i>	Cross-validation	99.2 %	97.5 %	99.2 %
	Test set validation	98.3 %	99.2 %	98.3 %
<i>Pentacebas</i>	Cross-validation	99.2 %	100 %	99.2 %
	Test set validation	99.2 %	99.2 %	99.2 %

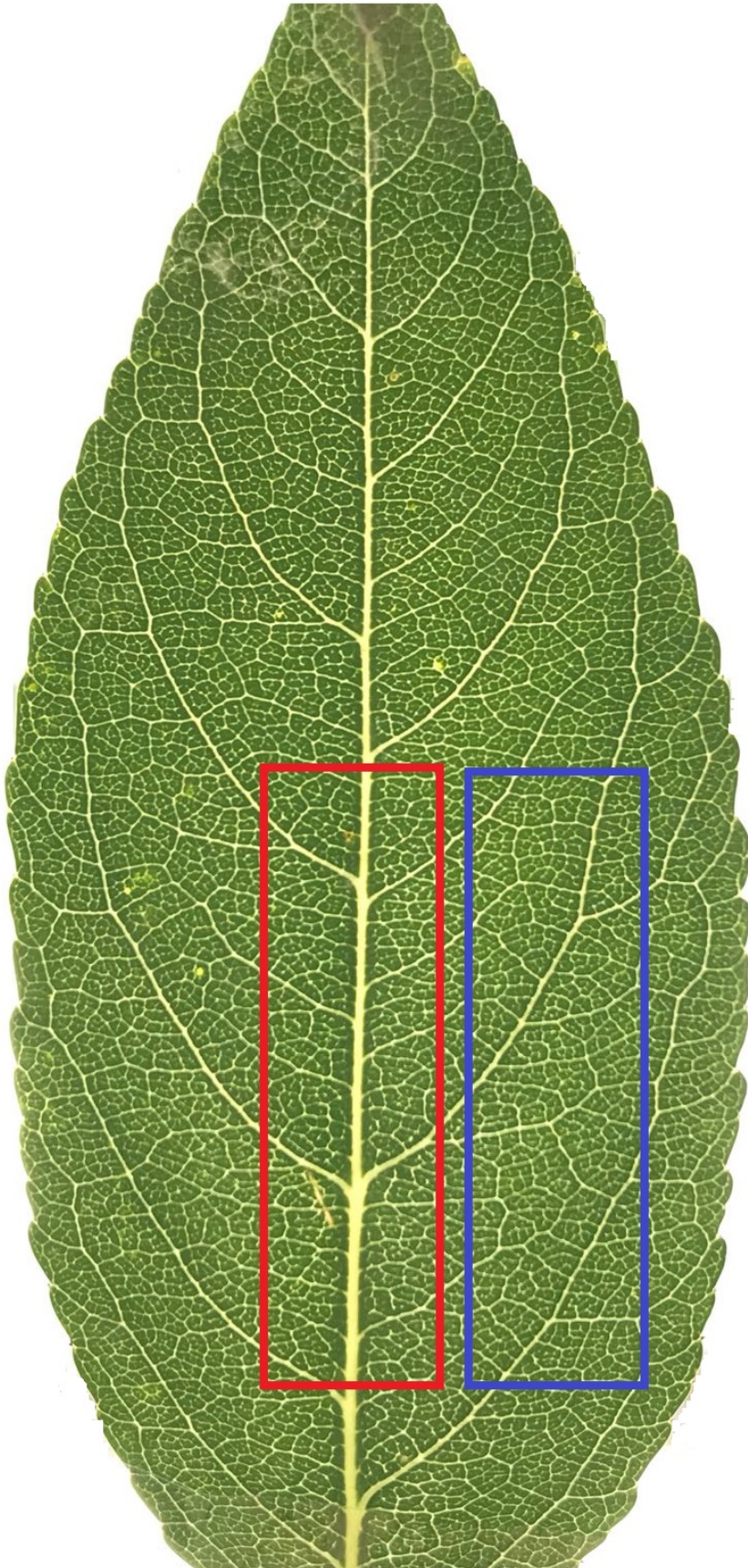
555

556

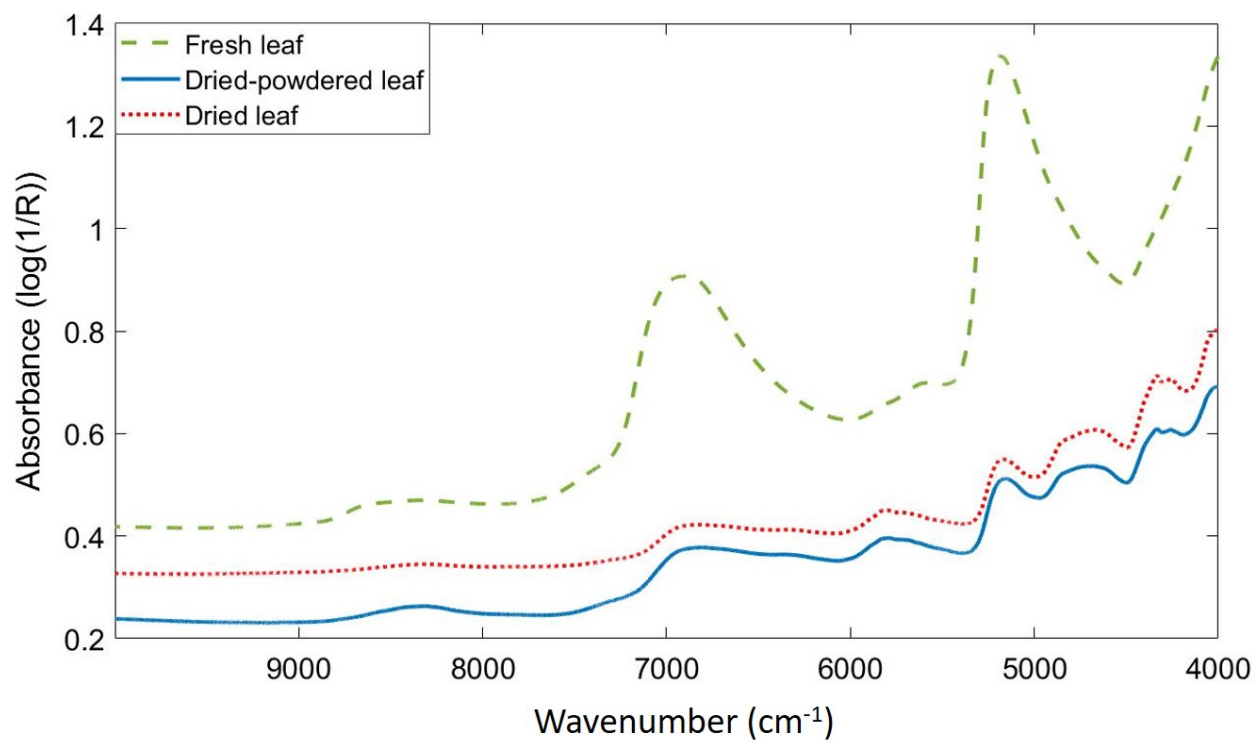
557

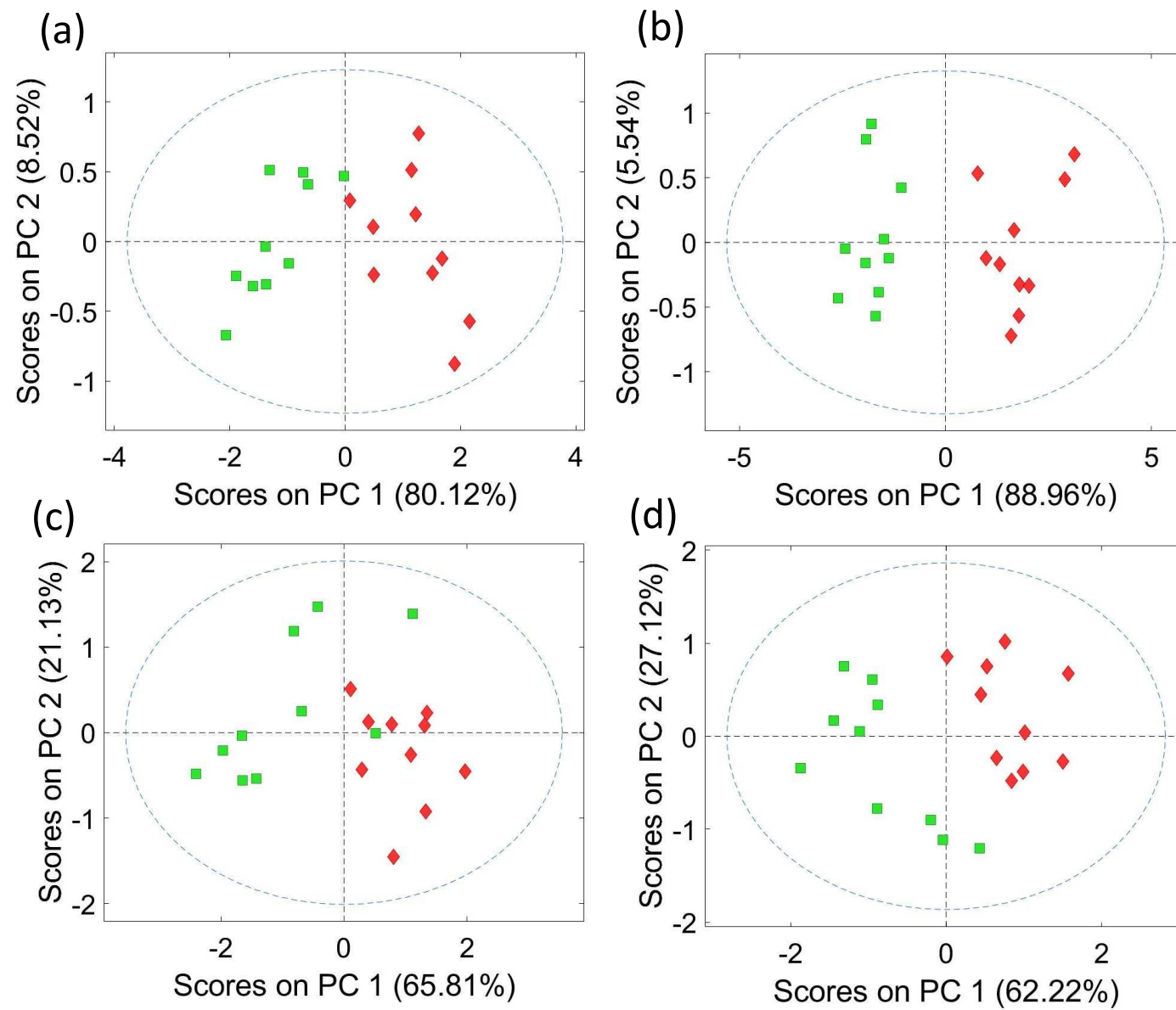
558

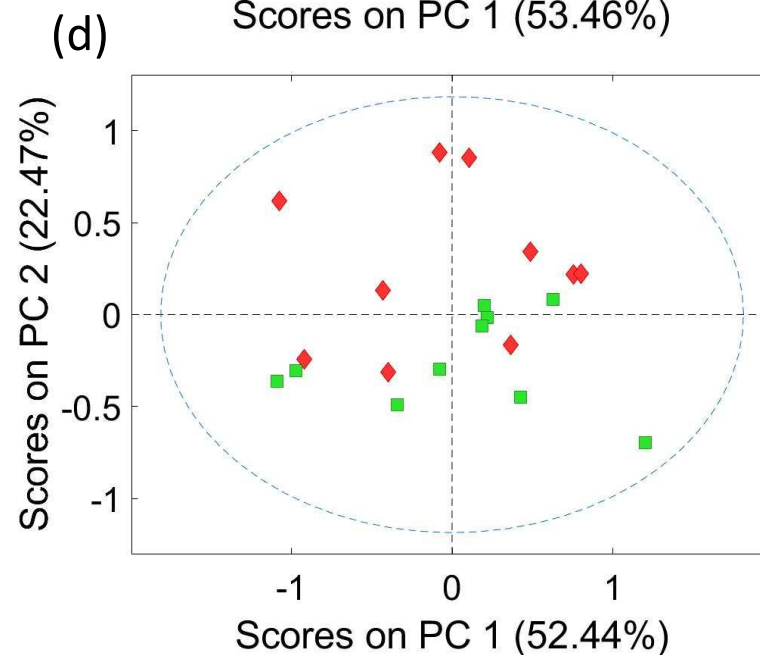
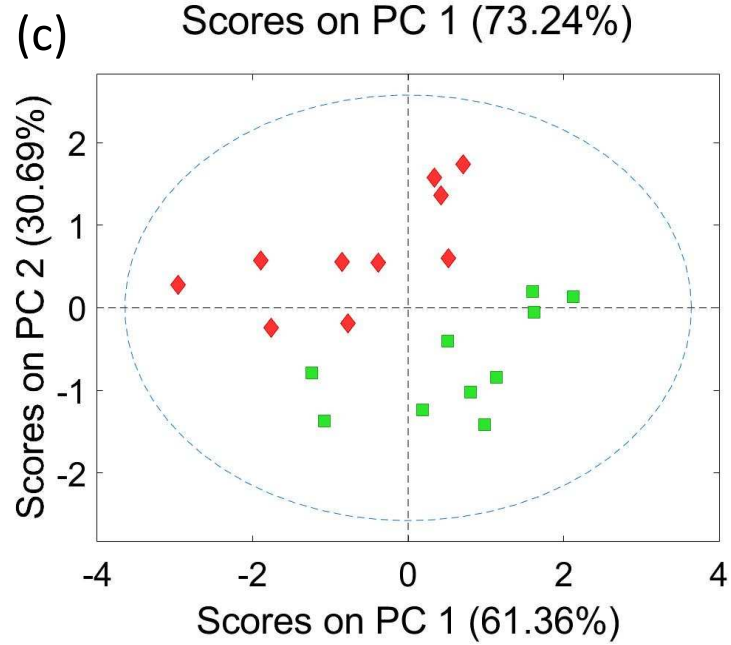
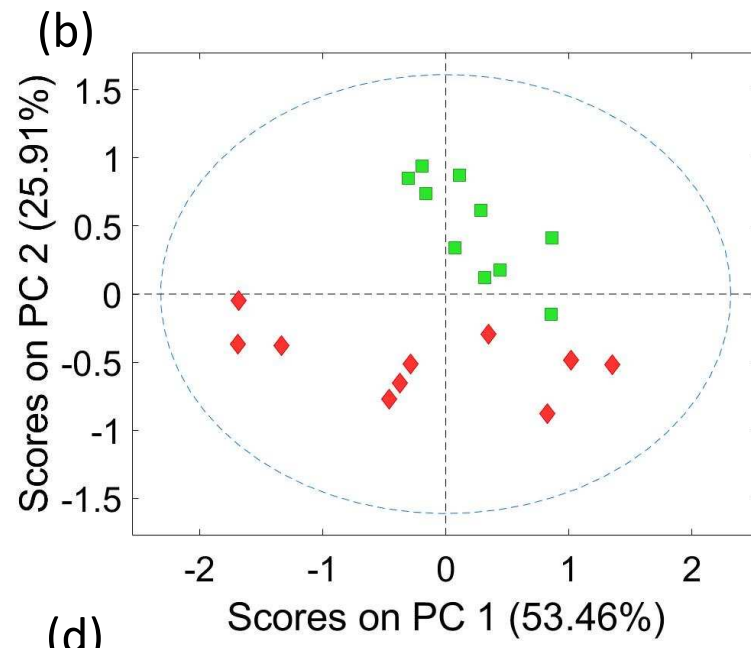
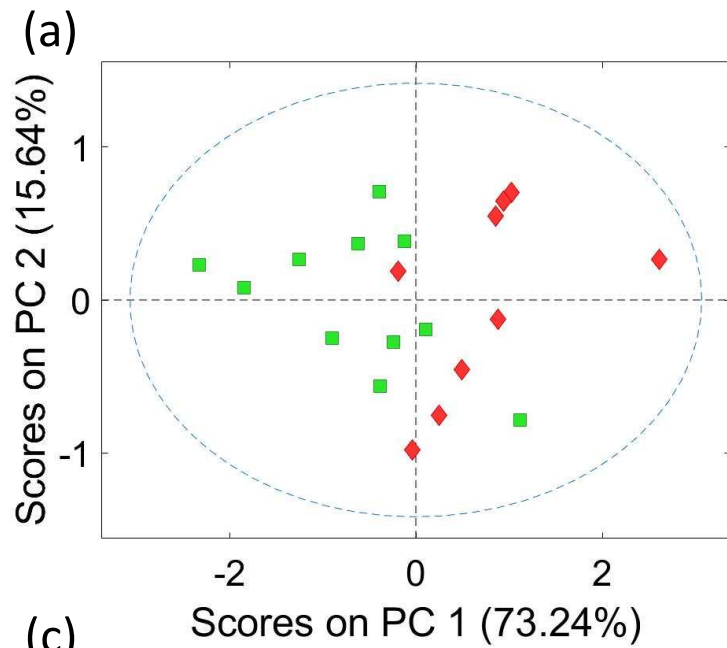
559

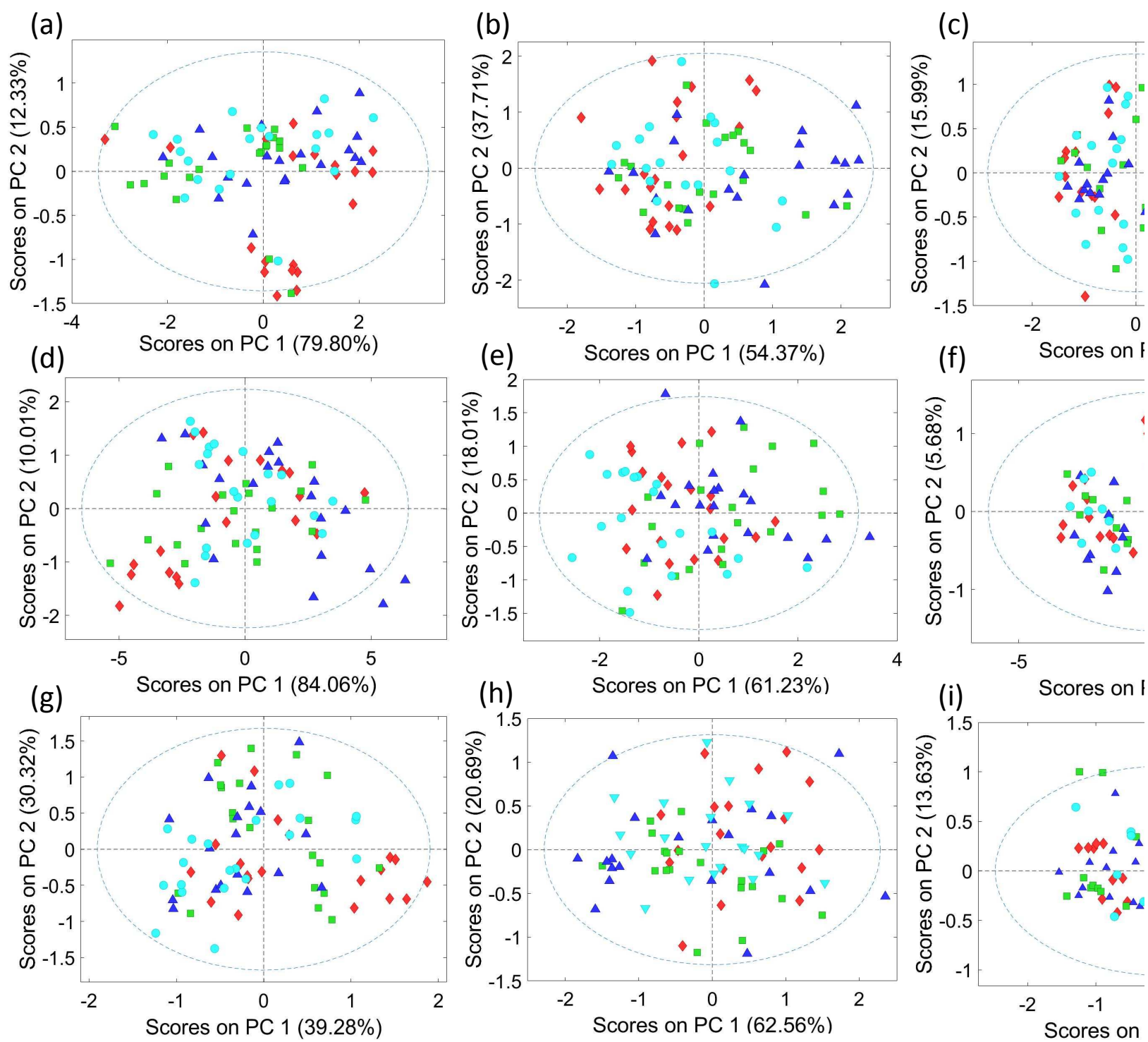


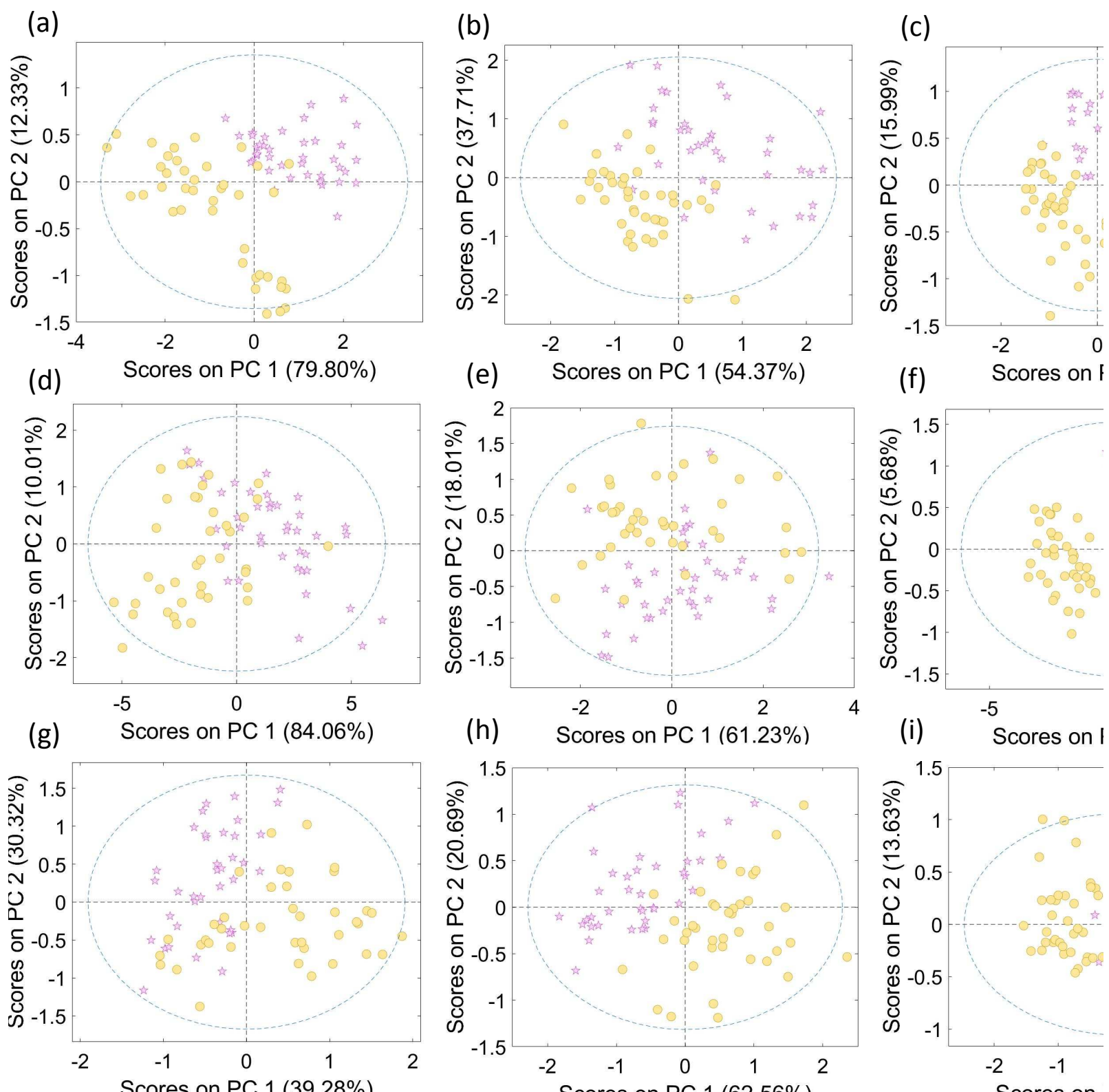
MANUSCRIPT











Highlights

- NIRS was used for discriminating between three *Prunus dulcis* varieties.
- Several spectral pre-treatment strategies were investigated.
- A combination of SNV, SG first derivative, and mean centering methods was optimal.
- Tree variety and leaf age were the most important classification factors for PLS-DA.
- NIRS is a rapid and economical method for nursery plant classification.