

## **HIGHLIGHTS**

- Computer vision is feasible for pine nuts species recognition
- Image analysis ensures excellent results through shape and size measurements
- NIR spectroscopy shows classification performances similar to image analysis
- Classification accuracy is higher than 95% for each tested model
- The suggested approach provides the basis for a rapid online detection system

1 **PINE NUT SPECIES RECOGNITION USING NIR SPECTROSCOPY AND IMAGE**  
2 **ANALYSIS**

3 Roberto Moschetti<sup>a</sup>, Daniel Hagos Berhe<sup>a,b</sup>, Mariagrazia Agrimi<sup>a</sup>, Ron P. Haff<sup>c</sup>, Peishih Liang<sup>c</sup>,  
4 Serena Ferri<sup>d</sup>, Danilo Monarca<sup>d</sup>, Riccardo Massantini<sup>a\*</sup>

5 <sup>a</sup> Department for Innovation in Biological, Agro-food and Forest system, University of Tuscia,  
6 Via S. Camillo de Lellis snc, 01100 Viterbo, Italy

7 <sup>b</sup> Department of Natural Resources Management, Adigrat University, P. O. Box 50, Adigrat,  
8 Tigray, Ethiopia

9 <sup>c</sup>United States Department of Agriculture, Agricultural Research Service, Western Regional  
10 Research Center, 800 Buchanan St, Albany, CA, 94710, United States.

11 <sup>d</sup> Department of Agricultural and Forestry Sciences, University of Tuscia, Via S. Camillo de  
12 Lellis snc, 01100 Viterbo, Italy

13 \*Corresponding author: Tuscia University, Department for Innovation in Biological, Agro-food  
14 and Forest system, S. Camillo de Lellis snc, 01100 Viterbo, Italy. Tel.: +39 0761 357496; e-mail  
15 address: massanti@unitus.it (Massantini, R.)

16 **ABSTRACT**

17 NIR spectroscopy and physical properties derived from image analysis were evaluated as  
18 potential features for the classification of seed kernels from two pine nut species (*P. pinea* L. and  
19 *P. sibirica* Du Tour) using Partial Least Squares Discriminant Analysis (PLS-DA). Model  
20 performances were evaluated in terms of specificity, sensitivity and accuracy. Data pre-  
21 treatments were essential for achieving excellent performances (accuracy rate > 95%) in all tests.  
22 The interval PLS-DA highlighted that the most important features for (1) the NIR method were  
23 the absorption bands at 1640-1658, 1720-1738 and 1880-1998 nm, while for (2) the image  
24 analysis were kernel eccentricity, kernel major axis length, kernel **lightness** (L\*) and kernel  
25 perimeter. The results demonstrate potential of both techniques for discriminating the two pine  
26 nut species.

27

## 28 INTRODUCTION

29 Pine nuts are edible seeds harvested from pinecones (family *Pinaceae*, genus *Pinus*)  
30 commonly used in cuisines worldwide. Although generally referred to as nuts, they are in fact  
31 classified as seeds with the edible part (the embryo) surrounded by a hard shell. While unshelled  
32 pine nuts have a long shelf life if kept dry, shelled nuts deteriorate rapidly and are susceptible to  
33 rancidity. Pine nuts have been harvested (both wild and domesticated) for human consumption  
34 since prehistoric times (Awan and Pettenella, 2017). Both naturally occurring stands and  
35 cultivated plantations of pine species are found in Asia, Europe, the Near East and North  
36 America (Nergiz and Dönmez, 2004). About 20 species are harvested in significant quantities,  
37 the most important for human consumption being *P. pinea* L. (Mediterranean stone pine), *P.*  
38 *koraiensis* Siebold & Zucc. (Korean pine), and *P. sibirica* Du Tour (Siberian pine). Nuts from  
39 the pine species *Pinus pinea* L. (commonly named Mediterranean stone pine or Italian stone  
40 pine) have long been an important component of the Mediterranean diet.

41 Pine nuts are a small (< 1 %) but rapidly growing segment of the global tree nut market. In  
42 fact, over the three-year period between 2015/2016 and 2018/2019 global production nearly  
43 doubled from around 18,600 metric tons to 34,000 metric tons (INC, 2017; Statista, 2019).  
44 Because the demand always exceeds the supply; prices vary widely by region (Sharashkin and  
45 Gold, 2004), and have been reported as high as 100 €/kg (retail) (Calama et al., 2016).  
46 Production costs are high in Europe due to multiple factors. Intensive manual labour is required  
47 for collection of cones and for separating the nuts from the shells, as no automated methods have  
48 yet been developed. Climate change has contributed to declines in pine nut production, partially  
49 through the spread of the western conifer insect (*Leptoglossus occidentalis* Heidemann) with the  
50 fungus *Diplodia sapinea* in some Mediterranean countries. Furthermore, local supplies are also  
51 insufficient to meet demand in North America (Awan and Pettenella, 2017; Parks, 2017;  
52 Vanhanen and Savage, 2013). Thus, most pine nuts are imported from China, Korea and  
53 Pakistan, with China being the main exporting country with 78% of total exports (INC, 2017).

54 The chronic global shortage of pine nuts from the traditional species, especially *P. pinea*,  
55 has led to introduction into the marketplace of nuts of other species, mostly but not exclusively  
56 *P. sibirica* and *P. koraiensis* from China (Ballin, 2012; Loewe et al., 2017). The latter nuts are  
57 sold in local markets at lower prices and are not easily distinguishable from *P. pinea* by  
58 consumers, even though they have different flavour, shape and size (Evaristo et al., 2010). In

59 fact, while significant differences exist in price and quality of nuts from different origins, this is  
60 not transparent to consumers since they are not reported on product labels (Awan and Pettenella,  
61 2017; Mutke et al., 2012). Tracking information for nut production and distribution streams is  
62 not well maintained as processing centres lose track of points of origin (Awan and Pettenella,  
63 2017; Mutke et al., 2012). Most nuts exported from China have their origins in other countries  
64 (Sharashkin and Gold, 2004). Shipping documents often lack product status details such as  
65 whether they are cones, unshelled pine nuts, or shelled kernels (Awan and Pettenella, 2017;  
66 Sharashkin and Gold, 2004), making the tracing of product origins difficult. Given the difficulty  
67 of product tracking coupled with visual similarities between species, the likelihood of product  
68 adulteration (undeclared mixing of lower value product with high value product) is high. For  
69 pine nuts, the issue of adulteration has been recently elevated to a health issue as well as an  
70 economic issue because of the occurrence dysgeusia, or Pine Nut Syndrome (PNS) in which  
71 consumers can experience a strong metallic or acid aftertaste that lasts up to 14 days in most  
72 cases but can last as much as 42 days in extreme cases (Awan and Pettenella, 2017). PNS has  
73 been linked to *P. armandii* Franch (Awan and Pettenella., 2017), its subspecies and other  
74 varieties, such as *P. sibirica* and *P. koraiensis* (Matthäus et al., 2018; Mikkelsen et al., 2014).

75 Given the economic and potential health impacts of undeclared mixing of pine nut species  
76 for consumer consumption it would be advantageous to develop methods to identify nut origins  
77 at all points of distribution (Loewe et al., 2017). Traceability has the potential to improve product  
78 quality and safety, provide geographic identification, and hinder black market trade. Finally, it  
79 would facilitate organic certification labelling (Mutke et al. 2012).

80 The use of NIR spectroscopy to detect adulteration in food products is widely reported and  
81 reviews of this research are available in the literature (Valand et al., 2020). NIR spectroscopy has  
82 also been reported as one potential method for tracing product origins, although available  
83 literature, especially for pine nuts, is limited. Tigabu et al. (2005) used NIR to trace the origins of  
84 seeds of Scots pine (*Pinus. sylvestris* L.) for reforestation purposes. Loewe et al. (2017) used  
85 NIR spectroscopy for the discrimination of Mediterranean stone pine nuts from Chilean  
86 plantations. To the best of our knowledge, these are the only reported studies using NIR  
87 spectroscopy to discriminate the origin of pine nuts, while no literature on the image analysis of  
88 pine nuts is available.

89 The objective of this research is to assess the potential use of both NIR spectroscopy and  
90 image analysis for the discrimination of pine nuts species of different origin and management.

## 91 MATERIALS AND METHODS

### 92 2.1 Sample Preparation

93 Whole, dried and decorticated pine nuts kernels of species certified as *P. pinea* L. and *P.*  
94 *sibirica* Du Tour were obtained either from local markets (batches 1 to 5) or provided by the  
95 BIOSIC company (Viterbo, Central Italy) (batches 6 to 9), respectively (Table1). Batches from  
96 different production areas were selected so that trained classification models would be  
97 insensitive to geographical origins, as well as to minimize the selection bias and secure chemical  
98 and physico-chemical representativeness to a large extent. One hundred samples (i.e. kernels)  
99 from each batch were randomly selected and assigned as either ‘Class I’ (*P. pinea* L.) or ‘Class  
100 II’ (*P. sibirica* Du Tour). Selection was performed manually by removing kernels with  
101 discoloration, infestation, infection or other visually apparent flaws. The ‘Class I’ and ‘Class II’  
102 groups were thus composed of 500 and 400 kernels, respectively.

### 103 2.2 Spectral acquisition

104 A Luminar 5030 Acousto-Optic Tunable Filter-Near Infrared (AOTF-NIR) Miniature  
105 ‘Hand-held’ Analyzer coupled with the bundled ‘SNAP! 2.04’ software (Brimrose Corp.,  
106 Baltimore, USA) was used to acquire NIR spectra. The analyser was equipped with reflectance  
107 post-dispersive optics, a pre-aligned dual beam lamp assembly (5 W halogen lamp), and an  
108 indium gallium arsenide (InGaAs) array (range 1100–2300 nm, 2 nm resolution) with an  
109 integrating time of 60 ms. The instrument was allowed to warm up for at least an hour to reach a  
110 stable state before the use. The reference spectrum was automatically measured by the second  
111 detector of the instrument, which was a dual-beam spectrophotometer (Moscetti et al., 2016).  
112 Because each NIR measurement of the instrument covers only a single point, each sample was  
113 measured on opposite sides of the kernel in order to acquire more representative data. Diffuse  
114 reflectance spectra were transformed into absorbance ( $A = \log_{10} R^{-1}$ ) and the average spectrum  
115 of each kernel was used for further computations.

### 116 2.3 Spectral pre-treatments

117 NIR spectra quite often suffer from problems of unwanted spectral variations and baseline  
118 shifts. In the present study the most common baseline correction methods (i.e. standard normal  
119 variation, SNV; multiplicative scatter correction, MSC; and Savitzky-Golay derivative) were

120 tested to evaluate the presence of light scattering, which can be detrimental to  
121 quantitative/qualitative analysis and lead to inaccurate results. In a solid matrix, light scattering  
122 can be related to the variability in refractive index, morphology (e.g., surface roughness) and  
123 density of sample. SNV and MSC methods mathematically differ but are similar in their  
124 outcomes. The Savitzky-Golay first, second or third derivative (i.e.  $D1f$ ,  $D2f$  and  $D3f$ ,  
125 respectively) were used alone or in combination with SNV/MSC for its capability of improving  
126 the baseline correction as well as resolving overlapping bands, as resolution enhancement  
127 method. Each derivative pre-treatment employed a second or third order polynomial fitted over a  
128 window of five, seven, nine, or eleven features (Moscetti et al., 2017) with the aim of improving  
129 the signal to noise ratio of the spectra. The narrowest window size (i.e. filter length of 3 features)  
130 was excluded to circumvent noise inflating of the original spectrum due to derivative calculation.  
131 The highest window size was set to 11 to avoid rounding-off of the peaks and troughs, which  
132 happens when the filter length is too wide. Every possible combination of these spectral pre-  
133 treatments was tested and only the best model, in terms of performance metrics, was retained.  
134 Finally, spectral data were mean centered.

## 135 **2.4 Image acquisition and image segmentation**

136 Allowing for the exploratory nature of this research activity, a flatbed scanner was chosen as  
137 the imaging acquisition system as it is better suited for scanning still images than a digital  
138 camera, which is more suitable for high-speed applications (Sun, 2020). Thus, digital colour  
139 images of the kernels were acquired one batch at a time (i.e. 100 kernels per scan) using a  
140 CM2320nf (Hewlett-Packard, Palo Alto, USA) flatbed scanner with VueScan 9.2.11 Professional  
141 Edition software (Hamrick Software, Phoenix, AZ, USA). The scanner was profiled using the  
142 Colour Checker Passport target (X-Rite Ltd., U.K.). Each batch was scanned three times and  
143 averaged for the final image. The scan was performed only on one side of the sample because (1)  
144 all selected nuts were free from apparent visually flaws and discolouration and (2) to avoid  
145 redundant data from spatial features. Each scan was conducted with the following parameters:  
146 resolution of  $2466 \times 3498$  pixels (240 dpi); 48-bit colour intensity resolution (16 bits per RGB  
147 channel), and; Digital Negative (DNG) raw file format. Each DNG image was colour corrected  
148 by applying the scanner profile through the Camera Raw 6.0 software (Adobe Systems Inc., San  
149 Jose, CA) then saved in Tagged Image File format. Finally, image segmentation was performed  
150 by distinguishing kernels from the image background using a custom script written in Python

151 3.7.0 coupled with OpenCV 3.4.1, employing Otsu's binarization for clustering-based image  
152 thresholding (OpenCV, 2020).

### 153 *CIELab colour measurements*

154 The custom script also extracted the average colour data from each kernel. Specifically, the  
155 'cvtColor()' function of the OpenCV library was applied for RGB-to-CIELab colour space  
156 conversion. CIELab was chosen as a uniform colour space and then as a coordinate system in  
157 which perceived colour differences correspond to Euclidean distances. Thus, results were  
158 expressed in terms of **lightness** ( $L^*$ ), red/green colour ( $a^*$ ), yellow/blue colour ( $b^*$ ), hue angle  
159 ( $h$ ) and chroma ( $C^*$ ) (Moscetti et al., 2013a). Moreover, the CIE 1976 color difference ( $\Delta E^* =$   
160  $[\Delta L^{*2} + \Delta a^{*2} + \Delta b^{*2}]^{1/2}$ ) was computed to evaluate how much the two species deviated from  
161 each other in terms of pine nut colour. Its value was expressed according to the evaluation scale  
162 used by Cecchini et al. (2011): imperceptible ( $\Delta E^* < 1$ ), minimal ( $1 \leq \Delta E^* < 2$ ), just perceptible  
163 ( $2 \leq \Delta E^* < 3$ ), perceptible ( $3 \leq \Delta E^* < 5$ ), strong difference ( $5 \leq \Delta E^* < 12$ ) and different colour  
164 ( $\Delta E^* \geq 12$ ).

### 165 *Shape and size measurements*

166 Kernel perimeter (mm), major and minor axis lengths (mm), surface area ( $\text{mm}^2$ ) and  
167 eccentricity were extracted from each image using the 'regionprops()' function (Matlab R2015a  
168 'Image Processing' toolbox). The kernel perimeter was computed as the sum of the distances  
169 between adjacent pixels around the edge of the kernel. Major and minor axis lengths were  
170 computed on the ellipse with the same normalized second central moments (rotational inertia) of  
171 the kernel region. Area was acquired by calculating the number of pixels in the kernel, and  
172 eccentricity was computed as the ratio of the distance between the ellipse foci and the major axis  
173 length. Size and shape parameters were converted from pixel to metric unit (mm) based on the  
174 reference embedded in the Colour Checker Passport target (X-Rite Ltd., U.K.). In total, 10  
175 features based on imaging data were extracted from each kernel image as features for  
176 classification models (Table 2).

177 Imaging features were scaled before model development.

## 178 **2.5 Classification models development**

179 Spectral-based and imaging-based classification models were developed separately using  
180 Partial Least Squares Discrimination Analysis (PLS-DA) through the SIMPLS algorithm (de  
181 Jong, 1993). PLS-DA seeks linear combinations of the original independent variables with the

182 ability to discriminate data classes while discarding irrelevant and unstable information and  
183 solving collinearity issues. In addition, Interval PLS-DA (iPLS-DA) was used to select subsets of  
184 features which could still achieve good prediction results (Xing and Guyer, 2008). The iPLS-DA  
185 algorithm used the stepwise forward mode to select a maximum of 10 intervals of 10 features  
186 each for spectral model and 10 intervals of 1 feature each for imaging model. The number of the  
187 interval which provides the lowest RMSECV was selected by the algorithm, allowing to  
188 circumvent under-/overfitting problems. Both selectivity ratio and  $\beta$ -coefficient were computed  
189 to assess the relative contribution of each feature or subset of features to the performance of each  
190 model (Rajalahti et al., 2009). The larger the selectivity ratio and/or the absolute value of  $\beta$ -  
191 coefficient, the more useful the given feature is in classification.

192 In a first step, samples (spectra or images) for each class were randomly split between  
193 calibration set (75% or 675 kernels) and prediction set (25% or 225 kernels). As PLS regression  
194 performs a dimensionality reduction, it is essential to test the model and select the correct  
195 number of latent variables to find the optimal trade-off between under-fitting and over-fitting.  
196 Thus, a model optimization was conducted using venetian blinds cross-validation with 10 data  
197 splits. Root Mean Square Error for calibration, cross-validation and prediction calculations  
198 (RMSEC, RMSECV and RMSEP, respectively) were used to evaluate each discriminant model  
199 with the purpose of selecting the optimal number of latent variables and, thus, to circumvent  
200 unrealistic results (Moscetti et al., 2015). In a second step, the prediction set was used to give an  
201 independent assessment of the accuracy, precision and robustness of the calibration model,  
202 following the criteria contained in the adopted validation guidelines (Broad et al., 2006). The  
203 Hotelling's t-squared and Q residuals were used in combination to identify potential outliers of  
204 each given model. Chemometrics was performed using Matlab software R2015a coupled with  
205 PLS\_Toolbox software v8.1 (Eigenvector Research Inc., WA, USA).

206 The classification performance of selected models was tested in terms of sensitivity (Eq. 1),  
207 selectivity (Eq. 2) and accuracy (Eq. 3) (Dejaegher et al., 2011). While accuracy rates are  
208 correlated with model predictivity, sensitivity and specificity rates are generally used to evaluate  
209 robustness (i.e. the capability of the model to resist to small changes in test conditions).

210 (1) 
$$\text{Sensitivity rate} = \frac{\text{True Positives}}{\text{True Positives} + \text{False Negatives}}$$



211 (2) 
$$\text{Selectivity rate} = \frac{\text{True Negatives}}{\text{False Positives} + \text{True Negatives}}$$

212 (3) 
$$\text{Accuracy rate} = \frac{\text{True Positives} + \text{True Negatives}}{\text{Total Positives} + \text{Total Negatives}}$$

## 213 2.6 Moisture content analysis

214 Moisture content was measured following the official oven-drying method ‘Moisture in  
215 Dried Fruits’ - AOAC 934.06. Results were obtained at batch level and expressed as a  
216 percentage by mass (grams per 100 grams) (Horwitz, 2005). The analysis was performed with  
217 the aim of (i) assessing the potential difference in terms of moisture content between classes and  
218 (ii) evaluating the impact of moisture content on developed NIR models, as it is well known that  
219 NIR spectroscopy is sensitive to the water content of a food matrix.

## 220 2.7 Data handling and data analysis

221 One-way analysis of variance (ANOVA) was performed to evaluate statistical differences  
222 between classes in terms of moisture content, colour as well as shape and size. The Tukey’s  
223 pairwise comparison method was performed, and the Honestly Significant Difference (HSD) was  
224 calculated for an appropriate level of interaction ( $P \leq 0.05$ ). Results were reported as the mean  
225 and standard error of the mean. Data handling and ANOVA were both performed using R v3.3.3  
226 software in combination with ‘dplyr’ v0.5.0 and ‘agricolae’ v1.2-4 R-packages (CRAN, 2017).

# 227 RESULTS AND DISCUSSION

## 228 3.1 Data overview

### 229 Spectral data

230 Although visual inspection of the spectra is not sufficient to distinguish kernels from  
231 different species, their graphical overview may help understanding spectral differences between  
232 classes. Figure 1 thus illustrates the mean absorbance spectra of pine nuts from *P. pinea* and *P.*  
233 *sibirica* in the full NIR spectral range from 1100 to 2300 nm both without (a) and with (b) the  
234 best spectral pre-treatments. The best classification performance was obtained following spectral  
235 pre-treatments of SNV (standard normal variate) in combination with 1<sup>st</sup> order derivative of 2<sup>nd</sup>  
236 order polynomial over window 9, regardless of the classification algorithm used (i.e. PLS-DA or  
237 iPLS-DA). It means that the classification model was negatively affected by light scattering,  
238 while improved by noise filtering and unravelling overlapping bands. The PLS-DA model was

239 based on the first two latent variables that were able to capture much of the variability between  
240 species. In the same Figure, the 10-feature intervals selected by the iPLS-DA algorithm are also  
241 showed. Each interval represents a window of 10 adjacent features which gave superior  
242 prediction over those obtained using all variables in the spectral data set.

### 243 *Imaging data*

244 Summary statistics and ANOVA results for image-based features are shown in Table 2.  
245 Lightness, hue angle, chroma, yellow/blue colour and red/green colour of kernels showed  
246 significant variation between species ( $P \leq 0.05$ ) with higher values of lightness and red/green  
247 colour for *P. pinea* and higher values of hue angle, chroma and yellow/blue colour for *P.*  
248 *sibirica*. In other words, *P. sibirica* showed a yellower and vivid, but darker colour than *P. pinea*.  
249 This colour difference between the two species can be detectable by human eye as corroborated  
250 by the  $\Delta E^*$ , which assumed an average value of 4.68 (i.e.,  $3 \leq \Delta E^* < 5$ , perceptible difference of  
251 colour) and ranged from a minimum of 1.55 (i.e.,  $1 \leq \Delta E^* < 2$ , minimal difference of colour) to a  
252 maximum of 26.37 (i.e.,  $\Delta E^* \geq 12$ , different colours).

253 Kernel shape and size were significantly different ( $P \leq 0.05$ ) in terms of perimeter, major  
254 axis, eccentricity and area, although the minor axis did not significantly differ between the two  
255 species. Thus, *P. pinea* was bigger and more elliptical (i.e. higher eccentricity) than *P. sibirica*.

256 While discrimination based on imaging techniques have been widely reported for many  
257 commodities including walnuts (Calama et al., 2017; Ercisli et al., 2012; Huang et al., 2016; Kuo  
258 et al., 2016; Liu et al., 2015; Lurstwut and Pornpanomchai, 2018; Menesatti et al., 2008;  
259 Pallottino et al., 2009; Rodríguez-Pulido et al., 2012; Wu et al., 2018; Zhang et al., 2016), none  
260 have been reported involving pine nuts making comparison to this study problematic.

## 261 **3.2 Classification models**

### 262 *Prediction models based on NIR spectral data*

263 Figure 2a and 2b show score plots derived from the PLS-DA and iPLS-DA models,  
264 respectively, for the prediction set. Good separation between classes was observed for both  
265 models (species). No outliers were detected for either classification model. Models affected by  
266 outliers were discarded due to low classification performances which were the result of the lack  
267 of proper selection of data pre-treatments and number of latent variables. Table 3 showed the  
268 performance metrics of both models. PLS-DA yielded 10 misclassifications versus 6 for iPLS-  
269 DA, corresponding to very good accuracy rates equal to 96 % and 98 %, respectively. For PLS-

270 DA, 80% of misclassifications were for *P. pinea* and this increased to 100% for the iPLS-DA  
271 model. This suggests that derived models intending to evaluate the purity of *P. pinea* product  
272 might tend to have higher false negative (*P. pinea* misclassified; 1- sensitivity) than false  
273 positive results (*P. sibirica* misclassified; 1- specificity).

274 Misclassifications might be related to variances among features within each class or noise.  
275 Higher accuracy for iPLS-DA is not unexpected since it uses feature regions rather than discrete  
276 features (i.e. wavebands rather than wavelengths) thereby focusing on important spectral regions  
277 and removing interference from other regions (Xiaobo et al., 2010).

278 The PLS-DA model was characterized by 2 latent variables, while iPLS-DA only by one.  
279 The observed decrease in number of latent variables is quite common after variable selection,  
280 further reducing the complexity of models and making it easier to interpret. Therefore, it is also  
281 important to consider that a single latent variable may not be enough to allow a model to exploit  
282 all the available information, leading to the risk of underfitting. In our case, the choice for a  
283 single latent variable for the iPLS-DA model was driven by cross-validation, which proved that  
284 the additional latent variables did not significantly improve model performance.

285 To evaluate the contribution of each individual feature on the PLS-DA model two diagnostic  
286 methods were computed and compared: (i) the selectivity ratio (Figure 4a) and (ii) the  $\beta$   
287 coefficients (Figure 4c). Conventionally, features with  $\beta$  coefficients larger than  $\pm 2$  times  
288 standard deviation of the regression vector were selected as important wavebands. Based on this,  
289 spectral regions around 1644, 1722 and 1858 nm were ranked as the most useful. However, this  
290 approach may lead to selecting or ignoring features with low or high contribution to the model,  
291 respectively (Rajalahti et al., 2009). Consequently, the selectivity ratio was also computed for its  
292 capability to highlight important features by combining predictive power ( $\beta$ -coefficients) with  
293 explanatory power (variance/covariance among features). Using this new diagnostic method, the  
294 cut-off ratio between the explained and the residual variance was computed and features above  
295 the cut-off were recognised as important for the model. In fact, the results indicated that the  
296 absorption bands in the regions around 1270, 1410-1450, 1550-1630, 1720 and 1870 nm had the  
297 highest selectivity ratio (up to 6.42), i.e. highest contribution to the model. However, below  
298 ~1250 nm, between 1724 and 1866 nm, as well as beyond ~1900 nm observed selectivity ratios  
299 were lower than the cut-off, suggesting a lower contribution to the model.

300 The iPLS-DA algorithm iteratively selected 10-feature intervals which provided the lowest  
301 model root-mean-square error of cross-validation (RMSECV). Accordingly, the absorption  
302 bands were: (i) 1640-1658 nm, (ii) 1720-1738 nm and (iii) 1880-1898 nm. These bands partially  
303 corresponded to the marker wavelengths selected through the selectivity ratio and, in general, are  
304 associated with the (i) lipids, (ii) proteins and (iii) carbohydrates, as well as moisture (Loewe et  
305 al., 2017; Tigabu et al., 2005; Workman and Weyer, 2008).

306 Absorption peaks have been reported in different NIR spectral regions of agricultural  
307 commodities and attributed to a variety of molecular processes. For instance, Loewe et al. (2017)  
308 reported that absorption peaks observed at 1200, 1500, 1720, 1760, 1940 and 2350 nm in spectra  
309 of Mediterranean pine nuts grown in Chile corresponded to the most important wavebands for  
310 discrimination of product origin. The importance of these bands in derived classification models  
311 was attributed to known absorption peaks for proteins, amino acids and moisture (1200 and 1940  
312 nm), lipids (1500; 1720; 1760 and 2350 nm), and polysaccharides and carbohydrates (2350 nm).  
313 Tigabu et al. (2005) identified strong NIR absorption peaks in the spectra of pine nut seeds from  
314 different sources at 1422 and 1930 nm along with less visible bands around 1500 and 2200 nm  
315 which were the most important features for discriminant models. Moreover, Moscetti et al.,  
316 (2013b) reported complex bands in hazelnut spectra attributable to various molecular vibrations  
317 of functional groups, some linked with fatty acid content, unsaturated fatty acids with cis double  
318 bonds, protein acid and ester stretching vibrations, and others to polysaccharides, carbohydrates  
319 and lipids.

320 The moisture content of both kernel classes was measured to determine whether water bands  
321 had the potential to affect the NIR-based model development. The results showed a significantly  
322 different moisture content between *P. pinea* ( $5.36\pm 0.06$  %) and *P. sibirica* ( $3.22\pm 0.09$  %).  
323 Results agreed with the literature, which showed that kernels from *P. pinea* had a higher  
324 moisture content ( $\sim 5.10$ - $5.60$ %) than kernels belonging to *P. sibirica* ( $< 3.50$ %) (Evaristo et al.,  
325 2010; Nergiz and Dönmez, 2004). Therefore, both NIR-based models (i.e. PLS-DA and iPLS-  
326 DA) were probably affected by moisture content. In fact, the selectivity ratio showed that  
327 features close to the  $\sim 1480$ - and  $\sim 1900$  nm water bands had high importance for both PLS-DA  
328 and iPLS-DA models. Considering that “*pine nut kernels should have a moisture content not*  
329 *exceeding 3.5 per cent, except for Pinus pinea, which should not exceed 6.0 per cent and Pinus*  
330 *gerardiana, which should not exceed 7.0 per cent*” (UNECE, 2013), the moisture content was

331 expected to significantly improve discriminant performances of models and, thus, the water  
332 bands were not excluded.

### 333 *Prediction models based on imaging data*

334 Figure 3 reports the score plots of the PLS-DA and iPLS-DA of image analysis for the  
335 prediction set. Both models showed very good and similar performance metrics ranging from 95  
336 to 98 % (Table 3). Misclassifications were slightly higher for iPLS-DA (11 versus 9) with both  
337 heavily weighted towards misclassification of *P. pinea* as was also observed above for NIR-  
338 based models. **Similar to the observation for the NIR-based models, the imaging-based iPLS-DA**  
339 **models was characterized by only one latent variable.** In all cases, the sensitivity, specificity and  
340 accuracy rates showed similarity towards calibration, cross-validation and prediction, indicating  
341 that the models were robust. As observed for the NIR-based models, outliers were not detected  
342 for the selected imaging-based models.

343 Figure 4b and 4d show the selectivity ratio and the  $\beta$ -coefficients bar plots of the PLS-DA  
344 model, respectively. Regarding the  $\beta$ -coefficients, **lightness** ( $L^*$ ) was the only feature  
345 characterized by a value larger than  $\pm 2$  standard deviation over the regression vector. However,  
346 similar to the observation for the NIR-based model, the selectivity ratio suggested a higher  
347 number of marker features. This suggests that  $\beta$ -coefficients may underestimate the feature's  
348 contribution in the model under the studied experimental conditions. In fact, the selectivity ratio  
349 indicated that among the 10 imaging features, **lightness**, perimeter, eccentricity and major axis  
350 length had the strongest contributions in the PLS-DA model, with eccentricity and major axis  
351 length scoring highest.

352 The features selected by the iPLS-DA algorithm were mostly the same as those observed as  
353 marker wavelengths for the PLS-DA model. In fact, they consisted of 3 colour attributes (i.e.  
354 **lightness**, red/green colour and hue angle) and 2 spatial properties (i.e. eccentricity and major  
355 axis length). Size features (i.e. area and perimeter) were discarded, suggesting that shape  
356 recognition played a major role in the classification of the two pine nut species considered in this  
357 study.

## 358 **CONCLUSIONS**

359 Pine nuts are widely consumed in domestic and foreign markets and, despite their  
360 importance, the industry faces market challenges related to adulteration and subsequent potential

361 for PNS. This study addresses these challenges by demonstrating the potential use of NIR  
362 spectroscopy and image analysis to distinguish pine nuts from different geographic origins.

363 NIR spectra of two species of pine nuts were subjected to different pre-treatments and  
364 presented as input features for PLS-DA and iPLS-DA discrimination. Absorption bands at 1640  
365 – 1658 nm, 1720-1738 nm and 1880-1998 nm were found to be the most important for  
366 classification purposes.

367 Various features derived from image analysis, including CIELab colour space and measured  
368 physical properties, were also tested for their ability to distinguish classes. The dominant features  
369 in the resulting models were eccentricity, major axis length, and perimeter, based on calculations  
370 of selectivity ratios. PLS-DA model performance based on four latent variables was above 95%  
371 accuracy in classifying the pine nuts. The iPLS-DA model required only one latent variable to  
372 achieve greater than 95% accuracy for both calibration and prediction.

373 Based on the present study findings, it can be concluded that either NIR spectroscopy or  
374 image analysis coupled with chemometrics have potential for the classification of pine nuts  
375 species. Use of these techniques could improve the traceability of pine nuts, which is essential  
376 for controlling quality and the incidence of pine nut syndrome (PNS). However, prior to the  
377 implementation of this approach in industry further study is recommended to; (i) elucidate the  
378 major chemical constituents, i.e. moisture, fat, and protein content, and fatty acid profile, related  
379 to the spectral ranges on which classification models rely; and (ii) to validate each model with  
380 larger sample sizes and different regions, production years, agro-pedo-climatic conditions, and  
381 species. In addition, the effect of (i) fluctuations in the moisture content of fruit, (ii) excluding  
382 NIR water bands, and/or (iii) using spectral and spatial data in combination (i.e. multi-  
383 /hyperspectral imaging) should also be investigated in the future.

## 384 **ACKNOWLEDGMENTS**

385 The authors gratefully acknowledge (1) the ‘Departments of excellence 2018’ program (i.e.  
386 ‘Dipartimenti di eccellenza’) of the Italian Ministry of Education, University and Research for  
387 the financial support through the ‘Landscape 4.0 food, wellbeing and environment’ (DIBAF  
388 department of University of Tuscia); (2) the BIOSIC srl (Viterbo, Central Italy) for providing  
389 samples; and (3) MSc Gianpaolo Moschetti and Dr. Swathi Sirisha Nallan Chakravartula for the  
390 English language revision of the manuscript.

391 **REFERENCES**

- 392 Awan, H.U.M., Pettenella, D., 2017. Pine nuts: A review of recent sanitary conditions and  
393 market development. *Forests* 8. <https://doi.org/10.3390/f8100367>
- 394 Ballin, N.Z., 2012. Investigating cases of taste disturbance caused by pine nuts in Denmark, in:  
395 *Case Studies in Food Safety and Authenticity*. Elsevier, pp. 318–325.  
396 <https://doi.org/10.1533/9780857096937.6.318>
- 397 Broad, N., Graham, P., Hailey, P., Hardy, A., Holland, S., Hughes, S., Lee, D., Prebble, K.,  
398 Salton, N., Warren, P., Leiper, K., 2006. Guidelines for the Development and Validation of  
399 Near-Infrared Spectroscopic Methods in the Pharmaceutical Industry. *Handb. Vib.*  
400 *Spectrosc.* <https://doi.org/10.1002/0470027320.s8303>
- 401 Calama, R., Fortin, M., Pardos, M., Manso, R., 2017. Modelling spatiotemporal dynamics of  
402 *Pinus pinea* cone infestation by *Dioryctria mendacella*. *For. Ecol. Manage.* 389, 136–148.  
403 <https://doi.org/10.1016/j.foreco.2016.12.015>
- 404 Calama, R., Gordo, J., Madrigal, G., Mutke, S., Conde, M., Montero, G., Pardos, M., 2016.  
405 Enhanced tools for predicting annual stone pine (*Pinus pinea* L.) cone production at tree and  
406 forest scale in Inner Spain. *For. Syst.* 25. <https://doi.org/10.5424/fs/2016253-09671>
- 407 Cecchini, M., Contini, M., Massantini, R., Monarca, D., Moscetti, R., 2011. Effects of controlled  
408 atmospheres and low temperature on storability of chestnuts manually and mechanically  
409 harvested. *Postharvest Biol. Technol.* 61, 131–136.
- 410 CRAN, 2017. Comprehensive R archive Network.
- 411 de Jong, S., 1993. SIMPLS: An alternative approach to partial least squares regression. *Chemom.*  
412 *Intell. Lab. Syst.* 18, 251–263.
- 413 Dejaegher, B., Dhooghe, L., Goodarzi, M., Apers, S., Pieters, L., Vander Heyden, Y., 2011.  
414 Classification models for neocryptolepine derivatives as inhibitors of the  $\beta$ -haematin  
415 formation. *Anal. Chim. Acta* 705, 98–110.
- 416 Ercisli, S., Sayinci, B., Kara, M., Yildiz, C., Ozturk, I., 2012. Determination of size and shape  
417 features of walnut (*Juglans regia* L.) cultivars using image processing. *Sci. Hortic.*  
418 (Amsterdam). 133, 47–55. <https://doi.org/10.1016/j.scienta.2011.10.014>
- 419 Evaristo, I., Batista, D., Correia, I., Correia, P., Costa, R., 2010. Chemical profiling of  
420 Portuguese *Pinus pinea* L. nuts. *J. Sci. Food Agric.* 90, 1041–1049.  
421 <https://doi.org/10.1002/jsfa.3914>

422 Horwitz, W., 2005. Official methods of analysis of AOAC international, 18th Edition.

423 Huang, M., Tang, J., Yang, B., Zhu, Q., 2016. Classification of maize seeds of different years  
424 based on hyperspectral imaging and model updating. *Comput. Electron. Agric.* 122, 139–  
425 145. <https://doi.org/10.1016/j.compag.2016.01.029>

426 INC, 2017. INTERNATIONAL NUT AND DRIED FRUIT COUNCIL: Nuts & amp; Dried  
427 Fruits Statistical Yearbook 76.

428 Kuo, T.Y., Chung, C.L., Chen, S.Y., Lin, H.A., Kuo, Y.F., 2016. Identifying rice grains using  
429 image analysis and sparse-representation-based classification. *Comput. Electron. Agric.*  
430 127, 716–725. <https://doi.org/10.1016/j.compag.2016.07.020>

431 Liu, D., Ning, X., Li, Z., Yang, D., Li, H., Gao, L., 2015. Discriminating and elimination of  
432 damaged soybean seeds based on image characteristics. *J. Stored Prod. Res.* 60, 67–74.  
433 <https://doi.org/10.1016/j.jspr.2014.10.001>

434 Loewe, V., Navarro-Cerrillo, R.M., García-Olmo, J., Riccioli, C., Sánchez-Cuesta, R., 2017.  
435 Discriminant analysis of Mediterranean pine nuts (*Pinus pinea* L.) from Chilean plantations  
436 by near infrared spectroscopy (NIRS). *Food Control* 73, 634–643.  
437 <https://doi.org/10.1016/j.foodcont.2016.09.012>

438 Lurstwut, B., Pornpanomchai, C., 2018. Image analysis based on color, shape and texture for rice  
439 seed (*Oryza sativa* L.) germination evaluation. *Agric. Nat. Resour.* 51, 383–389.  
440 <https://doi.org/10.1016/j.anres.2017.12.002>

441 Matthäus, B., Li, P., Ma, F., Zhou, H., Jiang, J., Özcan, M.M., 2018. Is the Profile of Fatty  
442 Acids, Tocopherols, and Amino Acids Suitable to Differentiate *Pinus armandii* Suspicious  
443 to Be Responsible for the Pine Nut Syndrome from Other *Pinus* Species? *Chem. Biodivers.*  
444 15. <https://doi.org/10.1002/cbdv.201700323>

445 Menesatti, P., Costa, C., Paglia, G., Pallottino, F., D'Andrea, S., Rimatori, V., Aguzzi, J., 2008.  
446 Shape-based methodology for multivariate discrimination among Italian hazelnut cultivars.  
447 *Biosyst. Eng.* 101, 417–424. <https://doi.org/10.1016/j.biosystemseng.2008.09.013>

448 Mikkelsen, A.T., Jessen, F., Ballin, N.Z., 2014. Species determination of pine nuts in  
449 commercial samples causing pine nut syndrome. *Food Control* 40, 19–25.  
450 <https://doi.org/10.1016/j.foodcont.2013.11.030>

451 Moscetti, R., Carletti, L., Monarca, D., Cecchini, M., Stella, E., Massantini, R., 2013a. Effect of  
452 alternative postharvest control treatments on the storability of “Golden Delicious” apples. *J.*



453 Sci. Food Agric. 93, 2691–2697.

454 Moschetti, R., Haff, R.P., Aernouts, B., Saeys, W., Monarca, D., Cecchini, M., Massantini, R.,  
455 2013b. Feasibility of Vis/NIR spectroscopy for detection of flaws in hazelnut kernels. J.  
456 Food Eng. 118, 1–7.

457 Moschetti, R., Haff, R.P., Ferri, S., Raponi, F., Monarca, D., Liang, P., Massantini, R., 2017.  
458 Real-Time Monitoring of Organic Carrot (var. Romance) During Hot-Air Drying Using  
459 Near-Infrared Spectroscopy. Food Bioprocess Technol. 10, 2046–2059.  
460 <https://doi.org/10.1007/s11947-017-1975-3>

461 Moschetti, R., Haff, R.P., Monarca, D., Cecchini, M., Massantini, R., 2016. Near-infrared  
462 spectroscopy for detection of hailstorm damage on olive fruit. Postharvest Biol. Technol.  
463 120, 204–212. <https://doi.org/10.1016/j.postharvbio.2016.06.011>

464 Moschetti, R., Saeys, W., Keresztes, J.C., Goodarzi, M., Cecchini, M., Danilo, M., Massantini, R.,  
465 2015. Hazelnut Quality Sorting Using High Dynamic Range Short-Wave Infrared  
466 Hyperspectral Imaging. Food Bioprocess Technol. 8, 1593–1604.

467 Mutke, S., Calama, R., González-Martínez, S.C., Montero, G., Gordo, F.J., Bono, D., Gil, L.,  
468 2012. Mediterranean stone pine: Botany and horticulture. Hortic. Rev. (Am. Soc. Hortic.  
469 Sci). 39, 153–201. <https://doi.org/10.1002/9781118100592.ch4>

470 Navarro-cerrillo, R.M., García-olmo, J., Riccioli, C., 2017. Discriminant analysis of  
471 Mediterranean pine nuts ( *Pinus pinea* L .) from Chilean plantations by near infrared  
472 spectroscopy ( NIRS ) 73, 634–643. <https://doi.org/10.1016/j.foodcont.2016.09.012>

473 Nergiz, C., Dönmez, İ., 2004. Chemical composition and nutritive value of *Pinus pinea* L. seeds.  
474 Food Chem. 86, 365–368. <https://doi.org/10.1016/j.foodchem.2003.09.009>

475 OpenCV, n.d. OpenCV: Image Thresholding [WWW Document]. URL  
476 [https://docs.opencv.org/master/d7/d4d/tutorial\\_py\\_thresholding.html](https://docs.opencv.org/master/d7/d4d/tutorial_py_thresholding.html) (accessed 5.10.20).

477 Pallottino, F., Menesatti, P., Costa, C., Paglia, G., Salvador, F.R., Lolletti, D., 2009. Image  
478 Analysis Techniques for Automated Hazelnut Peeling Determination. Food Bioprocess  
479 Technol. 3, 155–159. <https://doi.org/10.1007/s11947-009-0211-1>

480 Parks, S., n.d. Is the U.S. Pine Nut Industry on the Brink of Extinction? | Civil Eats [WWW  
481 Document]. URL [https://civileats.com/2017/06/01/is-the-u-s-pine-nut-industry-on-the-](https://civileats.com/2017/06/01/is-the-u-s-pine-nut-industry-on-the-brink-of-extinction/)  
482 [brink-of-extinction/](https://civileats.com/2017/06/01/is-the-u-s-pine-nut-industry-on-the-brink-of-extinction/) (accessed 5.10.20).

483 Rajalahti, T., Arneberg, R., Berven, F.S., Myhr, K.-M., Ulvik, R.J., Kvalheim, O.M., 2009.

484 Biomarker discovery in mass spectral profiles by means of selectivity ratio plot. *Chemom.*  
485 *Intell. Lab. Syst.* 95, 35–48. <https://doi.org/10.1016/j.chemolab.2008.08.004>

486 Rodríguez-Pulido, F.J., Gómez-Robledo, L., Melgosa, M., Gordillo, B., González-Miret, M.L.,  
487 Heredia, F.J., 2012. Ripeness estimation of grape berries and seeds by image analysis.  
488 *Comput. Electron. Agric.* 82, 128–133. <https://doi.org/10.1016/j.compag.2012.01.004>

489 Sharashkin, L., Gold, M., 2004. Pinenuts: Species, Products, Markets, and Potential for U.S.  
490 Production, in: Northern Nut Growers Association 95th Annual Report. Proceeding for the  
491 95th Annual Meeting, Columbia, Missouri, August 16-19, 2004.

492 Statista, n.d. • Nuts: global production by type 2019 | Statista [WWW Document]. URL  
493 <https://www.statista.com/statistics/1030790/tree-nut-global-production-by-type/> (accessed  
494 5.10.20).

495 Sun, D.-W., 2020. *Computer Vision Technology for Food Quality Evaluation - 2nd Edition.*

496 Tigabu, M., Oden, P.C., Lindgren, D., 2005. Identification of seed sources and parents of *Pinus*  
497 *sylvestris* L. using visible – near infrared reflectance spectra and multivariate analysis.  
498 <https://doi.org/10.1007/s00468-005-0408-5>

499 UNECE, 2013. UNECE Standard for Pine Nuts (DDP-12). **Last access: 02 July 2020.**

500 Valand, R., Tanna, S., Lawson, G., Bengtström, L., 2020. A review of Fourier Transform  
501 Infrared (FTIR) spectroscopy used in food adulteration and authenticity investigations.  
502 *Food Addit. Contam. - Part A Chem. Anal. Control. Expo. Risk Assess.*  
503 <https://doi.org/10.1080/19440049.2019.1675909>

504 Vanhanen, L., Savage, G., 2013. Mineral Analysis of Pine Nuts (*Pinus* spp.) Grown in New  
505 Zealand. *Foods* 2, 143–150. <https://doi.org/10.3390/foods2020143>

506 Workman, J., Weyer, L., 2008. *Practical Guide to Interpretive Near-Infrared Spectroscopy.* CRC  
507 Press, London, UK. **ISBN: 978-1-57444-784-2**

508 Wu, Q., Xie, L., Xu, H., 2018. Determination of toxigenic fungi and aflatoxins in nuts and dried  
509 fruits using imaging and spectroscopic techniques. *Food Chem.* 252, 228–242.  
510 <https://doi.org/10.1016/j.foodchem.2018.01.076>

511 Xiaobo, Z., Jiewen, Z., Povey, M.J.W., Holmes, M., Hanpin, M., 2010. Variables selection  
512 methods in near-infrared spectroscopy. *Anal. Chim. Acta* 667, 14–32.  
513 <https://doi.org/10.1016/j.aca.2010.03.048>

514 Xing, J., Guyer, D., 2008. Comparison of transmittance and reflectance to detect insect

515 infestation in Montmorency tart cherry. *Comput. Electron. Agric.* 64, 194–201.  
516 Zhang, C., Guo, C., Liu, F., Kong, W., He, Y., Lou, B., 2016. Hyperspectral imaging analysis for  
517 ripeness evaluation of strawberry with support vector machine. *J. Food Eng.* 179, 11–18.  
518 <https://doi.org/10.1016/j.jfoodeng.2016.01.002>  
519

Figure 1

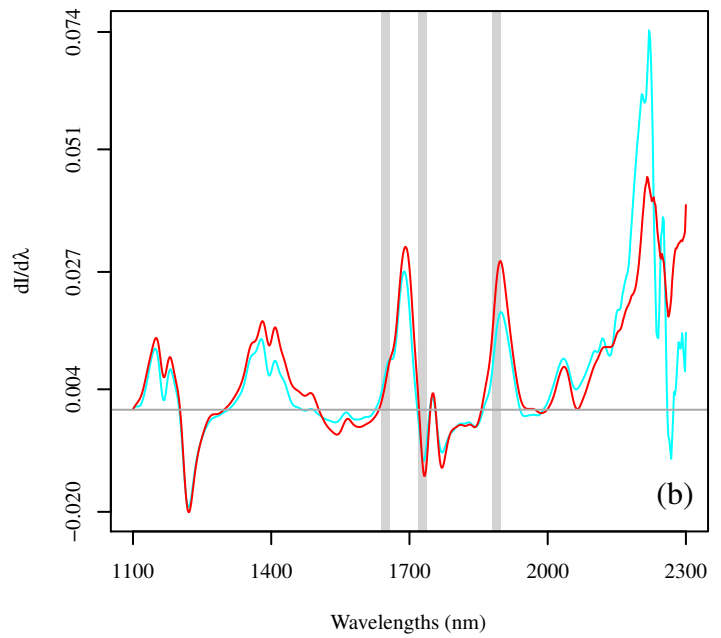
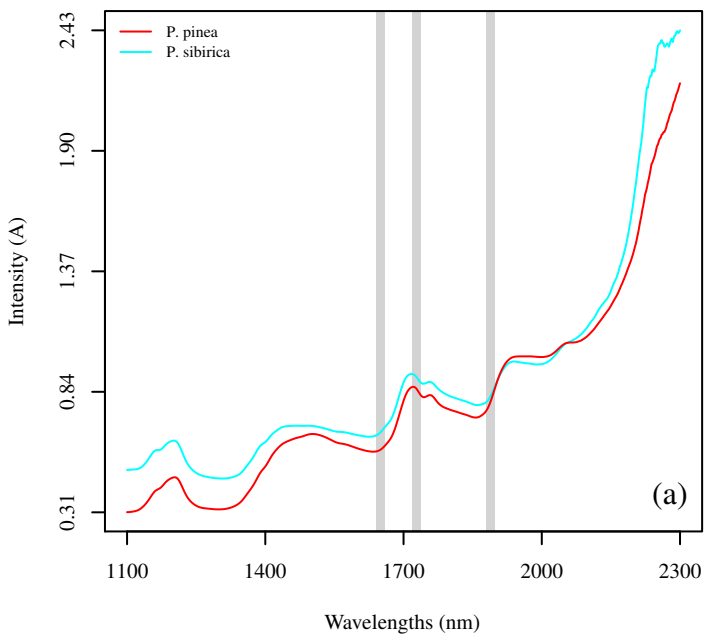


Figure 2

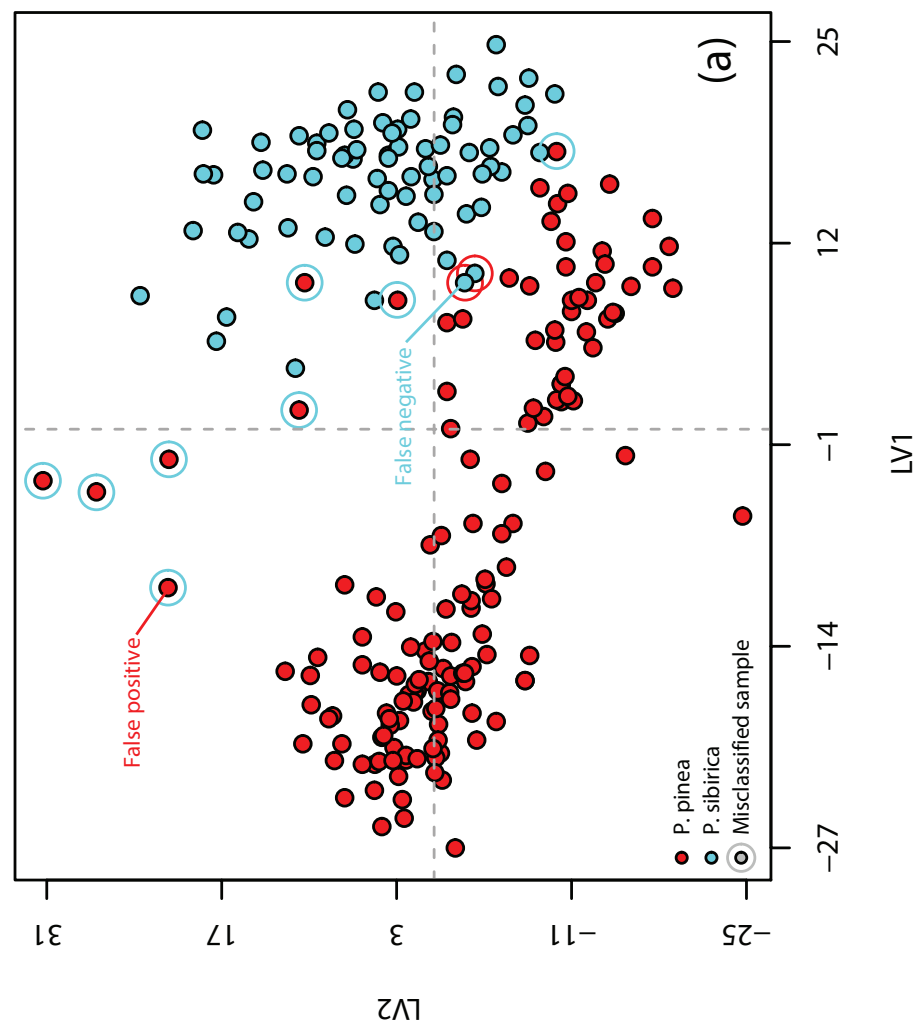
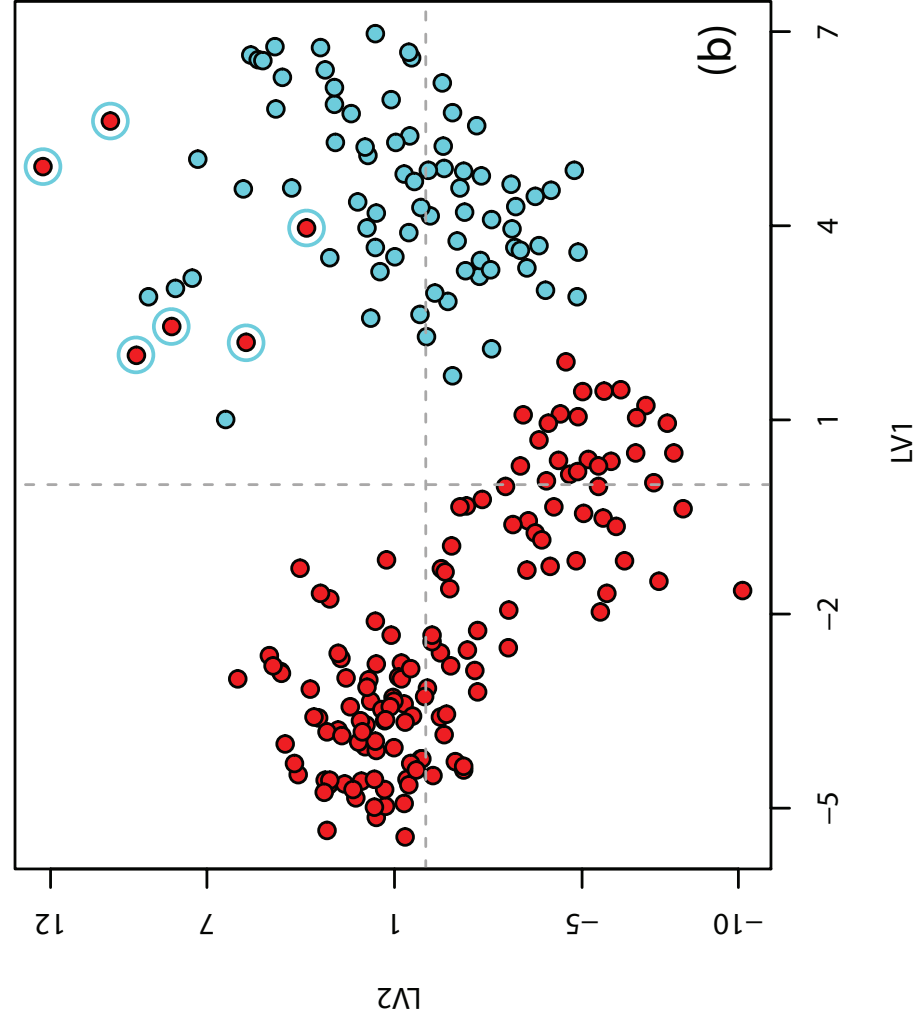


Figure 3

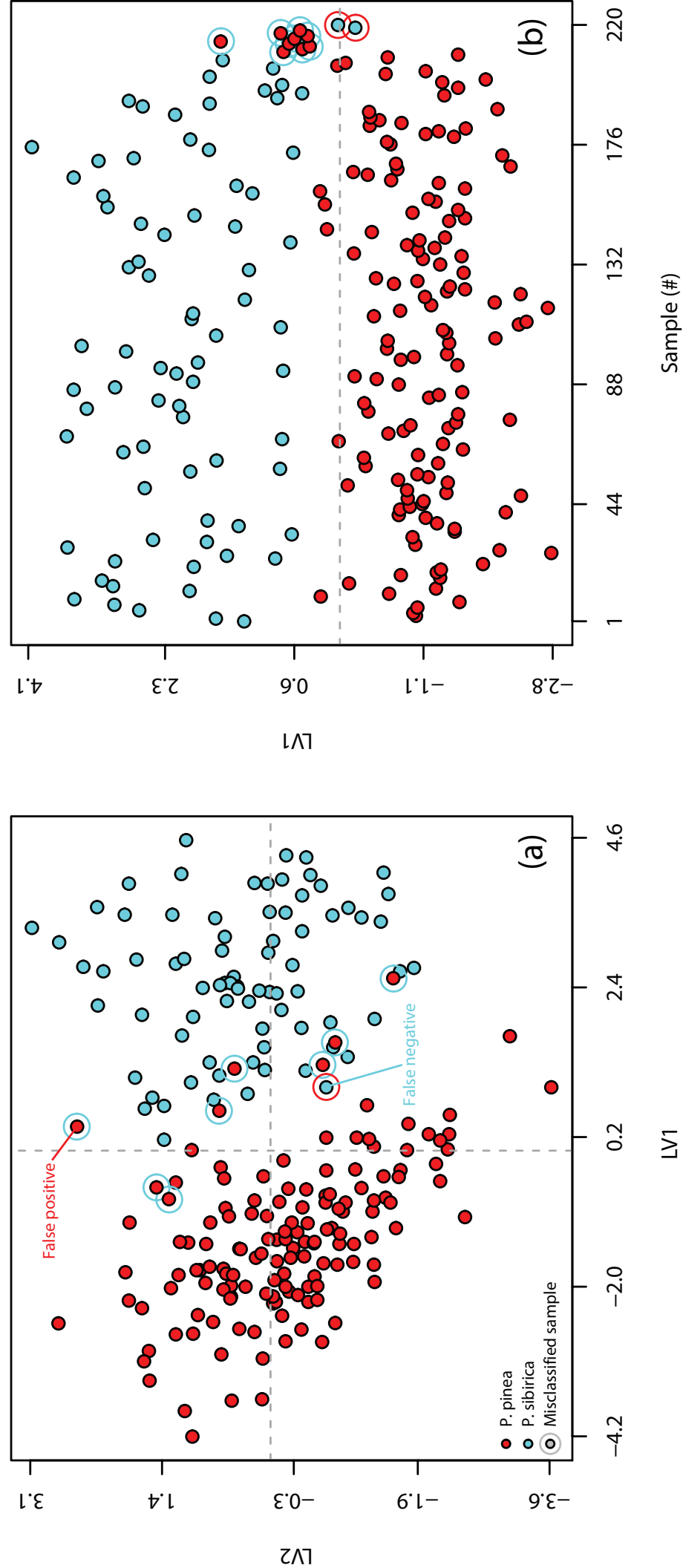
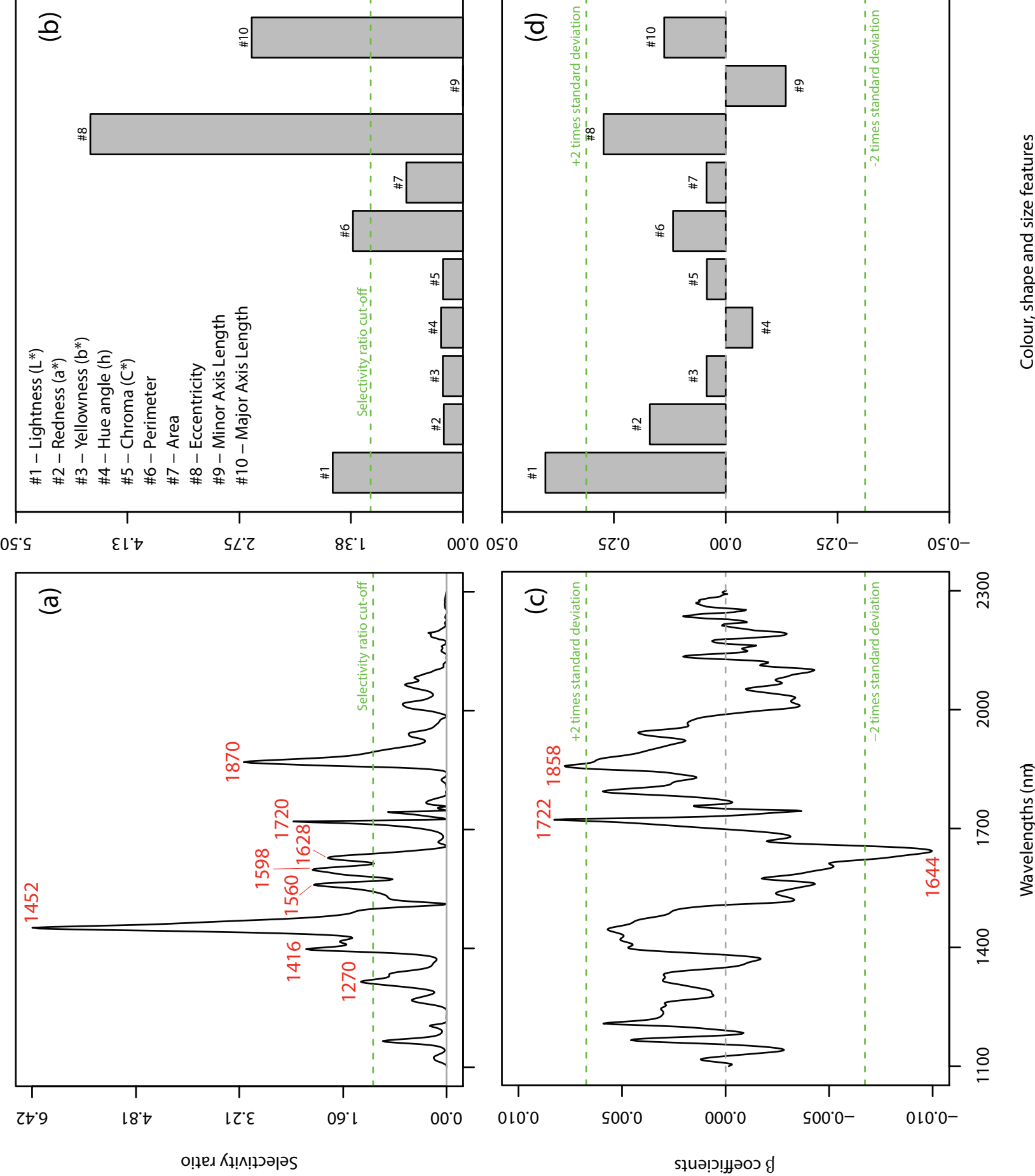


Figure 4



## FIGURE CAPTIONS

- Figure 1.** Mean raw (a) and pre-treated (b) absorbance spectra for both *P. pinea* L. and *P. sibirica* Du Tour species. Spectral pre-treatment consisted of the Standard Normal Variate scatter correction followed by the 1<sup>st</sup> derivative Savitzky-Golay filter with 9-smoothing points. The vertical straight stripes represent the 10-features intervals at 1640-1658 nm, 1720-1738 nm and 1880-1898 nm selected by the iPLS-DA algorithm.
- Figure 2.** Spatial distribution plots obtained from PLS-DA (a) and iPLS-DA (b) models developed using NIR spectral features of the prediction set. Red point with cyan outline corresponds to false positive error (i.e. *P. sibirica* sample erroneously classified as belonging to *P. pinea*). Cyan point with red outline corresponds to false negative error (i.e. *P. pinea* sample erroneously classified as belonging to *P. sibirica*).
- Figure 3.** Spatial distribution plots obtained from PLS-DA (a) and iPLS-DA (b) models developed using imaging features of the prediction set. Red point with cyan outline corresponds to false positive error (i.e. *P. sibirica* sample erroneously classified as belonging to *P. pinea*). Cyan point with red outline corresponds to false negative error (i.e. *P. pinea* sample erroneously classified as belonging to *P. sibirica*).
- Figure 4.** Selectivity ratio plots for the PLS-DA models based on NIR spectral features (a) and imaging features (b). The horizontal green-dashed line corresponds to cut-off ratio between the explained and the residual variance. The larger the selectivity ratio, the more useful the given feature was for the classification task.



TABLE 1

<b>Species</b>	<b>Class</b>	<b>Batch #</b>	<b>Packaging Country</b>
<i>Pinus pinea</i> L.	1	1	Italy
		2	Spain
		3	Unknown
		4	Italy
		5	Italy
<i>Pinus sibirica</i> Du Tour	2	6	China
		7	Russia, Altai
		8	Russia, far east
		9	Russia, Buryatia

TABLE 2

Factor	Minimum	Q1 <sup>a</sup>	Median	Mean		Q3 <sup>b</sup>	Maximum	SE <sup>c</sup>
<b>Lightness (<math>L^*</math>)</b>								
<i>P. pinea</i>	74.09	81.19	82.48	82.23	a	83.59	86.78	0.08
<i>P. sibirica</i>	72.54	77.87	78.66	78.62	b	79.41	83.64	0.08
<b>Red/green colour (<math>a^*</math>)</b>								
<i>P. pinea</i>	-1.88	0.39	0.70	0.71	a	1.05	3.36	0.02
<i>P. sibirica</i>	-1.81	-0.35	0.18	0.20	b	0.73	3.15	0.05
<b>Yellow/blue colour (<math>b^*</math>)</b>								
<i>P. pinea</i>	11.28	15.13	16.32	16.95	b	17.94	30.91	0.12
<i>P. sibirica</i>	11.38	16.68	18.95	19.89	a	23.16	35.34	0.25
<b>Hue angle (<math>h</math>)</b>								
<i>P. pinea</i>	81.89	86.43	87.54	87.58	b	88.64	94.74	0.07
<i>P. sibirica</i>	80.44	87.68	89.46	89.15	a	90.95	94.11	0.14
<b>Chroma (<math>C^*</math>)</b>								
<i>P. pinea</i>	11.29	15.16	16.34	16.97	b	17.97	30.92	0.12
<i>P. sibirica</i>	11.39	16.72	18.97	19.91	a	23.17	35.36	0.25
<b>Perimeter (mm)</b>								
<i>P. pinea</i>	17.83	27.82	29.90	29.94	a	31.99	39.36	0.13
<i>P. sibirica</i>	16.99	20.68	22.31	23.30	b	25.81	34.57	0.20
<b>Surface area (mm<sup>2</sup>)</b>								
<i>P. pinea</i>	21.82	45.83	52.54	52.98	a	60.10	82.56	0.42
<i>P. sibirica</i>	19.98	29.83	34.40	37.52	b	44.24	70.02	0.58
<b>Eccentricity</b>								
<i>P. pinea</i>	0.79	0.89	0.91	0.90	a	0.92	0.96	1.2E-03
<i>P. sibirica</i>	0.55	0.76	0.80	0.79	b	0.83	0.94	3.5E-03
<b>Minor axis length (mm)</b>								
<i>P. pinea</i>	3.73	4.90	5.30	5.32	ns	5.73	7.21	0.03
<i>P. sibirica</i>	3.71	4.85	5.21	5.34	ns	5.78	7.32	0.04
<b>Major axis length (mm)</b>								
<i>P. pinea</i>	7.07	11.69	12.72	12.75	a	13.79	18.02	0.07
<i>P. sibirica</i>	6.62	7.75	8.63	8.93	b	10.15	15.38	0.09

<sup>a</sup> First quartile, <sup>b</sup> Third quartile, <sup>c</sup> Standard error, ns = no significant difference

TABLE 3

METHOD	ALGORITHM	FEATURES	DATA PRE-TREATMENTS			LVs <sup>f</sup>	SENSITIVITY			SPECIFICITY			ACCURACY				
			SC <sup>c</sup>	Savitzky-Golay filter			n.	Variance (%)	C <sup>g</sup>	CV <sup>h</sup>	P <sup>i</sup>	C	CV	P	C	CV	P
				Derivative	Smoothing points												
NIR	PLS-DA <sup>a</sup>	Whole spectrum	SNV <sup>d</sup>	D1 <sup>f</sup>	9	2	54.28	0.96	0.96	0.95	0.99	0.98	0.97	0.97	0.97	0.96	
	iPLS-DA <sup>b</sup>	1) 1640-1658 nm 2) 1720-1738 nm 3) 1880-1898 nm	SNV	D1 <sup>f</sup>	9	1	90.29	0.98	0.98	0.96	0.99	0.99	1.00	0.98	0.98	0.98	
Imaging	PLS-DA	All imaging features				4	83.98	0.95	0.95	0.98	0.97	0.97	0.99	0.96	0.96	0.98	
	iPLS-DA	1) <b>Lightness (L*)</b>				1	52.89	0.95	0.95	0.97	0.97	0.97	0.97	0.96	0.96	0.97	
		2) Red/green (a*)															
		3) Hue angle (h)															
		4) Eccentricity															
		5) Major Axis Length															

<sup>a</sup>PLS-DA, Partial Least Squares Discriminant Analysis;

<sup>b</sup>iPLS-DA, Interval Partial Least Squares Discriminant Analysis;

<sup>c</sup>SC, Scatter Correction method;

<sup>d</sup>SNV, Standard Normal Variate;

<sup>e</sup>SP, Savitzky-Golay smoothing points;

<sup>f</sup>LVs, number of Latent Variables.

<sup>g</sup>C, calibration;

<sup>h</sup>CV, cross-validation;

<sup>i</sup>P, prediction.

## TABLE CAPTIONS

**Table 1.** *List of the batches arranged into classes (i.e. pine nut species), which were used for the experimental activity.*

**Table 2.** *Summary of descriptive statistics of the imaging features of pine nuts from *P. pinea* L. and *P. sibirica* Du Tour species. Mean values belonging to the same factor without common letters are statistically different according to HSD ( $P \leq 0.05$ ).*

**Table 3.** *Summary of performance metrics for classification algorithm (i.e. PLS-DA and iPLS-DA) complexity which gave the best results for both analytical methods used in the experimentation (i.e. NIR spectroscopy and image analysis). The pre-treatments associated to each model were applied in combination.*