1 **UNTARGETED CLASSIFICATION FOR PAPRIKA POWDER AUTHENTICATION**

2 **USING VISIBLE – NEAR INFRARED SPECTROSCOPY (VIS-NIRS)**

3 Olga Monago-Maraña[a*], Carl Emil Eskildsen[a], Teresa Galeano-Díaz[b,c], Arsenio Muñoz de la

4 Peña[b,c], Jens Petter Wold[a]

5 [a]Nofima AS – Norwegian Institute of Food, Fisheries and Aquaculture Research, PB 210, N-

6 1431, Ås, Norway

7 [b]Department of Analytical Chemistry, University of Extremadura, Badajoz 06006, Spain

8 [c]Research Institute on Water, Climate Change and Sustainability (IACYS), University of

9 Extremadura, Badajoz 06006, Spain

10

11 *corresponding author. E-mail: olgamonago@gmail.com

**Abstract**

This paper describes a non-destructive screening method for authentication of paprika belonging to the Spanish Protected Designation of Origin (PDO) *"Pimentón de La Vera"*. Different multivariate classification models were developed in order to differentiate PDO and non-PDO samples, using visible-near infrared spectra as fingerprint for each paprika sample. Sample treatment was not required. Principal component analysis (PCA) was applied in different spectral ranges: 400 - 2500, 400 - 800 and 800 - 2500 nm. In all spectral ranges, PCA was largely able to differentiate PDO from non-PDO samples. Partial least-squares - discriminant analysis (PLS-DA), PCA-linear discriminant analysis (LDA) and PCA-quadratic discriminant analysis (QDA) were used as classification methods in the different spectral ranges. All methods were able to differentiate PDO from non-PDO samples, with error rates (ER) lower than 0.15. The best models were those obtained with PLS-DA in the NIR range (800 - 2500 nm), showing ERs lower than 0.07 and error indexes ($I_{ERROR}$) (false positives) lower than 0.05.

1    **1. Introduction**

2    Paprika powder is used as a spice in many countries. In Spain, there are three traded types of

3    paprika, which differ in their drying process (air, sun and smoke drying). Air-dried paprika, using

4    heated air, is produced mainly in the south-east and central-east of Spain (Murcia), where the high

5    temperature conditions allow peppers to undergo rapid dehydration. Sun-dried paprika are

6    imported from South America and South Africa. Smoked paprika originates from La Vera region,

7    Extremadura in the south-west of Spain. Here, a traditional drying process is used, where oak logs

8    are burnt to heat the paprika to 40 ºC and give it a smoked flavor (Martín et al., 2017).

9    Smoked paprika is recognized under the quality seal Protected Designation of Origin (PDO)

10   *"Pimentón de La Vera"* by the European Union since 2006 (Unión Europea, 2006). This product

11   is considered a high-quality product obtained by drying the fruit of autochthonous varieties of

12   peppers *(Capsicum annum L.).* Moreover, the traditional drying process confers the paprika its

13   aroma, flavor, and color (Martín et al., 2017). Adulteration of smoked paprika *"Pimentón de La*

14   *Vera"* with foreign paprika of lower quality, primarily to increase profit margins, has been a

15   concern for many years  to the smoked paprika industry (Hernández, Martín, Aranda, Bartolomé,

16   & Córdoba, 2007). Therefore, inexpensive and high throughput screening tools to differentiate

17   paprika based on origin is interesting for the industry.

18   Recent reviews show how spectroscopic techniques, including near-infrared spectroscopy

19   (NIRS), can be used for detection of adulteration in herbs and spices (Kucharska-Ambrożej &

20   Karpinska, 2020; Marciano M. Oliveira, Cruz-Tirado, & Barbin, 2019). However, not many

21   studies about paprika powder adulteration were found. In the case of paprika or related products,

22   NIRS has been mainly used for quantification. For example, to quantify ASTA color, moisture

23   (Bae, Han, & Hong, 1998), capsaicinoids (Lim, Kim, Mo, & Kim, 2015; Park et al., 2008), arsenic

24   and lead (Moros et al., 2008), soluble solids content (SSC), firmness of peppers (Penchaiya,

25   Bobelyn, Verlinden, Nicolaï, & Saeys, 2009) and mycotoxins (Hernández-Hierro, García-

26   Villanova, & González-Martín, 2008). In addition, Vis-NIRS combined with multivariate

27   analysis has been used to determine total carotenoids, chlorophylls, as well as maturity stage of

28    intact peppers (Timea Ignat et al., 2013) and ascorbic acid (T. Ignat, Schmilovitch, Fefoldi,

29    Steiner, & Alkalai-Tuvia, 2012). Few works about the adulteration and/or authentication of

30    paprika powder using NIRS as analytical technique have been found in the literature. A recent

31    work about this topic was based on the detection of adulterants such us potato starch, annatto and

32    acacia gum in paprika powder samples from Spain (n = 3) and Brazil (n = 2) (M. M. Oliveira,

33    Cruz-Tirado, Roque, Teófilo, & Barbin, 2020). Detection and quantification of adulterants was

34    done using a portable NIR instrument in combination with partial least squares (PLS) regression

35    and PLS-Discriminant Analysis (PLS-DA). The results were promising with a specificity greater

36    than 90% and error rate lower than 2 % for the PLS-DA models.

37    In another study, paprika samples were clustered based on origin using NIRS and Principal

38    Component Analysis (PCA) (Molnár et al., 2018).  However, only six paprika samples from Spain

39    were included in the analysis, and PDO specifications were not taken into account.

40    Only few studies have investigated the possibility of differencing between paprika samples

41    belonging to the PDO *"Pimentón de La Vera"* and samples not belonging to the PDO.

42    Discrimination has been based on color measurements with visible spectrophotometry, being

43    samples, belonging to the PDO *"Pimentón de La Vera*" or not, correctly grouped in two groups

44    with PCA (Monago Maraña, Bartolomé García, & Galeano Díaz, 2016). Then, samples were

45    classified as different PDOs *("Pimentón de La Vera"* or *"Pimentón de Murcia")* with

46    classification efficiencies ranging from 92 to 95 % when visible spectra and multilayer

47    perceptrons artificial neural networks (MLP-ANN) were used (A. Palacios-Morillo, Jurado,

48    Alcázar, & Pablos, 2016).

49    Regarding to destructive methods, liquid chromatography has been widely used for the paprika

50    authentication. Classification and authentication have been done with different Spanish PDOs,

51    *"Pimentón de La Vera"*, *"Pimentón de Murcia"*, and Czech Republic paprika samples without

52    PDO. Employing ultra-high-performance liquid chromatography coupled with high-resolution

53    mass spectrometry (UHPLC-HRMS), samples were discriminated on a non-target way (Barbosa,

54    Saurina, Puignou, & Núñez, 2020) and based on the polyphenolic and capsaicinoid profiling

55  (Barbosa, Saurina, & Oscar, 2020) with classification results of 100%. On the other hand, HPLC-

56  UV was used to obtain the phenolic profile of paprika for their authentication, confirming that

57  was enough to discriminate between PDOs (Cetó, Sánchez, Serrano, Díaz-Cruz, & Núñez, 2020).

58  Also, the presence or absence of sub-products from the smoking process (Polycyclic Aromatic

59  Hydrocarbons, PAHs) (Monago-Maraña, Galeano-Díaz, & Muñoz de la Peña, 2017),

60  hydrophobic proteins (Hernández et al., 2007) or metallic content (Ana Palacios-Morillo, Jurado,

61  Alcázar, & De Pablos, 2014) have allowed differentiation of paprika at different conditions.

62  Although being very selective, discriminating on these compounds requires sample extraction

63  steps, which normally is time consuming. For this reason, high throughput screening methods are

64  interesting for practical use in the paprika industries.

65  In this study, Vis-NIR measurements will be used, which are cost effective, high throughput and

66  non-destructive, to discriminate paprika powder samples belonging to the PDO *"Pimentón de La*

67  *Vera"* from paprika powder samples not belonging to the PDO. To achieve this goal, we use

68  multivariate qualitative analytical methods for authenticating the PDO *"Pimentón de La Vera"*

69  paprika powder samples. Different methods for classification of multivariate data were compared

70  and ranked.

71  **2. Material and methods**

72  **2.1. Samples**

73  A total of 49 paprika powder samples under the PDO *"Pimentón de La Vera"* were included in

74  the study. These samples were from five different producers and were made over a period of ten

75  years (2010 – 2020). Samples from 2010 to 2017 were obtained in 2017 (n = 35) from producers

76  and measured in that year. Samples from 2017 – 2020 (n = 14) were acquired in Spanish markets

77  in 2020 and measured that year. The samples were made under smoked conditions, following the

78  traditional process from La Vera, in Extremadura, Spain. Among these samples, there were sweet,

79  sweet/hot and hot paprika samples.

80   A total of 50 samples not belonging to any PDO were acquired from different markets in Spain

81   and Norway. Samples acquired in Norway (n = 9) were bought and measured in 2017, but samples

82   acquired in Spanish markets (n = 23) were acquired in 2017 and 2020 (n = 18), and measured the

83   corresponding year of acquisition. The production processes of these samples are unknown as

84   well as the peppers used for their production due to the fact that it is not mandatory to include

85   that information in labels of paprika samples. Among these samples, there were sweet and hot

86   paprika samples.

87   **2.2. Spectroscopic acquisition**

88   The VIS-NIRS measurements were obtained in reflectance mode using a FOSS NIRS Systems

89   XDS Rapid Content™ Analyzer (FOSS Analytical A/S, Hillerød, Denmark). In order to

90   obey *Beer's law*, the NIR spectra were transformed from reflectance (R) units into absorbance-

91   like units (log(1/R)). An internal ceramic standard was used as reference. Spectra were obtained

92   from 400 to 2500 nm, with a resolution of 0.5 nm. Paprika powder samples were measured in

93   circular sample cups of approximately 79 cm$^2$ (FOSS Analytical A/S, Hillerød, Denmark).

94   Spectra from each sample were acquired in triplicate, mixing the powder for obtaining different

95   surfaces each time to obtain a representative sample spectrum. The average spectrum was used

96   for further analysis.

97   **2.3. Data processing and multivariate analysis**

98   **2.3.1. Principal component analysis**

99   Principal component analysis (PCA) was applied to explore the main variation over samples.

100  During PCA all samples were included. Prior to PCA the spectral measurements were

101  preprocessed by extended multiplicative signal corrected (EMSC) (Martens & Stark, 1991) and

102  mean centered variable-wise.

103  The objective of PCA is to compress the data, reducing it from the high dimensional variable

104  space into a lower dimensional principal component space. Each new principal component (PC)

105  is a linear combination of the original variables. The loadings describe the direction of each

principal component in the original X-space and the scores are the projections of the original data onto the loading vectors (Wold, Esbensen, & Geladi, 1987).

PCAs was performed separately for the entire spectral range, the visible range (from 400 to 800 nm) and the NIR (800 - 2500 nm) range.

**2.3.2. Classification analysis**

For the classificatory analysis, samples were divided in two sets (training and test). Approximately 60 % of the samples were used for training and the remaining 40 % of the samples were used for validation. Hence, the training set was composed by 59 samples (29 PDO and 30 non-PDO) and the test set was formed by 40 samples (20 PDO and 20 non-PDO). The split of samples was based on the recently published EuroLab Guide (TR No 01/2015, 2015), which recommends a minimum of 20 samples for each class in the test sets. The training and test samples were randomly chosen. Hence, this division was performed three times, and three different training and test sets were obtained and used for building different calibration models. As a result, the average results of three training and test sets were given with the corresponding standard deviation.

The following classification algorithms were tested for discrimination of the sample spectra: discriminant partial least-squares (PLS-DA) (Barker & Rayens, 2003), linear discriminant analysis based on the PC scores of the spectra (PCA-LDA) (Mohanty, John, Manmatha, & Rath, 2013), and quadratic discriminant analysis based on the PC scores of the spectra (PCA-QDA) (Tharwat, 2016).

PLS-DA involves performing a multivariate regression model to establish class limits and placing a numeric value to each object/sample first, and then classifying them into a specific class. As in PLS regression, the relation between instrumental response in X (spectra) and y (class coding) is established, and the optimal number of latent variables is chosen based on the error range by cross-validation.

131　To apply LDA or QDA, it is necessary to reduce the dimensionality of the spectral data. For that

132　PCA is used. After PCA, LDA is used when the decision line between the two groups can be

133　represented by a linear function. However, if a curved line is needed to separate the groups, then

134　QDA is more effective.

135　Prior to classification the spectral training data were preprocessed by EMSC and variable-wise

136　mean centered. Classification models were fitted on the training set using full-cross validation to

137　determine the optimal models. Then the models were tested with the external test set (pre-

138　processed with the EMSC model obtained for training previously). Data analysis was done using

139　a graphical interface (Ballabio & Consonni, 2013) in Matlab (R2016b, The MathWorks, Inc.,

140　Natick, MA, USA).

141　**2.3.3. Evaluation of the methodology**

142　In order to evaluate the screening methodology, the confusion matrices were obtained and the

143　performance parameters such as precision (PREC), sensitivity (SENS), error rate (ER), accuracy

144　(ACCU) and specificity (SPEC) were calculated.

145　The PREC is defined as the number of samples correctly assigned as belonging to the PDO (i.e.

146　true positives (TP)) over the total number of samples assigned as belonging to the PDO (i.e. the

147　total number of true positives and false positives (FP)) (Eq. 1). The SENS is the number of true

148　positives over the total number of samples belonging to the PDO (i.e. the total number of true

149　positives and false negatives (FN)) (Eq. 2). The ER is the number of samples incorrectly classified

150　by the model (i.e. the total number of false positives and false negatives) over the total number of

151　samples (Eq. 3). The ACCU is the number of samples correctly classified by the model (i.e. the

152　total number of true positives and true negatives (TN)) over the total number of samples (Eq. 4).

153　The SPEC is the number of samples correctly assigned as not belonging to the PDO (i.e. true

154　negatives) over the total number of samples not belonging to the PDO (Eq. 5).

155　$PREC = \frac{TP}{TP+FP}$ (1)

156　$SENS = \frac{TP}{TP+FN}$ (2)

157 $$ER = \frac{FN + FP}{TP+TN+FP+FN} \qquad (3)$$

158 $$ACCU = \frac{TP + TN}{TP+TN+FP+FN} \qquad (4)$$

159 $$SPEC = \frac{TN}{TN + FP} \qquad (5)$$

160 Where TP and TN are the number true positive and number of true negative, respectively, and FN

161 and FP are the number of false negative and number of false positive, respectively.

162 Furthermore, two recently proposed indexes, error index ($I_{ERROR}$) and loss index ($I_{LOSS}$), for

163 assigning a specification-based quality grade for a PDO label are calculated (Cuadros-Rodríguez,

164 Valverde-Som, Jiménez-Carvelo, & Delgado-Aguilar, 2020).

165 $I_{ERROR}$ is the probability of a sample being incorrectly assigned to the PDO class (Eq. 6). $I_{LOSS}$ is

166 the probability of obtaining false negatives and thus the risk of economic loss due to assignment

167 error.

168 $$I_{ERROR} = \frac{FP}{TP+TN+FP+FN} \qquad (6)$$

169 $$I_{LOSS} = \frac{FN}{TP+TN+FP+FN} \qquad (7)$$

170 **3. Results and discussion**

171 **3.1. VIS-NIRS spectral profiling**

172 Figure 1A shows the mean of the absorption spectra for both classes (PDO and non-PDO). The

173 mean spectrum of non-PDO shows higher intensity over the whole spectral range as compared

174 with the mean spectrum for PDO. More subtle differences can be seen after pre-processing by

175 EMSC (Figure 1B). The main difference in the visible range was observed at 670 nm, and in the

176 NIR range at 1450, 1940, 2305, 2346 and 2490 nm. The visible range was previously reported to

177 be useful for the quantification of total carotenoids and chlorophylls in intact bell pepper (Timea

178 Ignat et al., 2013). In the case of NIR bands, some of them might be due to water peaks (1450 and

179 1940 nm) and the other three main peaks (2305, 2350 and 2490 nm) do most likely originate from

180 fat (Núñez-Sánchez et al., 2016).

## 3.2. Exploratory analysis

In order to study the most important spectral variation for discriminating PDO and non-PDO samples, detect potential outliers and systematic artifacts in the samples, PCA was performed on the EMSC pre-processed spectra. All 99 samples were included in the analysis. As described above, PCAs were performed on different spectral ranges.

When including the whole spectral range, the first three principal components (PCs) explain 84 % of the total variation in the data set. The first principal component (PC1) explains 50 % of the variation, and the corresponding loading plot (not shown) reveals the most important peaks at approximately 480 and 600 nm in the visible range and at 1450 and 1940 nm in the NIR range (water peaks). However, this component does not differentiate PDO from non-PDO paprika samples.

The best discrimination is observed for scores of PC3 and PC5, explaining 12 and 4 % of the total variation, respectively (Figure 2A). Clearly, two groups are established according to PDO and non-PDO samples. However, the two groups are slightly overlapping. PC3 provides the clearest discrimination of the groups. The clear unsupervised clustering is a good basis for supervised classification.

The loadings for PC3 and PC5 are presented in Figure 2B. The main variables affecting the separation of the groups were 540 and 670 nm in the visible range and water peaks in the NIR range (Figure 2A). Score values for PC3 are generally high for the PDO samples, which means that positive loadings, representing certain chemical components, are positively related to PDO samples. The negative loadings observed at 1720 and 1760 nm are related with first overtone C-H stretching vibration of methyl ($-CH_3$), methylene ($-CH_2$) and ethenyl ($-CH=CH-$) groups. The loadings close to 1725 nm has been related to oleic acid and the band close to 1760 nm to saturated components. The bands at 2305 and 2350 nm have previously been assigned to combination of C-H stretches and deformations (Núñez-Sánchez et al., 2016; Pérez-Juan et al., 2010). Also, the small band at 1207 nm is related with fat. All bands related to fat are negative loadings, suggesting a relatively low concentration of fat in PDO samples.

208  Scores for PCA in the visible spectral range are presented in Figure 2C. PC4, explaining 6 % of

209  the variance, discriminates quite well between the two groups. Note that the overlap of the groups

210  is stronger when using only the visible range, compared to using the whole range. The main

211  variables affecting the clustering are those mentioned before (570 and 670 nm) as seen in the

212  loading for PC4 (Figure 2D).

213  Finally, for the NIR range, a quite good grouping of the samples is obtained in PC2 (Figure 2E)

214  due to variables corresponding to water and fat peaks. Interestingly, some peaks are more

215  pronounced in the loadings in this case. These peaks can be attributed to proteins bands: 2056 nm

216  (N-H stretching vibrations) and 2478 nm (-C-N-C stretching first overtone).

### 3.2. Classificatory analysis

218  As detailed in the section 2.4.2, samples were divided into training and test sets. This step was

219  performed three times and the classification model was obtained for each case. Average results

220  for confusion matrices from different sets and the corresponding validation parameters are shown

221  in Tables 1 and 2, respectively. The numbers in parentheses correspond with the standard

222  deviations from the three sets assayed.

223  For PLS-DA, the best classification results were obtained for the NIR range in both training and

224  test samples. The ERs obtained for this range were overall lower than for other ranges.

225  Interestingly, from a quality-point of view, the $I_{ERROR}$ was lower for the NIR spectral range as

226  compared with the other spectral ranges, for both the training set and test set. This is important

227  for avoiding non-PDO samples being classified as PDO samples. The visible range gave slightly

228  less correct classifications than the whole range, but all models provided acceptable results, with

229  ERs lower than 0.11 and $I_{ERROR}$ lower than 0.10. According to (Cuadros-Rodríguez et al., 2020),

230  a good screening method should offer an $I_{ERROR}$ equal to or lower than 0.1 in order to minimize

231  the false-compliance error. Hence, the best choice with PLS-DA would be with the NIR range,

232  although in some cases that means that some samples would be false-negative and refused

233  categorized as PDO (PDO samples categorized as non-PDO samples).

234   Regarding the other performance parameters, SENS and SPEC present similar values (Table 2),

235   mainly in the NIR range. This means that the error is balanced, and there is not a clear trend in

236   the models for false positives, or vice versa. PREC values were higher for the NIR range, which

237   means that false positives were lower in these models, as observed in the $I_{ERROR}$ values as well.

238   The regression coefficients for each spectral range (Figure 3) were evaluated in order to elucidate

239   the main variables contributing to the classification. For the visible range, the main variables were

240   570 and 670 nm with negative values, and 540 nm with positive value. It might be expected that

241   the variation in the visible range would be related to total carotenoids, ASTA values (extractable

242   color), as other authors reported (A. Palacios-Morillo et al., 2016). In these samples, the ASTA

243   value was not so relevant since some PDO samples were old and therefore had low ASTA values

244   (between 25 - 70). Therefore, it was expected that some samples were incorrectly classified when

245   using the visible range. However, acceptable results for classification were obtained due to other

246   variables, not related to total carotenoids. The VIP scores (not shown) were also investigated.

247   Similar information was retrieved from the VIP scores and the regression vectors (Figure 3).

248   The absorption around 670 nm has previously been related with chlorophylls (Timea Ignat et al.,

249   2013) and could be also related with pheophytins formed from chlorophylls during ripening or

250   drying process (Bonaccorsi et al., 2016) . This peak has negative regression coefficients (Figure

251   3A and 3B), which suggests that non-PDO samples have lower content of chlorophyll compared

252   to PDO samples. This is also observed in Figure 1B.

253   Regarding the NIR range, the regression coefficient positive wavelength bands associated with

254   fat, such as, 1725, 2305, 2350 and 2490 nm, again suggesting a relatively high fat content in non-

255   PDO samples. A higher fat content can have different reasons. Different types of peppers used

256   for paprika production vary in the fatty acid composition depending on genotype and

257   environmental factors. Kim et al., 2019 recently reported this for some varieties of peppers and

258   this could be extended to other kind of peppers (Kim et al., 2019). Another reason may be related

259   with the addition of sunflower vegetal oil to give stronger brightness of the powder. In the case

260   of PDO *"Pimentón de La Vera"* the amount of oil is limited to 3 % (w/w) (Unión Europea, 2006).

261 However, there are not specifications reported about other kind of paprika samples, which are not

262 under the PDO. This could mean that other paprika samples contain a higher percentage of

263 sunflower oil to give more brightness. A third reason could be related to the addition of seeds

264 from peppers used in the paprika production, which would influence in the fatty acid composition.

265 This kind of addition is not allowed in PDO samples (Unión Europea, 2006).

266 PCA-LDA and PCA-QDA gave results in accordance with PLS-DA; better results were obtained

267 when the NIR range or whole range were used to classify samples, giving ERs lower than 0.15

268 and $I_{ERROR}$ lower than 0.11. Another important result was that PCA-QDA offered better results

269 than PCA-LDA in all cases. In the case of PCA-LDA and PCA-QDA, PREC, SENS and SPEC

270 values were slightly better for the NIR range. As in previous case, SENS and SPEC values were

271 similar, which proved that errors did not follow a clear trend.

272 Finally, it must be highlighted that these good results were obtained for three training/test sets,

273 which proved the robustness of the methods. To our knowledge, this is the first work where non-

274 destructive classification of PDO *"Pimentón de La Vera"* has been performed. The method is

275 easy and quick to use and could with some more development contribute to effective control in

276 the paprika industries.

277 **4. Conclusions**

278 Vis-NIR spectroscopy with different multivariate classification techniques have been proven to

279 discriminate between paprika samples belonging to the PDO *"Pimentón de La Vera"* and other

280 paprika samples. The variability of samples and the random choice of samples for training and

281 test, indicate that the models are quite robust. The visible range offered the good classification

282 due to chlorophylls or pheophytin compounds and NIR range showed slightly better classification

283 based on differences in absorbance of fat. PLS-DA offered somewhat better results than other

284 classification methods. It can be highlighted that all methods offered acceptable ERs and $I_{ERROR}$,

285 always lower than 0.15 and 0.11, respectively. This method is easy, rapid and non-destructive,

286 being an advantage in order to implement the method for industrial purposes.

287

**Acknowledgements**

297

298    **References**

299    Bae, M.-J., Han, E.-S., & Hong, S.-H. (1998). Use of near infrared spectroscopy in quality
300        control of red pepper powder. *Journal of Near Infrared Spectroscopy*, *6*, A333–
301        A337.

302    Ballabio, D., & Consonni, V. (2013). Classification tools in chemistry. Part 1: linear
303        models. PLS-DA. *Analytical Methods*, *5*, 3790–3978.

304    Barbosa, S., Saurina, J., & Oscar, N. (2020). Capsaicinoid profiling for the chemometric
305        characterization and classification of Paprika with Protected Designation of Origin
306        (PDO) attributes. *Molecules*, *25*, 1–16.

307    Barbosa, S., Saurina, J., Puignou, L., & Núñez, O. (2020). Classification and
308        authentication of paprika by UHPLC-HRMS fingerprinting and multivariate
309        calibration methods (PCA and PLS-DA). *Foods*, *9*, 1–10.

310    Barker, M., & Rayens, W. (2003). Partial least squares for discrimination. *Journal of
311        Chemometrics*, *17*, 166–173.

312    Bonaccorsi, I., Cacciola, F., Utczas, M., Inferrera, V., Giuffrida, D., Donato, P., …
313        Mondello, L. (2016). Characterization of the pigment fraction in sweet bell peppers
314        (Capsicum annuum L.) harvested at green and overripe yellow and red stages by
315        offline multidimensional convergence chromatography/liquid chromatography–
316        mass spectrometry. *Journal of Separation Science*, *39*(17), 3281–3291.
317        https://doi.org/10.1002/jssc.201600220

318    Cetó, X., Sánchez, C., Serrano, N., Díaz-Cruz, J. M., & Núñez, O. (2020). Authentication
319        of paprika using HPLC-UV fingerprints. *LWT - Food Science and Technology*, *124*,
320        109153.

321    Cuadros-Rodríguez, L., Valverde-Som, L., Jiménez-Carvelo, A. M., & Delgado-Aguilar,
322        M. (2020). Validation requirements of screening analytical methods based on
323        scenario-specified applicability indicators. *TrAC - Trends in Analytical Chemistry*,
324        *122*.

325    Hernández-Hierro, J. M., García-Villanova, R. J., & González-Martín, I. (2008). Potential
326        of near infrared spectroscopy for the analysis of mycotoxins applied to naturally
327        contaminated red paprika. *Analytica Chimica Acta*, *2*, 189–194.

328    Hernández, A., Martín, A., Aranda, E., Bartolomé, T., & Córdoba, M. de G. (2007).
329        Application of temperature-induced phase partition of proteins for the detection of
330        smoked paprika adulteration by free zone capillary electrophoresis (FZCE). *Food*

15

*Chemistry*, *105*, 1219–1227.

Ignat, T., Schmilovitch, Z., Fefoldi, J., Steiner, B., & Alkalai-Tuvia, S. (2012). Non-destructive measurement of ascorbic acid content in bell peppers by VIS-NIR and SWIR spectrometry. *Postharvest Biology and Technology*, *74*, 91–99.

Ignat, Timea, Schmilovitch, Z., Feföldi, J., Bernstein, N., Steiner, B., Egozi, H., & Hoffman, A. (2013). Nonlinear methods for estimation of maturity stage, total chlorophyll, and carotenoid content in intact bell peppers. *Biosystems Engineering*, *114*, 414–425.

Kim, E. H., Lee, S. Y., Baek, D. Y., Park, S. Y., Lee, S. G., Ryu, T. H., … Oh, S. W. (2019). A comparison of the nutrient composition and statistical profile in red pepper fruits (Capsicums annuum L.) based on genetic and environmental factors. *Applied Biological Chemistry*, (1), 62–48.

Kucharska-Ambrożej, K., & Karpinska, J. (2020). The application of spectroscopic techniques in combination with chemometrics for detection adulteration of some herbs and spices. *Microchemical Journal*, *153*, 104278.

Lim, J., Kim, G., Mo, C., & Kim, M. (2015). Design and fabrication of a real-time measurement system for the capsaicinoid content of Korean red pepper (Capsicum annuum L.) powder by visible and Near-Infrared Spectroscopy. *Journal of Biosystems Engineering*, *15*, 47–60.

Martens, H., & Stark, E. (1991). Extended multiplicative signal correction and spectral interference subtraction: New preprocessing methods for near infrared spectroscopy. *Journal of Pharmaceutical and Biomedical Analysis*, *9*(8), 625–635.

Martín, A., Hernández, A., Aranda, E., Casquete, R., Velázquez, R., Bartolomé, T., & Córdoba, M. G. (2017). Impact of volatile composition on the sensorial attributes of dried paprikas. *Food Research International*, *100*, 691–697.

Mohanty, N., John, A. L. S., Manmatha, R., & Rath, T. M. (2013). Shape-based image classification and retrieval. *Handbook of Statistics*, *31*, 249–267. https://doi.org/10.1016/B978-0-444-53859-8.00010-2

Molnár, H., Kónya, É., Zalán, Z., Bata-Vidács, I., Tömösközi-Farkas, R., Székács, A., & Adányi, N. (2018). Chemical characteristics of spice paprika of different origins. *Food Control*, *83*, 54–60. https://doi.org/10.1016/j.foodcont.2017.04.028

Monago-Maraña, O., Galeano-Díaz, T., & Muñoz de la Peña, A. (2017). Chemometric Discrimination Between Smoked and Non-Smoked Paprika Samples. Quantification of PAHs in Smoked Paprika by Fluorescence-U-PLS/RBL. *Food Analytical*

365       *Methods*, *10*, 1128–1137.

366 Monago Maraña, O., Bartolomé García, T. de J., & Galeano Díaz, T. (2016).
367       Characterization of Spanish Paprika by Multivariate Analysis of Absorption and
368       Fluorescence Spectra. *Analytical Letters*, *49*, 1184–1197.

369 Moros, J., Llorca, I., Cervera, M. L., Pastor, A., Garrigues, S., & de la Guardia, M. (2008).
370       Chemometric determination of arsenic and lead in untreated powdered red paprika
371       by diffuse reflectance near-infrared spectroscopy. *Analytica Chimica Acta*, *613*,
372       196–206.

373 Núñez-Sánchez, N., Martínez-Marín, A. L., Polvillo, O., Fernández-Cabanás, V. M.,
374       Carrizosa, J., Urrutia, B., & Serradilla, J. M. (2016). Near Infrared Spectroscopy
375       (NIRS) for the determination of the milk fat fatty acid profile of goats. *Food*
376       *Chemistry*, *190*, 244–252.

377 Oliveira, M. M., Cruz-Tirado, J. P., Roque, J. V., Teófilo, R. F., & Barbin, D. F. (2020).
378       Portable near-infrared spectroscopy for rapid authentication of adulterated paprika
379       powder. *Journal of Food Composition and Analysis*, *87*, 103403.

380 Oliveira, Marciano M., Cruz-Tirado, J. P., & Barbin, D. F. (2019). Nontargeted analytical
381       methods as a powerful tool for the authentication of spices and herbs: a review.
382       *Comprehensive Reviews in Food Science and Food Safety*, *18*, 670–689.

383 Palacios-Morillo, A., Jurado, J. M., Alcázar, A., & Pablos, F. (2016). Differentiation of
384       Spanish paprika from Protected Designation of Origin based on color measurements
385       and pattern recognition. *Food Control*, *62*, 243–249.

386 Palacios-Morillo, Ana, Jurado, J. M., Alcázar, Á., & De Pablos, F. (2014). Geographical
387       characterization of Spanish PDO paprika by multivariate analysis of multielemental
388       content. *Talanta*, *128*, 15–22.

389 Park, T. S., Candidate, P. D., Bae, Y. M., Researcher, S., Sim, M. J., & Student, G. (2008).
390       Analysis of Capsaicinoids from Hot Red Pepper Powder by Near-Infrared
391       Spectroscopy. *ASABE Annual International Meeting*, (January 2008), 1–7.
392       https://doi.org/10.13031/2013.25077

393 Penchaiya, P., Bobelyn, E., Verlinden, B. E., Nicolaï, B. M., & Saeys, W. (2009). Non-
394       destructive measurement of firmness and soluble solids content in bell pepper using
395       NIR spectroscopy. *Journal of Food Engineering*, *94*, 267–273.

396 Pérez-Juan, M., Afseth, N. K., González, J., Díaz, I., Gispert, M., Furnols, M. F. i., …
397       Realini, C. E. (2010). Prediction of fatty acid composition using a NIRS fibre optics
398       probe at two different locations of ham subcutaneous fat. *Food Research*

399      *International*, *43*(5), 1416–1422.

400    Tharwat, A. (2016). Linear vs. quadratic discriminant analysis classifier: a tutorial.

401      *International Journal of Applied Pattern Recognition*, *3*(2), 145.

402      https://doi.org/10.1504/ijapr.2016.079050

403    TR No 01/2015. (2015). *Guide to NMR Method Development and Validation – Part II :*

404      *Multivariate data analysis*.

405    Unión Europea, U. (2006). *Reglamento (CE) Nº 510/2006 del Consejo. S. Diario Oficial*

406      *de la Unión Europea* (Vol. C 287/2).

407    Wold, S., Esbensen, K. I. M., & Geladi, P. (1987). Principal Component Analysis.

408      *Chemometrics and Intelligent Laboratory Systems*, *2*, 37–52.

409

411 **Figure 1. (A)** Average of absorption spectra (B) Average of EMSC pre-processed spectra. Black

412 lines correspond to the PDO samples and red lines correspond to the non-PDO samples.
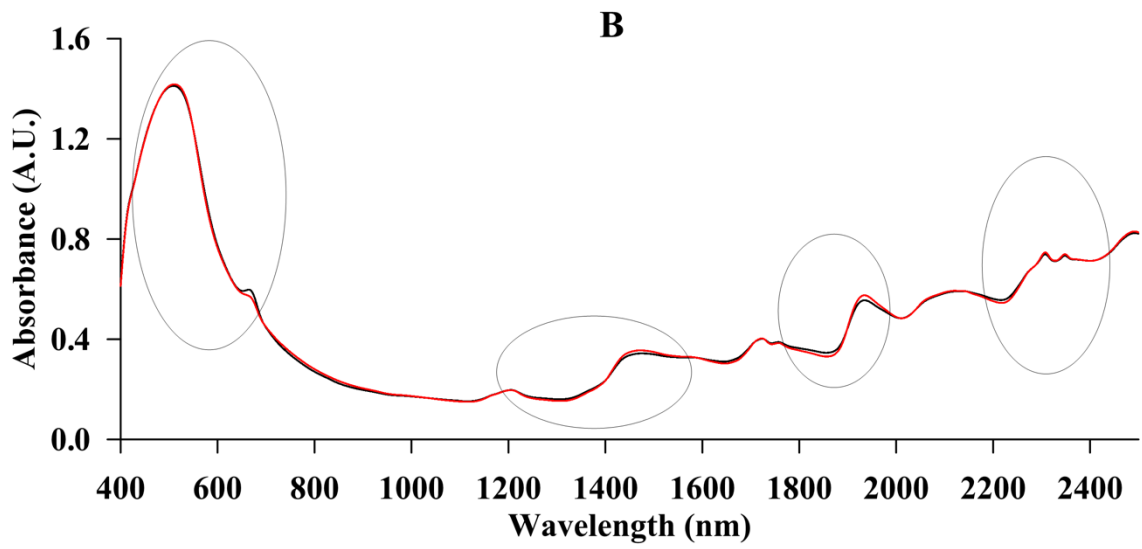
413

414 **Figure 2.** Loadings (B, D, F) and scores values (A, C, E) obtained from PCA of the spectra in

415 wavelength ranges: 400 - 2500 nm, 400 - 800 nm and 800 - 2500 nm.

416

417 **Figure 3.** Regression coefficients for non-PDO samples obtained for the PLS-DA models for the

418 different spectral ranges studied.

419

420

421

422                                      **Figure 1**

**A** Range: 400 - 2500 nm **B**

**C** Range: 400 - 800 nm **D**

**E** Range: 800 - 2500 nm **F**

423

424

425 **Figure 2**

Range: 400 - 2500 nm

Range: 400 - 800 nm

Range: 800 - 2500 nm

426

427

428　　　　　　　　　　　　　　**Figure 3**
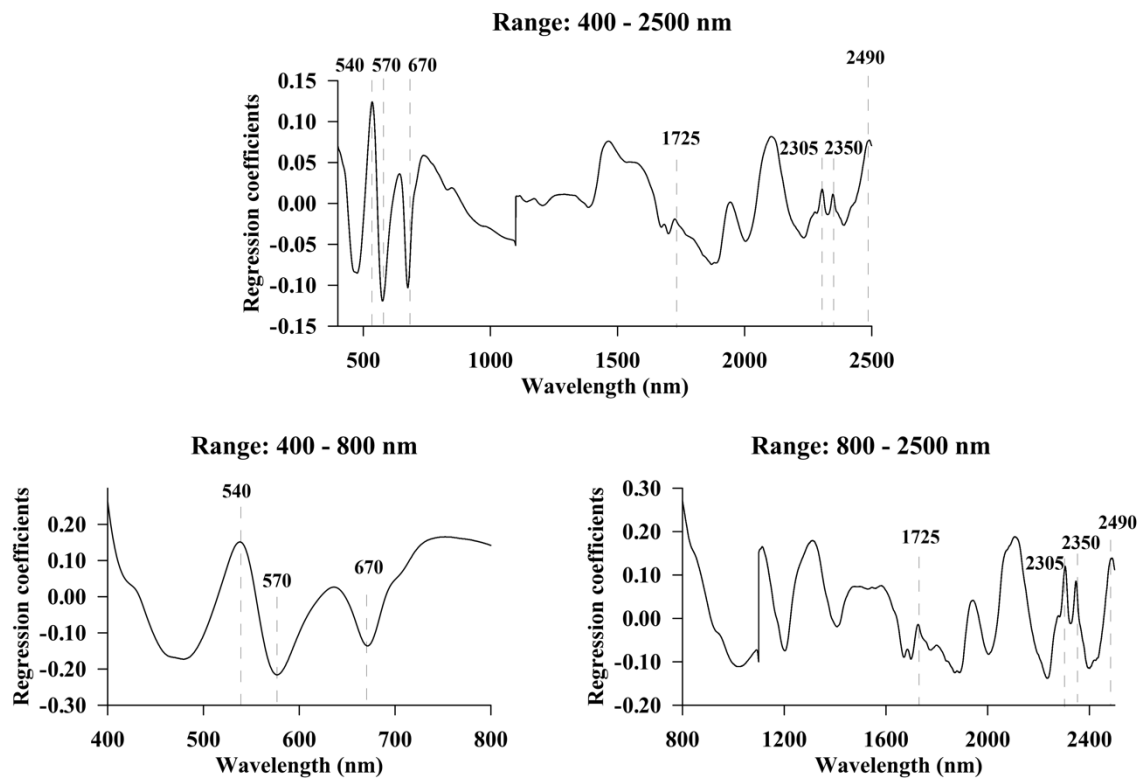
**Table 1.** Confusion matrices for the different algorithms and ranges studied in the training and test sets.

| Algorithm | Range (nm) | Nº comp | %EV (X) | | PDO (CV) | NON- PDO (CV) | PDO (val) | NON-PDO (val) |
|---|---|---|---|---|---|---|---|---|
| | | | | | **Training set** | | **Test set** | |
| **PLS-DA** | 400 - 2500 | 6 | 96 (1) | **PDO** | 28 (1) | 1 (1) | 19 (1) | 1 (1) |
| | | | | **NON-PDO** | 5 (3) | 25 (3) | 2 (2) | 18 (2) |
| | 400 - 800 | 5 | 99 (0) | **PDO** | 28 (1) | 1 (1) | 19 (1) | 1 (1) |
| | | | | **NON-PDO** | 6 (2) | 24 (2) | 3 (1) | 17 (1) |
| | 800 - 2500 | 6 | 98 (1) | **PDO** | 28 (0) | 1 (0) | 19 (1) | 1 (1) |
| | | | | **NON-PDO** | 3 (2) | 27 (2) | 1 (2) | 19 (2) |
| **PCA-LDA** | 400 - 2500 | 5 | 96 (0) | **PDO** | 27 (1) | 2 (1) | 19 (1) | 1 (1) |
| | | | | **NON-PDO** | 7 (2) | 23 (2) | 2 (2) | 18 (2) |
| | 400 - 800 | 5 | 99 (0) | **PDO** | 28 (1) | 1 (1) | 19 (1) | 1 (1) |
| | | | | **NON-PDO** | 6 (2) | 24 (2) | 3 (1) | 17 (1) |
| | 800 - 2500 | 5 | 98 (1) | **PDO** | 25 (2) | 4 (2) | 17 (2) | 3 (2) |
| | | | | **NON-PDO** | 5 (2) | 25 (2) | 2 (2) | 18 (2) |
| **PCA-QDA** | 400 - 2500 | 5 | 97 (1) | **PDO** | 27 (1) | 2 (1) | 17 (2) | 2 (2) |
| | | | | **NON-PDO** | 3 (1) | 27 (1) | 2 (2) | 18 (2) |
| | 400 - 800 | 5 | 99 (0) | **PDO** | 26 (2) | 3 (2) | 18 (1) | 2 (1) |
| | | | | **NON-PDO** | 3 (1) | 27 (1) | 2 (2) | 18 (2) |
| | 800 - 2500 | 5 | 96 (1) | **PDO** | 26 (1) | 3 (1) | 16 (3) | 4 (3) |
| | | | | **NON-PDO** | 2 (1) | 28 (1) | 2 (2) | 18 (2) |

*CV: cross-validation; numbers in parentheses correspond to the standard deviation of three sets assayed.

429

430

**Table 2.** Validation parameters calculated for the target class (PDO class) in the different classification methods.

| Algorithm | Range (nm) | Training set | | | | | | | Test set | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | SPEC | SENS | PREC | ER | ACCUR | $I_{ERROR}$ | $I_{LOSS}$ | SPEC | SENS | PREC | ER | ACCUR | $I_{ERROR}$ | $I_{LOSS}$ |
| PLS-DA | 400 - 2500 | 0.85 (0.11) | 0.98 (0.02) | 0.86 (0.09) | 0.09 (0.06) | 0.91 (0.06) | 0.08 (0.05) | 0.01 (0.05) | 0.92 (0.10) | 0.97 (0.03) | 0.93 (0.09) | 0.06 (0.06) | 0.94 (0.05) | 0.04 (0.05) | 0.02 (0.01) |
| | 400 - 800 | 0.81 (0.08) | 0.97 (0.04) | 0.83 (0.07) | 0.11 (0.05) | 0.89 (0.06) | 0.10 (0.04) | 0.02 (0.02) | 0.87 (0.06) | 0.97 (0.06) | 0.88 (0.05) | 0.08 (0.06) | 0.92 (0.06) | 0.07 (0.03) | 0.02 (0.03) |
| | 800 - 2500 | 0.90 (0.07) | 0.97 (0.0) | 0.91 (0.06) | 0.07 (0.04) | 0.93 (0.04) | 0.05 (0.03) | 0.02 (0.00) | 0.93 (0.08) | 0.97 (0.03) | 0.94 (0.07) | 0.05 (0.04) | 0.95 (0.04) | 0.03 (0.04) | 0.02 (0.01) |
| PCA-LDA | 400 - 2500 | 0.78 (0.07) | 0.94 (0.02) | 0.80 (0.05) | 0.14 (0.04) | 0.86 (0.04) | 0.11 (0.03) | 0.03 (0.00) | 0.88 (0.08) | 0.97 (0.06) | 0.89 (0.06) | 0.08 (0.04) | 0.92 (0.04) | 0.06 (0.04) | 0.02 (0.03) |
| | 400 - 800 | 0.80 (0.06) | 0.97 (0.04) | 0.82 (0.05) | 0.12 (0.03) | 0.88 (0.03) | 0.10 (0.04) | 0.02 (0.01) | 0.87 (0.06) | 0.97 (0.06) | 0.88 (0.05) | 0.08 (0.06) | 0.92 (0.06) | 0.07 (0.03) | 0.02 (0.03) |
| | 800 - 2500 | 0.82 (0.05) | 0.87 (0.05) | 0.82 (0.05) | 0.15 (0.04) | 0.85 (0.04) | 0.09 (0.03) | 0.04 (0.02) | 0.90 (0.09) | 0.87 (0.10) | 0.90 (0.07) | 0.12 (0.06) | 0.88 (0.06) | 0.05 (0.04) | 0.07 (0.05) |
| PCA-QDA | 400 - 2500 | 0.91 (0.02) | 0.92 (0.02) | 0.91 (0.02) | 0.08 (0.01) | 0.92 (0.00) | 0.04 (0.01) | 0.04 (0.01) | 0.92 (0.08) | 0.87 (0.10) | 0.92 (0.07) | 0.11 (0.02) | 0.89 (0.01) | 0.04 (0.04) | 0.07 (0.05) |
| | 400 - 800 | 0.90 (0.03) | 0.89 (0.05) | 0.90 (0.04) | 0.11 (0.04) | 0.89 (0.04) | 0.04 (0.02) | 0.05 (0.03) | 0.92 (0.08) | 0.90 (0.05) | 0.92 (0.07) | 0.09 (0.02) | 0.91 (0.01) | 0.05 (0.05) | 0.06 (0.01) |
| | 800 - 2500 | 0.92 (0.02) | 0.90 (0.04) | 0.92 (0.02) | 0.09 (0.03) | 0.91 (0.03) | 0.04 (0.01) | 0.05 (0.02) | 0.92 (0.10) | 0.78 (0.16) | 0.92 (0.10) | 0.15 (0.07) | 0.85 (0.07) | 0.04 (0.05) | 0.11 (0.08) |

Numbers in parentheses correspond to the standard deviation of three sets assayed.

431

432

24