# Raman spectroscopy coupled with chemometric modeling approaches for authentication of different paprika varieties at physiological maturity

Stefan Kolašinac [a,*], Ilinka Pećinar [a], Dario Danojević [b], Zora Dajić Stevanović [a]

[a] *Faculty of Agriculture, Department of Agricultural Botany, Nemanjina 6, 11000, Belgrade-Zemun, Serbia*
[b] *Institute of Field and Vegetable Crops, Maksima Gorkog 30, 21000, Novi Sad, Serbia*

## ARTICLE INFO

## ABSTRACT

Five Balkan paprika varieties at physiological maturity were investigated by means of Raman spectroscopy in order to discriminate the differences which stemmed from their genetic variability since the plants were grown under the same experimental conditions. The spectra were obtained using the 532 nm wavelength. In an effort to find the best classification power, several pre-processing methods were applied: 1) baseline correction, unit vector normalization; 2) baseline correction, unit vector normalization and first Savitzky-Golay derivative; 3) baseline correction, unit vector normalization and second Savitzky-Golay derivative; 4) baseline correction, unit vector normalization and third Savitzky-Golay derivative. All of the pre-processing methods were followed by making PCA-LDA (Principal Component Analysis-Linear Discriminant Analysis), QDA (Quadratic Discriminant Analysis), and PLS-DA (Partial Least Square - Discriminant Analysis) classification models. QDA showed the best discrimination power (83.87–100% and 89.47–100% for the training and the test data, respectively), then PCA-LDA (0.00–100 and 0.00–100% for the training and the test data, respectively) and PLS-DA (19.35–100% and 0.00–100.00% for the training and the test data, respectively). The results pointed out the applicability of chemometric modeling associated with Raman spectroscopy in the assessment of nutritionally similar samples, such as the studied red paprika varieties.

## 1. Introduction

Paprika (*Capsicum* spp.) is among the most cultivated vegetable crops. According to the recent FAO reports, the global production of fresh and dried green chilies and peppers in 2019 was estimated at 38,024,154 and 4,255,050 tonnes, respectively (www.fao.org). Paprika is characterized by a high number of varieties differing in shape and color, starting from green, yellow or white (unripe fruits), and turning to red, dark red or brown. Significant quantities of red paprika are produced in Hungary, Serbia, Croatia and North Macedonia as it is one of the favorite and most frequently used vegetables in traditional cuisine (Vinković et al., 2018). Because of its high nutritive value, and especially due to the presence of different bioactive compounds, mainly fibers, L-ascorbic acid (vitamin C), carotenoids and polyphenols, sweet paprika is recommended for daily consumption in human diet (Greco, Riccio, Bergero, Del Re, & Trevisan; Kolašinac, Dajić-Stevanović, Kilibarda, & Kostić, 2021). Different carotenoids of paprika are responsible for different fruit coloration, as well as for the significant health-related beneficial effects, which are mainly attributed to capsanthin,

capsorubin, beta carotene and lutein (Deli, Molnár, Matus, & Tóth, 2001). During the ripening period, the color of red paprika transforms from green to deep red due to the accumulation and transformation of carotenoids, which mostly refers to the conversion of lutein and beta carotene into zeaxanthin and capsanthin (Deli et al., 2001).

The level of bioactive compounds in paprika is influenced by the ripening stage, the genotype and cultivation practice (Hallmann & Rembiałkowska, 2012). The content of polyphenols, carotenoids and ascorbic acid increases during the fruit maturation reflecting increased antioxidant activity (Kim, Ahn, Ha, Rhee, & Kim, 2011).

Different analytical tools are used for food quality assessment, as well as for the confirmation of the geographic origin of target products and their possible adulteration. Regarding the quality assessment of different paprika varieties and their products, a standard technique and some advanced techniques are used, including NMR (Nuclear magnetic resonance) (Ramírez-Meraz et al., 2020), HPLC (High-performance liquid chromatography) (Cetó et al., 2018), VIS-NIRS (Monago-Maraña, Eskildsen, Galeano-Díaz, Muñoz de la Peña, & Wold, 2021), HPLC-FLD (High-performance liquid chromatography with fluorescence

detection) (Campmajó, Rodríguez-Javier, Saurina, & Núñez, 2021), UHPLC (Mudrić et al., 2017), FT-MIR (Fourier transform mid-infrared) (Horn, Esslinger, Pfister, Fauhl-Hassek, & Riedl, 2018), UHPLC–APCI–HRMS (ultra-high-performance liquid chromatography coupled to high-resolution mass spectrometry using atmospheric pressure chemical ionization) (Arrizabalaga-Larrañaga et al., 2021), and some others.

Today, there is a great need for the application of non-destructive and rapid analytical methods in food quality control, as well as food product authentication and adulteration, since the standard instrumental methods are high-priced, time-consuming and require special sample preparation. In general, all of these methods target the detection of a narrow group of chemical compounds (such as proteins, carbohydrates, lipids, polyphenols, etc.), as well as some specific (target) compounds or even the isotopic forms of an element (such as $^{18}$O, $^{13}$C) (Danezis, Tsagkaris, Camin Brusic, & Georgiou, 2016). Moreover, high chemical complexity of plant and food samples makes quality control analysis difficult. The compound-based approach in the chemical characterization of a complex sample is focused on the precise targeting of one or a small number of desired compounds (Cetó, Sánches, Serrano, Díaz-Cruz, & Núñez, 2020). Consequently, a relatively new approach, which shifted from a component-based to a pattern-based approach, should be applied (Esteki, Shahsavari, & Simal-Gandara, 2019). This approach allows simultaneous recording of a number of compounds and together with chemometrics (pattern recognition), it could provide respectable information about the biological sample.

Raman spectroscopy is a molecular vibrational spectroscopic technique based on the interaction between monochromatic light (from VIS and UV region) and a sample. The observation of the Raman scattering signal depends on the change of its polarizability during a particular mode of vibration (Larkin, 2011). During this collision, the vibrational (rotational) energy of the molecule is changed, and the scattered radiation is shifted to a different wavelength, a Raman shift (Yang & Ying, 2011). Obtained Raman spectra contain information about the chemical bonds in the sample (Yang & Ying, 2011). Raman microspectroscopy combined with multivariate chemometric analysis has already been used as a promising tool for the discrimination and authentication of different carotenoid-rich food samples. In such approaches, a number of classification methods are used, including unsupervised (Principal Component Analysis (PCA) and Hierarchical Cluster Analysis (HCA)) and supervised classification models (Partial least square discrimination analysis (PLS-DA), Linear discriminant analysis (LDA), k-nearest neighbors (KNN), and Soft independent modeling of class analogy (SIMCA) (Akpolat et al., 2020; Ivleva, Niessner, & Panne, 2005; Kolašinac, Pećinar, Danojević, Aćić, & Dajić-Stevanović, 2021b; Monago-Maraña et al., 2019). Carotenoids are especially detectable by Raman (micro)spectroscopy due to their polyene molecular structure (Schulz, Baranska, & Baranski, 2005). In addition, carotenoids are among the key nutrients in red paprika, contributing to both sensory, i.e. visual attractiveness (as they are responsible for fruit and product coloration) and nutritive quality of the fruit (Pugliese et al., 2014).

There are several reports addressing the use of Raman spectroscopy in the characterization of carotenoids in paprika, including the studies on the carotenoid content during ripening (e.g. Sharma, Sarika Bharti, Singh, & Uttam, 2019), the distribution of carotenoids along the fruit pericarp (fruit wall) (e.g. Baranski, Baranska, & Schulz, 2005) and the determination of particular carotenoid compounds in chili peppers (Sharma et al., 2019). However, to the best of our knowledge, no reports addressed the application of Raman spectroscopy as a tool in the discrimination of paprika varieties based on the differences between the components appearing within the target spectral range, mostly the total carotenoids and some individual carotenoid compounds.

The main objective of this paper is to establish the best multivariate classification model for the discrimination and therefore, the authentication of different paprika varieties at their final maturity stage, considering the needs for defining the best harvest practices, the quality assessment and the control of final food products.

## 2. Material and methods

### 2.1. Plant material

The samples of the traditional Balkan red paprika varieties (Amfora, Una, Kurtovska Kapija and Vrtka), as well as one inbreed line (derived from a cross between cultivar Amfora and Una) were used in the experiment (Fig. 1). The plants were grown under the same agro-ecological conditions, at the experimental field plots of the Institute of Field and Vegetable Crops, Novi Sad, the Republic of Serbia. The samples were harvested in the last week of September depending on the full maturity phase of each variety, i.e. as soon as the deep red color appeared across the entire fruit's surface. The fruits were previously marked to ensure the same harvest conditions for all of the samples. Five fruits were collected from each variety. After washing the fruits with deionized water, several fragments of about 2 cm$^2$ were taken from each sample for Raman spectroscopy recording. Spectra were recorded from the epidermal layer, i.e the fruits' surface was studied.

### 2.2. Raman instrumentation and spectra recording

Horiba Raman spectrometer system (Horiba Jobin Yvon, France) equipped with the Olympus BX 41 microscope was used in this study. This system possesses 532-nm and 785-nm laser sources. During the spectral recording, the 532-nm laser was focused onto the sample on the microscope stage through a 50 LWD (long working distance) objective (Olympus, Tokyo, Japan). The spectrometer is equipped with 1200 lines/mm grating. Raman scattering signals were detected by a charge-coupled-device (CCD) detector, the detection range from 900 to 1800 cm$^{-1}$ in the extended mode. The measurement was conducted with a 5s integration time, with 10 spectral accumulations, and 20–25 mW maximum output laser power. The spectral data were collected with LabSpec 6 (Horiba, France). The spectral resolution was about 3 cm$^{-1}$ and the calibration was checked by a 520.47 cm$^{-1}$ line of silicon. In order to consider a possible sample inhomogeneity, the experiments were performed in five replicates of each paprika variety. Ten spectra were recorded per each replicate making a total of 50 spectra per variety, i.e. 250 spectra in total. The assignment of the major bands was carried out using the literature data.

### 2.3. Chemometric analysis

Obtained raw spectral data were arranged in a matrix with 250 rows (objects-each variety 50 repetitions) and 334 columns (variables-wavenumbers). Raw spectra were pre-processed by applying the following procedures: 1) baseline correction, unit vector normalization; 2) baseline correction, unit vector normalization and first Savitzky-Golay derivative; 3) baseline correction, unit vector normalization and second Savitzky-Golay derivative; 4) baseline correction, unit vector normalization and third Savitzky-Golay derivative. After derivation, the spectra were smoothed with 13 smoothing points (Fig. 2).

The final stage requires the introduction of the Principal Component Analysis (PCA), prior to the application of different classification methods. PCA is required since a large number of independent variables usually cause misclassification problems (Abdul Rashid, Siti Esah Che Hussain, Razak Ahmad, & Norazami Abdullah, 2019). PCA is a method for reducing a large number of data sets of possibly correlated variables to a smaller number of uncorrelated variables (property of orthogonality) called principal components (PCs). Acquired PCs serve as the input for tested classification methods, including: PCA-LDA (Principal Component Analysis-Linear Discriminant Analysis), QDA (Quadratic Discriminant Analysis), and PLS-DA (Partial Least Square - Discriminant Analysis).

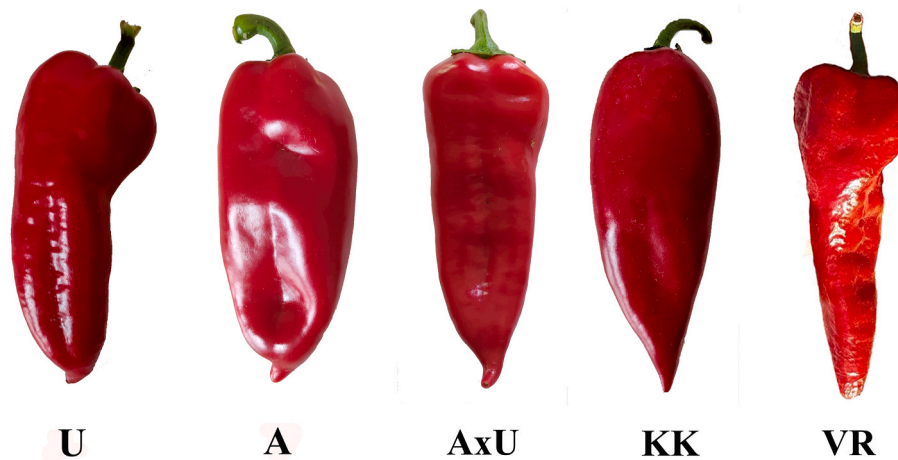LDA method serves as a model for separation into two classes

**Fig. 1.** Varieties used in the experiment: U-Una, A-Amfora, AxU-Amfora x Una, KK-Kurtovska kapija, VR-Vrtka.
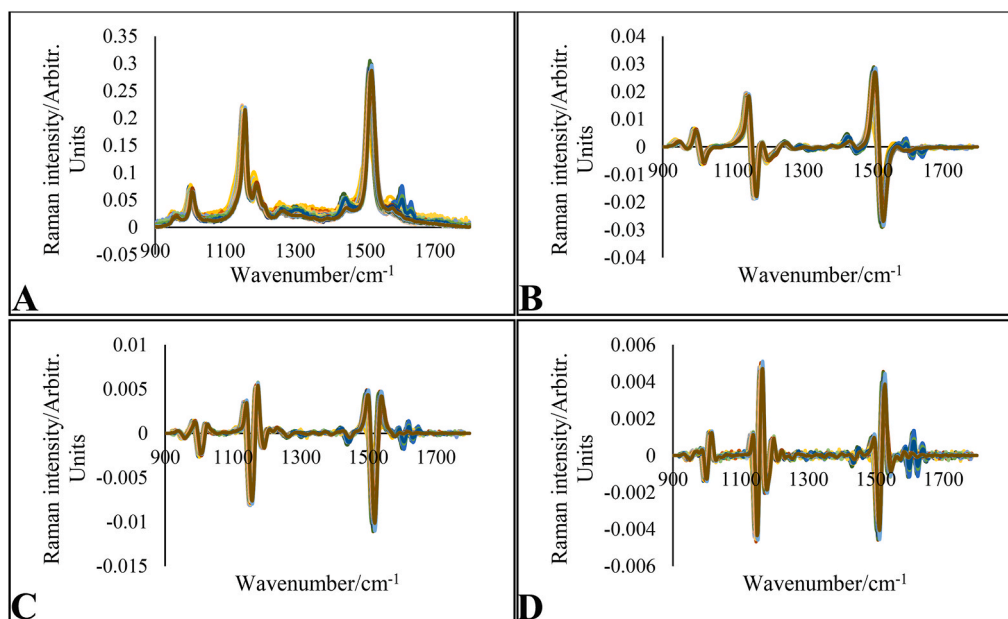


**Fig. 2.** Pre-processed spectra: A) baseline correction, unit vector normalization; B) baseline correction, unit vector normalization and the first Savitzky-Golay derivative; C) baseline correction, unit vector normalization and second Savitzky-Golay derivative; D) baseline correction, unit vector normalization and the third Savitzky-Golay derivative.

(Pomerantsev, 2014). The aim is to find the ideal projections and carry out the discrimination on the projected subspace. Distances between classes are maximized through projection while keeping minimum distance between the objects in the same class (Chen & Jiang, 2018). LDA creates linear boundaries by dividing the variable space into regions with a straight line or hyperplane (Dixon & Brereton, 2009).

QDA is a multiclass method, which can be used for the discrimination of several classes simultaneously. In contrast to LDA, QDA obtains quadratic boundaries, where a quadratic curve divides the variable space into regions (Dixon & Brereton, 2009). PLS-DA classification method is based on the search for the components or latent variables which serve to discriminate two or more different groups. Discrimination is implemented according to their maximum covariance with a target class (Uarrota et al., 2014).

To conduct multivariate classification methods, all data were divided into two sets: the training set (75% of the samples) and the test set (25% of the samples). The training set was used to establish the classification rules while the test set was used to validate them. Goodness of models was evaluated based on the highest value of sensitivity (SE) specificity (SP) and precission (P) in both training and test sample sets. SE refers to the percentage of samples of a given class that the model correctly recognises as belonging to that class:

$$SE = TP/(TP + FN)$$

SP refers to samples that do not belong to a given class and are correctly rejected by the model

$$SP = TN/(TN + FP)$$

Precision tells what fraction of predictions as a positive class were actually positive. To calculate precision, use the following formula:

$$P = TP/(TP + FP)$$

Computation of sensitivity, specificity and precision for PCA-LDA model is based on $2 \times 2$ confusion matrix (Table 1).

Unlike binary classification, PCA-QDA and PLS-DA were performed as multiclass discrimination models. Hence, metrics needed to calculate SE, SP and P for each variety were computed according more complex

**Table 1**

Confusion matrix. The table compares the prediction against the correct group assignment.

| | | Actual values | |
|---|---|---|---|
| | | positive | negative |
| Predicted values | positive | TP | FP |
| | negative | FN | TN |

*TP- true positive; TN- true negative; FP-false positive; FN-false negative.

confusion matrix (Table 2) and following equations:

## 3. Results and discussion

### 3.1. Raman spectra

Investigated paprika fruit samples were very similar in shape and color at the final maturity stage (Fig. 1). According to the obtained Raman spectra, it was clear that some subtle differences between the samples existed (Fig. 3). Fig. 2 shows the results of different pre-processing methods. Accordingly, 2B, 2C and 2D show that the changes in the Raman intensity and the bands' position (especially in the low wavenumber region) are more clearly visible using the 1st, 2nd and 3rd order derivatives. However, since there was a lot of noise coming from fluorescence, at that point, useful information was not obtained.

The literature data indicate that the red paprika varieties can slightly differ in their chemical composition at the final (deep red) physiological state, which is mainly related to some differences in the concentration of total polyphenols (Mudrić et al., 2017), flavonoids (Tundis et al., 2011) and carotenoids (Collera-Zúñiga, García Jiménez, & Meléndez Gordillo, 2005). Discrimination analysis in spectroscopy is a model of classification of observed objects (samples) according to their spectra. Generally, classification methods are divided into two groups: supervised, where the data are portioned according to their similarity into pre-defined groups; and unsupervised, where there is no prior information about the groups (Pomerantsev, 2014). So far, Raman spectroscopy coupled with chemometric classification methods has been used as a tool for discrimination of several fennel chemotypes (Gudi, Krähmer, Krüger, Hennig, & Schulz, 2014) and coffee varieties (Keidel, von Stetten, Rodrigues, Máguas, & Hildebrandt, 2010; Luna, da Silva, da Silva, Lima, & de Gois, 2019. In these studies, the Hierarchical Cluster Analysis (HCA), Soft Independent Modeling of Class Analogy (SIMCA), Linear Discriminant Analysis (LDA) and Quadratic Discriminant Analysis (QDA) were applied.

Significant bands associated with the pericarp essential constituents were observed at ~1511–1519 cm$^{-1}$ (very strong), ~1149–1151 cm$^{-1}$ (medium) and ~998–1006 cm$^{-1}$ (low intensity) (Fig. 3B), representing the characteristic bands of carotenoids. Although the carotenoids were present in the sample as minor components (in a quantity less than 1 mg kg$^{-1}$), they were very sensitive in the Raman visible region due to the signal enhancement caused by the pre-resonance effect of the analyte (Baranski et al., 2005). It is thought that the observed bands were assigned to the stretching vibration of the C=C, C−C bonds and the C−CH$_3$ in-plane group rocking vibrations, respectively (de Oliveira, Castro, Edwards, & de Oliveira, 2009; Schulz et al., 2005). These bands

**Table 2**

Confuson matrix of multiclass classification problem.

| | | Actual values | | | | |
|---|---|---|---|---|---|---|
| Predicted values | | U | A | AU | KK | VR |
| | U | a | b | c | d | e |
| | A | f | g | h | i | j |
| | AU | k | l | m | n | o |
| | KK | p | q | r | s | t |
| | VR | u | v | w | x | y |

suggest the presence of carotenoids as major secondary metabolites in a wide range of plant species including the paprika fruit (Sharma et al., 2019). According to Schulz et al. (2005), the main bands of the red paprika fruit are observed at 1517, 1158 and 1004 cm$^{-1}$ which is assigned to the capsanthin, the main carotenoid in red paprika at the physiological maturity stage. de Oliveira et al. (2009) reported that the bands at 1527, 1156 and 1006 cm$^{-1}$ can be allocated to beta carotene.

In addition to the bands corresponding to carotenoids, the weak bands were identified below 1000, at 1572–1575, 1440–1442 and ~1260 cm$^{-1}$, and in the region of 1600–1625 cm$^{-1}$ (Fig. 3B). As the parenchyma cells of paprika pericarp are very rich in various carbohydrates, the bands observed below 1000 cm$^{-1}$ probably refer to glycosidic linked stretches (Synytsya, Čopí;ková, Matějka, & Machovič, 2003). The bands observed at ~1445 cm$^{-1}$ (related to δ(CH$_2$) vibrational mode) and at ~1256 cm$^{-1}$ are most probably linked to the polygalacturonic acid (Chylińska, Szymańska-Chargot, & Zdunek, 2014) (Fig. 3B). The weak band at 1575 cm$^{-1}$ is associated with the presence of benzene ring (Trebolazabala, Maguregui, Morillas, de Diego, & Madariaga, 2017) which is the main constituent of phenolic compounds. In addition to this one, the weak band at ~1600 cm$^{-1}$ can also be attributed to phenolic compounds (Prats Mateu, Hauser, Heredia, & Gierlinger, 2016). It is known that the red pepper is rich in flavonoid glycosides and phenolic acids such as quercetin-3-glycoside (Lekala et al., 2019).

### 3.2. Multivariate classification analysis

It has already been discussed that the target compound classification analysis could sometimes be useless because of the fact that there is rarely a specific compound present only in one single type of a sample, which could consequently be used to separate (classify) samples according to such fine chemical differences. The non-targeted analysis also has some disadvantages, mainly because it requires many descriptors, some of which might be irrelevant, thus interfering with the data processing stage and possibly causing degradation of model performance. In our experiment, all classification rules are built upon the spectral region between 900 and 1800 cm$^{-1}$, which is mostly associated with the bands corresponding to the carotenoids.

Since the classification models are unpredictable due to unknown discrimination power, it is necessary to perform several models based on accepted pre-processing methods (Devos, Downey, & Duponchel, 2014). In all tested classification models, five principal components were used, explaining 99% of the entire variability.

The precision of PCA-LDA, PLS-DA and PCA-QDA in the training data was 0–100, 19.35–100 and 83.87–100%, respectively, while in the test data it was 0.00–100.00 0.00–100.00 and 89.47–100.00%, respectively. The results varied depending on the performed pre-processing method. In general, the applied pre-processing methods did not change the classification precision but PCA-LDA and PLS-DA showed uncertainty in the case where the number of corrected classified samples was zero. PCA-QDA showed the best classification power and consistency in all of the pre-processing methods used (for the graphical representation of the classification results see Supplementary material).

The main problem of the PLS-DA model was to classify the Amfora variety and Amfora x Una inbreed line. The effect could be connected with the similar genetic structure of these two samples, since they share a high number of common genes, including those possibly determining the metabolic pathways of synthesis of carotenoids, polyphenols, sugars and other metabolites. PCA-LDA showed better performance in the classification of these two samples, while PCA-QDA did not show weakness in the classification of genetically related samples. PLS-DA and PCA-LDA are based on Mahalanobis distance, but LDA presumes a single variance-covariance matrix over all classes. However, QDA presumes different variance-covariance matrices for each class separately, creating a more powerful classification rule (Dixon & Brereton, 2009) which can be the reason for the best discrimination power. In the event of incorrect classification, the varieties were randomly assigned to a
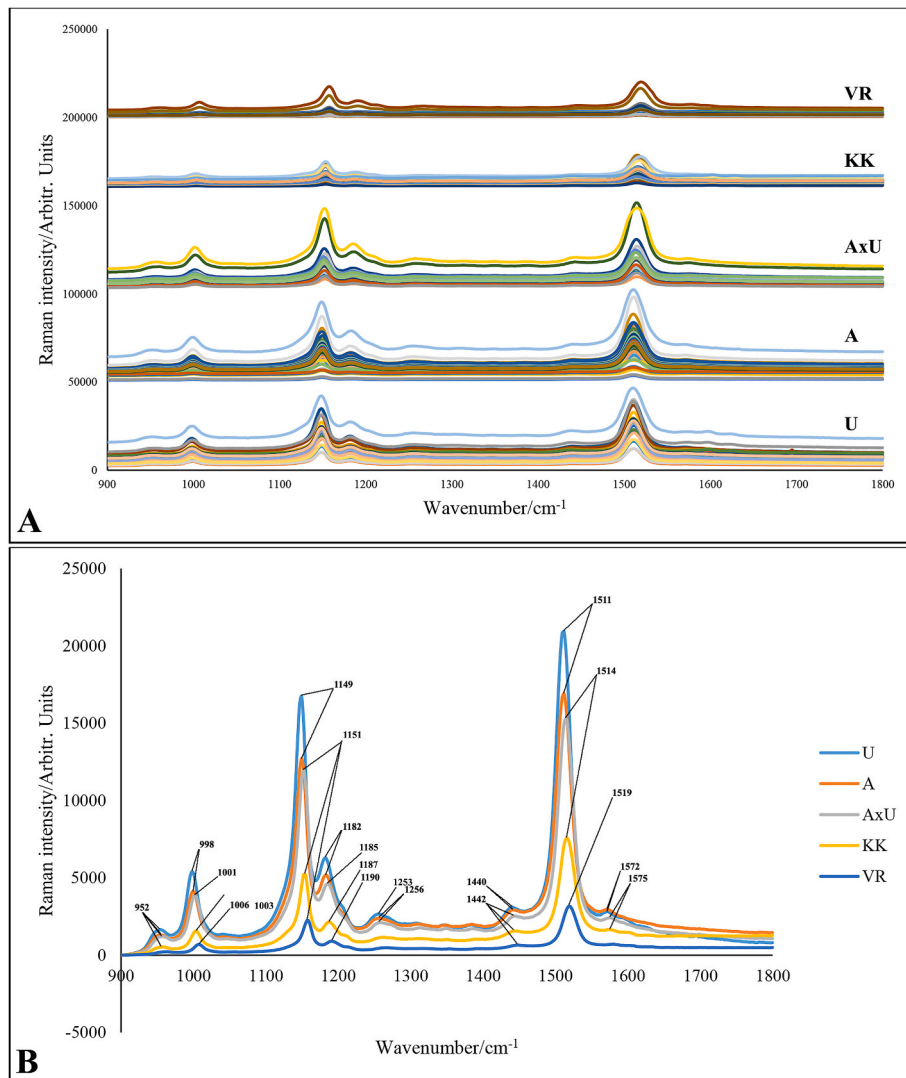
**Fig. 3.** Raw (A) and averaged (B) Raman spectra of paprika varieties: U-Una; A-Amfora; AxU-Amrofa x Una; KK-Kurtovska kapija; VR-Vrtka.

certain class.

The summary of the discrimination of each model and each pre-processing model combination is shown in Tables 3 and 4.

QDA is a very simple algorithm, and as opposed to LDA, it computes the variance structures for each class separately, creating a more powerful discrimination rule for classes with different covariance matrices, such as for biological spectra sets in which the variability within the class is the key issue.

**Table 3**
Classification results of Training sets of PCA-LDA, PLSDA and QDA models.

| Method | Variety | Training set | | | | | | | | | | | | Overall |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | BC + N | | | BC + N+ 1st D | | | BC + N+ 2nd D | | | BC + N+ 3rd D | | | |
| | | P | SE | SP | P | SE | SP | P | SE | SP | P | SE | SP | |
| | U | 100.00 | 96.77 | 100.00 | 100.00 | 87.10 | 79.03 | 100.00 | 87.10 | 97.58 | 93.10 | 87.10 | 98.39 | 93.85 |
| | A | 50.94 | 100.00 | 100.00 | 50.94 | 100.00 | 100.00 | 90.00 | 100.00 | 100.00 | 100.00 | 96.77 | 100.00 | 90.72 |
| | AU | 100.00 | 100.00 | 100.00 | 100.00 | 83.87 | 100.00 | 0.00 | 0.00 | 100.00 | 100.00 | 9.68 | 100.00 | 74.46 |
| | KK | 0.00 | 100.00 | 100.00 | 0.00 | 0.00 | 100.00 | 0.00 | 0.00 | 100.00 | 50.00 | 100.00 | 0.00 | 45.83 |
| PCA-LDA | VR | 0.00 | 100.00 | 100.00 | 0.00 | 0.00 | 100.00 | 0.00 | 0.00 | 100.00 | 0.00 | 0.00 | 100.00 | 41.67 |
| | U | 54.54 | 83.87 | 80.64 | 60.00 | 87.10 | 79.84 | 75.00 | 88.57 | 84.68 | 72.73 | 90.32 | 81.45 | 78.23 |
| | A | 52.00 | 19.35 | 95.97 | 51.92 | 19.35 | 96.77 | 58.70 | 34.29 | 96.77 | 54.90 | 25.81 | 97.58 | 58.62 |
| | AU | 96.87 | 100.00 | 99.19 | 62.22 | 90.32 | 86.29 | 100.00 | 83.87 | 100.00 | 100.00 | 83.87 | 100.00 | 91.89 |
| | KK | 100.00 | 100.00 | 100.00 | 82.35 | 51.61 | 97.58 | 86.11 | 100.00 | 95.97 | 86.11 | 100.00 | 95.97 | 91.31 |
| PLS-DA | VR | 100.00 | 100.00 | 100.00 | 100.00 | 100.00 | 100.00 | 100.00 | 100.00 | 100.00 | 100.00 | 100.00 | 100.00 | 100.00 |
| | U | 93.94 | 93.55 | 100.00 | 85.71 | 83.87 | 99.19 | 88.57 | 87.10 | 100.00 | 91.18 | 90.32 | 100.00 | 92.79 |
| | A | 100.00 | 100.00 | 98.39 | 96.30 | 96.77 | 95.97 | 100.00 | 100.00 | 96.77 | 100.00 | 100.00 | 97.58 | 98.48 |
| | AU | 100.00 | 100.00 | 100.00 | 100.00 | 100.00 | 100.00 | 100.00 | 100.00 | 100.00 | 100.00 | 100.00 | 100.00 | 100.00 |
| | KK | 100.00 | 100.00 | 100.00 | 100.00 | 100.00 | 100.00 | 100.00 | 100.00 | 100.00 | 100.00 | 100.00 | 100.00 | 100.00 |
| PCA-QDA | VR | 100.00 | 100.00 | 100.00 | 100.00 | 100.00 | 100.00 | 100.00 | 100.00 | 100.00 | 100.00 | 100.00 | 100.00 | 100.00 |

**Table 4**
Classification results of Test sets of PCA-LDA, PLSDA and QDA models.

| Method | Variety | Test set | | | | | | | | | | | | Overall |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | BC + N | | | BC + N+ 1st D | | | BC + N+ 2nd D | | | BC + N+ 3rd D | | | |
| | | P | SE | SP | P | SE | SP | P | SE | SP | P | SE | SP | |
| PCA-LDA | U | 0.00 | 0.00 | 50.00 | 100.00 | 100.00 | 86.84 | 100.00 | 100.00 | 92.10 | 100.00 | 100.00 | 94.74 | 76.97 |
| | A | 0.00 | 0.00 | 33.33 | 65.51 | 84.21 | 100.00 | 76.00 | 94.74 | 100.00 | 82.61 | 94.74 | 100.00 | 69.26 |
| | AU | 33.33 | 100.00 | 0.00 | 86.36 | 100.00 | 95.16 | 100.00 | 100.00 | 100.00 | 100.00 | 100.00 | 100.00 | 84.57 |
| | KK | 50.00 | 100.00 | 0.00 | 100.00 | 100.00 | 100.00 | 100.00 | 100.00 | 100.00 | 100.00 | 100.00 | 100.00 | 87.50 |
| | VR | 50.00 | 100.00 | 0.00 | 100.00 | 100.00 | 100.00 | 100.00 | 100.00 | 100.00 | 100.00 | 100.00 | 100.00 | 87.50 |
| PLS-DA | U | 0.00 | 0.00 | 75.00 | 61.29 | 100.00 | 84.21 | 61.29 | 100.00 | 84.21 | 59.37 | 100.00 | 82.90 | 67.36 |
| | A | 20.00 | 0.00 | 80.26 | 100.00 | 26.32 | 100.00 | 100.00 | 26.32 | 100.00 | 84.04 | 21.05 | 100.00 | 63.17 |
| | AU | 0.00 | 5.26 | 94.74 | 70.83 | 89.47 | 90.79 | 77.27 | 89.47 | 93.42 | 94.74 | 89.47 | 96.05 | 74.29 |
| | KK | 0.00 | 0.00 | 76.32 | 87.50 | 73.68 | 97.37 | 83.33 | 78.95 | 96.05 | 92.63 | 84.21 | 94.74 | 72.07 |
| | VR | 0.00 | 0.00 | 50.00 | 100.00 | 94.74 | 100.00 | 100.00 | 94.74 | 100.00 | 98.95 | 94.74 | 100.00 | 77.76 |
| PCA-QDA | U | 100.00 | 100.00 | 100.00 | 90.00 | 89.47 | 100.00 | 100.00 | 89.47 | 100.00 | 94.44 | 89.47 | 100.00 | 96.07 |
| | A | 100.00 | 94.74 | 100.00 | 100.00 | 94.74 | 97.37 | 90.47 | 100.00 | 97.37 | 90.00 | 94.74 | 97.37 | 96.40 |
| | AU | 95.00 | 100.00 | 98.68 | 95.00 | 100.00 | 98.68 | 100.00 | 100.00 | 100.00 | 100.00 | 100.00 | 100.00 | 98.95 |
| | KK | 100.00 | 100.00 | 100.00 | 100.00 | 100.00 | 100.00 | 100.00 | 88.00 | 100.00 | 89.47 | 89.47 | 100.00 | 97.25 |
| | VR | 100.00 | 100.00 | 100.00 | 100.00 | 100.00 | 100.00 | 100.00 | 100.00 | 100.00 | 100.00 | 100.00 | 100.00 | 100.00 |

BC + N - baseline correction, unit vector normalization; BC + N+1st D - baseline correction, unit vector normalization and first Savitzky-Golay derivative; BC + N+2nd D - baseline correction, unit vector normalization and second Savitzky-Golay derivative; BC + N+3rd D - baseline correction, unit vector normalization and third Savitzky-Golay derivative; P-precission; SE-sensitivity; SP-specificity.

Kolašinac, Pećinar, Danojević, Aćić, and Stevanović (2021) investigated different chemometric models in the classification of red pepper cultivars at different maturity stages. According to the training data in this research, PCA-LDA and PLS-DA indicated precision between 95% and 100%, respectively, while in the test data they showed 90–100% and 100% precision, respectively. Luna et al. (2019) used PCA-LDA, QDA and PLS-DA in combination with different pre-processing methods coupled with Raman spectroscopy for the classification of different coffee varieties. The results showed that precision was equal to 100% for all tested models except for QDA (97.8%), using MSC (multiplicative scatter correction) as a pre-processing method. On the other hand, when MC (mean centering) was used, the correct classification of the samples was 62.7% for PCA-LDA, 62.7% for QDA and 61.3% for PLS-DA.

The importance of this research is reflected in the potential use of Raman spectroscopy in the discrimination of very similar plant samples, such as some crop varieties grown at the same conditions and harvested at the same maturity phase. The paprika samples investigated in this paper are mostly used as raw material to obtain various food products that are available at the market. The most preferable and the most consumed product is the traditional Balkan dish called "ajvar", in addition to roasted and ground sweet red paprika. "Ajvar" is acknowledged Balkan food, prepared from just a few paprika varieties, known as "ajvar" or "ajvarusa" paprika types in Serbian, Bosnian and North-Macedonian languages. Such varieties as "Kurtovska kapija", "Amfora", "Crvena roga" and some others, have typical fleshy and taste-rich mesocarp (paprika "flesh"), shiny deep red color and elongated, cylindrical and somewhat flattened regular shape, suitable for efficient roasting. Best quality ajvar is made from just roasted paprika, vegetable oil and salt, where only one variety is used depending on the region, availability on the market and the price. In addition to homemade ajvar, there is significant industrial production offered by several food processing companies in the Balkans (e.g. "Bakina tajna", "Premija", "Podravka", "Dijamant", etc.), due to high demand and consumer preferences. Alongside ajvar, fleshy deep red paprika varieties are also used fresh, roasted in different products, as well as in the form of dried paprika and sweet ground paprika product. Therefore, future research should be directed towards the quality assessment of paprika products and the differences which could be accurately estimated by a non-destructive and rapid method such as Raman spectroscopy associated with proper chemometric models. Moreover, Raman imaging could be recommended for studies of alterations in nutrient content and composition in different parts of the fruit allowing a deeper understanding of carotenoid transformation during the ripening period and therefore performing the best harvest practices for different capsicum species and varieties.

## 4. Conclusion

Results obtained in this study show that the Raman spectroscopy coupled with multivariate classification analysis can be used as a tool for the discrimination of different paprika varieties. Overall, this paper aims to demonstrate the advantages derived from the use of chemometric methods as an alternative to a component-based approach. The pattern-based approach has proved to be precise, fast, low-cost and time-saving. The best classification results were acquired using QDA in all pre-treatment procedures, especially with baseline-correction and normalization. Finally, our study confirmed the possible applicability of the chemometrics-based Raman spectroscopy in food quality analyses, as well as in the authentication of food products, food adulteration and the assessment of the product's geographic origin, since the method allows determination of subtle alterations in physico-chemical traits of different products and at different processing phases.

## CRediT authorship contribution statement

**Stefan Kolašinac:** Conceptualization, Data curation, Methodology, Investigation, Writing – original draft. **Ilinka Pećinar:** Formal analysis, Investigation. **Dario Danojević:** Resources. **Zora Dajić Stevanović:** Conceptualization, Methodology, Formal analysis, Funding acquisition, Project administration, Writing – review & editing, Supervision.

## Acknowledgments

## Appendix A. Supplementary data

Supplementary data to this article can be found online at https://doi.org/10.1016/j.lwt.2022.113402.

BC + N - baseline correction, unit vector normalization; BC + N+1st D - baseline correction, unit vector normalization and first Savitzky-Golay derivative; BC + N+2nd D - baseline correction, unit vector normalization and second Savitzky-Golay derivative; BC + N+3rd D -