



Differentiating PDO Kalamata Olive Oil in Comparison with Crete Olive Oil Based on Statistical Data Analytics

Theodoros Anagnostopoulos^{a,b,}, Ioakeim Spiliopoulos^a*

^a Department of Food Science and Technology, University of Peloponnese, 24100, Kalamata, Greece

^b Department of Business Administration, University of West Attica, 12241, Athens, Greece

DOI : <https://doi.org/10.55248/gengpi.5.1024.2809>

ABSTRACT

Greece is a southeastern Mediterranean country, which cultivates olive oil since the ancient years' agriculture activities. There are several regional areas in Greece which produce high quality olive oil, such as Kalamata and Crete. Specifically, Kalamata is a city located in southeastern Greece in the Mediterranean basin and it is the capital of the Messenia regional unit. Concretely, Crete is an island located in the southern Greece in the Aegean Sea. Intuitively, both geographical areas are known for the famous olive oil quality they produce. However, Protected Designation of Origin (PDO) Kalamata olive oil, established by Council regulation (EC) No 510/2006, is considered an exceptional extra virgin olive oil variety produced in the Greece region. Subsequently, there is a need to distinguish PDO Kalamata olive oil from other olive oil varieties in Greece such as the Crete olive oil, since PDO Kalamata olive oil is the main variety exported in the global olive oil market. Differentiating PDO Kalamata olive oil compared with Crete olive oil requires experimentation with statistical data analytics models, which are able to distinguish variations of certain chemical characteristics observed in certain regions of Greece. In this paper, we use statistical data analytics models to differentiate the geographical origin of PDO Kalamata olive oil compared with Crete olive oil based on synchronous excitation-emission fluorescence spectroscopy of olive oils. Experimental evaluation as well as data visualization of the adopted statistical models are promising for differentiating the origin of PDO Kalamata olive oil with high values of prediction accuracy thus enabling companies to exploit extra virgin olive oil economic potentiality to global market place.

Keywords: PDO Kalamata olive oil, Crete olive oil, synchronous emission-excitation, fluorescence spectroscopy, statistical data analytics

1. Introduction

Cropping of plants useful for citizens is the main area of activities for smart farming (Boer & Erickson, 2019). Concretely, smart agriculture as a fundamental dimension of the smart city concept aims to define methods of efficient geographic cultivation in rural region areas (Paiho, et al., 2022). Intuitively, in the area of olive oil farmers in the Messenia region of Greece produce the protected designation of origin (PDO) extra virgin olive oil, established by Council regulation (EC) No 510/2006, with the name Kalamata olive oil in the rural areas of the Kalamata city. PDO Kalamata olive oil variety is extensively cultivated and produces the extra virgin olive oil with organoleptic properties (Trabelsi, et al., 2023), (Miao, et al., 2023). Exploiting region areas provide farmers the feasibility to gain more income since specific olive oil microclimates affect the quality of the selected variety (Issa, et al., 2023). To protect olive oil high quality and prevent its adulteration, global governmental agencies like the European Commission, International Olive Council, Codex Alimentarius, etc. have developed standards to regulate olive oil by establishing a set of physical, chemical, and organoleptic characteristics, (Aparicio, et al., 2013). The traditional chemical methods incorporated to ensure olive oil quality are focused on the identification and quantification of pre-defined compounds or classes of compounds of olive oil according to the regulations of the above-mentioned global governmental agencies. These methods are time-consuming as well as demand expensive apparatus. Subsequently, it holds the same for the detection of olive oil adulteration although these methods fail to detect the adulteration from certain adulterants.

Nowadays the non-targeted analysis has attracted much research and scientific attention. Adopted approach focuses on screening the olive oil without any prior knowledge of chemical composition. In this research approach, we used analytical techniques that produce a signal which is affected by all the compounds (i.e., metabolites) present in olive oil. Incorporated methods shorten the analysis process but a vast number of data sources are required to perform data analytics based on statistical data analytics models, (Ellis, et al., 2015). Concretely, to assess the quality of gathered olive oil there is a need to incorporate specific Internet of Things (IoT) devices, such as dedicated sensors and actuators (Rajak, et al., 2023).

Intuitively, a device that is commonly used for such a process is fluorescence spectroscopy, which is calibrated accordingly to perform differences of excitation and emission radiation to the olive oil sample (Krause, et al., 2021). Subsequently, fluorescence spectrometry has been used extensively in the past years due to its efficient precision in recognizing chemical components of olive oil samples thus exploiting its overall quality (Karoui & Blecker, 2011). Adopted technology can access input data from olive oil sample sources to measure optimally the chemical ingredients of a given olive

oil sample as well as to be able to discriminate the olive oil quality categories as well as its origin (Khani, Ghasemi, & Vanak, 2024). Specifically, fluorescence spectra technology can detect in high effectiveness adulteration of olive oil with other lower-quality oils, such as sunflower oil or soybean oil (Drakopoulou, et al., 2024). Collecting samples from different geographical origins of Greece agricultural regions enables the generation of different data sources (Schiano, et al., 2024). In addition, exploited data can be visualized and analyzed in deep detail by effective statistical data analytics models. Intuitively, the application of statistical classification models enables the classification of olive oil samples into certain categories able to differentiate the quality of each sample (Reda, et al., 2023).

In this paper, we input synchronous emission-excitation fluorescence spectra, provided by synchronous photoluminescence spectra IoT-enabled technology device, of PDO Kalamata olive oils as well as Crete olive oils. Each category of olive oil composes a unique class. Since the categories are two the classification problem is characterized as binary classification, where one class represents PDO Kalamata olive oils, while the other represents Crete olive oil. Specifically, PDO Kalamata olive oils were from the areas of: (1) Aris, (2) Thouria, (3) Verga, (4) Arfara, and (5) Meligalas. Intuitively, Crete olive oils were from the areas of: (1) Chania, (2) Rethymno, and (3) Heraklion. Subsequently, we input such data sets to certain statistical data analytics models to assess which of them has optimal results to recognize the different local cultivations. Adopted statistical data analytics models are evaluated with certain evaluation methods and metrics to observe an optimal classification of input olive oil samples. The outcome of the research effort is to be able to differentiate PDO Kalamata olive oil compared with Crete olive oil enabling companies to exploit extra virgin olive oil price potentiality to global market place.

The rest of the paper is organized as follows. In Section 2 it is presented the prior work in the research effort area. Section 3 defines the adopted statistical data analytics model. In Section 4 evaluation parameters are defined. In Section 5 experiments are performed and results are observed. Section 6 discusses the results observed in the proposed research effort, while Section 7 concludes the paper and proposes future work.

2. Prior work

Here Quality superiority of extra virgin and virgin olive oil have recently attracted consumer interest because of their quality, and its potential health benefits derived from their consumption. Concretely, high price of extra virgin olive oil and its reputation makes olive oil a target for fraudsters. Currently, significant research has been performed in the literature in the area of olive oils' analysis, classification, authentication, origin, and adulteration. Proposed spectroscopic techniques such as ultraviolet-visible (i.e., UV-Vis) absorption (Santos, et al., 2016), (Milanez, et al., 2017), fluorescence spectroscopy, (Meras, et al., 2018), mass spectrometry (Girolamo, et al., 2015), Raman spectroscopy, (Li, et al., 2018), nuclear magnetic resonance (Rotondo, Mannina, & Salvo, 2019) and FT-NIR (Mossoba, et al., 2017) have been proposed to classify and detect adulteration and origin of olive oil. Machine learning classification methods based on statistical data analytics is used to compare virgin olive oil quality in, (Zaroual, et al., 2022).

Intuitively, fluorescence spectroscopy is used along with principal component technology and factorial discriminant analysis for monitoring and classifying certain virgin olive oil varieties. Raman spectroscopy is incorporated in, (Wu, et al., 2022), to identify olive oil quality using classification techniques. Subsequently, the adopted method used a one-dimensional convolutional deep-learning neural network to observe optimal classification results. Portable Raman spectroscopy is used in, (Barros, et al., 2021), to provide quality assessment and control of several olive oil varieties. Subsequently, the proposed method adequately covers the cases of adulterated compound low-quality oils within the virgin olive oil.

Extensive machine learning classification and authentication techniques are incorporated in, (Cecchi, et al., 2020), to distinguish the origins of virgin olive oil. Concretely, it is proposed an authentication process is proposed to analyze volatile olive oil compounds and chemometrics to assess the quality of certain olive oil varieties within a local geographic area. Statistical data analytics models are incorporated in, (Stefas, et al., 2021), to classify specific olive oil varieties. Intuitively, the adopted classification method uses discrimination techniques to input machine learning models with spectroscopic data thus achieving effective prediction accuracy of olive oil behavior by exploiting fusion emission and absorption. Subsequently, fluorescence spectroscopy is incorporated in, (Rotich, et al., 2020), to classify the high quality of olive oil. Intuitively, the proposed method assesses a certain thermal oxidation technique, which exploits the potentiality of an Ultra Violet (UV) fluorescence spectroscopy system to perform specific imaging classification of extra virgin olive oil varieties.

A classification model based on time series analyses is incorporated in, (Bagnall, et al., 2012), which can distinguish several virgin olive oil varieties. Concretely, a statistical transformation of the generated input data sources is performed on each virgin olive oil variety to assess the ensemble classification schema thus observing optimal values of the prediction accuracy evaluation metric. Adulterated olive oil, in (Mata, et al., 2012), can be discriminated with the incorporation of Attenuated Total Reflection (ATR) and FTIR spectroscopy technologies. Intuitively, the proposed methods are capable of distinguishing pure samples of virgin olive oil from different oil blends by exploiting the potentiality of partial least squares discriminant analysis (PLS-DA) applied to given olive oil compounds. A multivariate classification analysis is incorporated in, (Tapp, Defernez, & Kemsley, 2003), which can distinguish extra virgin olive oils. Concretely, the adopted method is based on Fourier Transform Infrared Spectroscopy (FTIR) along with multivariate analysis to classify accurately virgin olive oils' geographic origins, which come from several producing countries.

Statistical methods and applications for distinguishing several extra virgin olive oils' local geographic origins are proposed in the literature, (Sikorska, Gorecki, & Khmelinskii, 2012). Specifically, the classification of olive oil geographic origins is based on certain chemometric data sources. Such chemometric data are generated from several olive oils compounds, which input the fluorescence spectroscopy decision-making models to achieve optimal prediction accuracy. Concretely, synchronous scanning of chemometric data sources produced by significantly detailed fluorescence spectroscopy measurements is also supported in certain research efforts, (Sikorska, Gorecki, et al., 2005). Such detailed knowledge is then exploited by specific statistical classification learning models, which can distinguish several varieties of edible extra virgin olive oils. Edible olive oils' premium

quality is assessed in the literature, (Sikorska, Swiglo, et al., 2005). Intuitively, such ability is achieved by the incorporation of synchronous fluorescence spectroscopy, which can differentiate the quantification of tocopherols from the input olive oil compounds.

Agriculture geographic origins of olive oil varieties, (Dupuy, et al., 2005), are feasible due to the incorporation of chemometric analysis. Subsequently, such advanced analytical methodology, which is applied to data sources can predict olive oil's registered designation with optimal precision taking into consideration synchronous excitation and emission of fluorescence spectra values. Rapid spectroscopic methods (Vis-NIR and FT-MIR) along with PLS analysis were applied to study thermal stress of virgin olive oils, (Maggio, et al., 2011). Specifically, due to the manipulation of generated data sources to certain statistical learning models, which can evaluate optimally spectroscopic and chemometric technologies. Pattern recognition is also incorporated in extra virgin olive oil varieties classification, (Bertran, et al., 2000).

Concretely, near-infrared spectrometry provides the technical methodology to assess the strengths of screening methods, which are then used to authenticate extra virgin olive oils from near local geographic origins. Shelf-Life olive oil varieties are monitored and then classified to certain geographic origins, (Prieto, et al., 2020). Subsequently, IoT sensors and actuators technology is exploited to enhance fluorescence spectroscopy characteristics thus being able to correctly assess the multiclass classification process, which is based on certain statistical learning models.

Currently, there are many research approaches that deal with the origins of olive oil based on statistical data analytics models. Promising research efforts incorporate generated data from several chemometric technologies. However, data manipulation requires improvement to distinguish data interconnections, which can provide efficient results. In this research paper, fluorescence spectroscopy is exploited by applying enhanced data preprocessing. Intuitively, such optimized data sources are then used by a statistical data analytics model to perform binary classification to distinguish between PDO Kalamata olive oil compared with Crete olive oil variety in cultivated within the Greece agricultural region.

3. Data model

Experimental data sources provided to perform analytics are synchronous emission-excitation fluorescence spectra. Specifically, these spectra were recorded on a Perkin Elmer LS55 spectrofluorometer using solution 1% w/v olive oil in n-hexane, where $\Delta\lambda$ (i.e., the difference between excitation and emission wavelength) was adjusted to 30 nm (Vokhminsev, et al., 2015). Intuitively, observed spectra in the current research paper were recorded at $\Delta\lambda = 30$. Concretely, the excitation and emission slit were tuned to 4 nm. The scan rate was 50nm/min. Intuitively, each olive oil sample was measured triplicate using the new freshly prepared solution. Subsequently, each measurement of an olive oil sample was statistically handled as a different sample of the same origin.

Provided data sources have a certain structure. Concretely, some of the observed data sources are collected from PDO Kalamata olive oil produced in local areas in the rural areas of the city of Kalamata in the Messenia region, while other provided data were collected from Crete olive oil produced in local areas of the Crete island. Since this is a binary classification problem there are two classes in total, namely Class 0 and Class 1. Let us define Class 0 assigned to PDO Kalamata olive oil while Class 1 assigned to Crete olive oil.

Specifically, data collected from PDO Kalamata olive oil are in total 29 olive oil samples from the local cultivation areas of: (1) Aris, (2) Thouria, (3) Verga, (4) Arfara, and (5) Meligalas. Intuitively, the distribution of collected data samples for PDO Kalamata Class 0 are as follows: (1) 2 samples from Aris, (2) 2 samples from Thouria, (3) 7 samples from Verga, (4) 15 samples from Arfara, and (5) 3 samples from Meligalas. Subsequently, data collected from Crete olive oil are in total 23 olive oil samples for the local cultivation areas of: (1) Chania, (2) Rethymno, and (3) Heraklion. Concretely, the distribution of collected data samples for Crete olive oil Class 1 are as follows: (1) 6 samples from Chania, (2) 13 samples from Rethymno, and (3) 4 samples from Heraklion. Intuitively, total number of data samples are 52, (i.e., 29 are from Class 0 and 23 are from Class 1).



Fig. 1 - Visualization of the two classes.

3.1 Data structure

Synchronous emission-excitation fluorescence spectra are composed of certain dimensional samples, where the first 5 dimensions denote the 5 predictive attributes, while the last 1 dimension denote the class attribute, such as (p_i, c_j) . Specifically, let us define the 5 predictive attributes as p_i , where $i \in [1, 5]$ is the identifier of each predictive attribute, where: $i = 1$ refers to tocopherols, $i = 2$ refers to phenolic compounds, $i = 3$ refers to oxidation products of triglycerides, $i = 4$ refers to oxidation products of tocopherols, and $i = 5$ refers to chlorophyll's predictive attributes. Intuitively, let us define the class attribute as c_j , where $j \in [0,1]$ is the identifier of the class attribute value, where: $j = 0$ refers to Class 0 value denoting PDO Kalamata olive oil, while $j = 1$ refers to Class 1 value denoting Crete olive oil.

3.2 Data visualization

There is a need to visualize provided data classes to observe in detail their structure. The distribution of the adopted classes (i.e., Class 0 for PDO Kalamata olive oil and Class 1 for Crete olive oil) should be plotted to be able to distinguish the different nature of the examined class attributes. Fig. 1 visualizes the two classes, where it can be observed that in the lower left side of Fig. 1 there are presented the 29 data samples assigned to Class 0. Specifically, PDO Kalamata olive oil sample instances could be distinguished since they are denoted with small blue 'x' marks. Concretely, it can be also observed that in the upper right side of Fig. 1 there are presented the 23 data samples assigned to Class 1. Intuitively, Crete olive oil sample instances could be differentiated since they are plotted as small red 'x' marks.

4. Evaluation parameters

Deep understanding of the results requires to incorporate certain parameters able to assess the performance of the adopted statistical data analytics models. Specifically, evaluation methods and evaluation metrics should be incorporated to perform specific experiments and observe derived results.

4.1 Evaluation method

Being able to evaluate a statistical data analytics model there are used certain evaluation methods. In the current research paper authors adopt one of the widely used evaluation methods, due to its simplicity and optimum results, which is 10-fold cross-validation, (Trevor, Tibshirani, & Friedman, 2009). Specifically, such an evaluation method divides the input dataset into 10 equal sized parts and then in a certain loop incorporates the first 9 parts to train the statistical learning classification algorithm and the remaining 1 to test the classifier. This process is repeated until all the parts are used for training and testing. The proposed evaluation method is adopted in the data analytics methodology since it provides effective results based on certain input data able to explain the observed data source's predictive analytics behavior.

4.2 Evaluation metrics

Provided the evaluation method, which is proposed to support the experimental setup there is a need to adopt specific evaluation metrics. There are incorporated several metrics, which are: (1) prediction accuracy, (2) correctly classified instances, and (3) confusion matrix that can assess the efficiency of a statistical data analytics model

4.2.1 Prediction accuracy

Synchronous eEfficiency of the adopted statistical data analytics model is assessed by incorporating prediction accuracy evaluation metric, $a \in [0, 1]$, which is defined in the following mathematical equation, (1):

$$a = \frac{tr_{pos} + tr_{neg}}{tr_{pos} + fl_{pos} + tr_{neg} + fl_{neg}} \quad (1)$$

Where, tr_{pos} , are the instances, which are classified correct as positives, and, tr_{neg} , are the instances, which are classified correct as negatives. In addition, fl_{pos} , are the instances, which are classified false as positives, and, fl_{neg} , are the instances, that are classified false as negatives. A low value of a means a weak classifier while a high value of a indicates an effective statistical pervasive data analytics classifier. Concretely, experimental assessment based on the defined statistical quantities of: (1) tr_{pos} , (2) tr_{neg} , (3) fl_{pos} , and (4) fl_{neg} , which compose the prediction accuracy evaluation metric's experimental value, achieve to express the data sources' dynamics and explain the observed optimal results.

Table 1 - Confusion matrix for binary classification.

	Class 0	Class 1	← Classified as
A		B	Class 0
C		D	Class 1

4.2.2 Correctly classified instances

In statistical data analytics, it is common to express prediction accuracy as a percentage thus observed results being more easily interpreted and presented. Concretely, it is used the term correctly classified instances, $c \in [0\%, 100\%]$, which is defined according to the following mathematical equation, (2):

$$c = a\% \quad (2)$$

Where, a value close to 0% means that the classification model is not efficient, while a value close to 100% indicates that the statistical data analytics model is able to classify instances optimally.

4.2.3 Confusion matrix

Adopted statistical classification algorithm is also evaluated with the confusion matrix evaluation metric. Confusion matrix is a special form of matrix, which in the case of a binary classification of two classes, (i.e., Class 0: PDO Kalamata olive oil, and Class 1: Crete olive oil) has the following encoded form, as described in Table 1.

Where, "A" quantity depicts the number of Class 0 instances, which are classified correctly as instances of Class 0. "B" quantity depicts the number of Class 0 instances, which are falsely classified as instances of Class 1. "C" quantity depicts the number of Class 1 instances, which are falsely classified as instances of Class 0, while "D" quantity depicts the number of Class 1 instances, which are correctly classified as instances of Class 1. A given classification model is considered efficient if it maximizes the elements of the main diagonal of the confusion matrix (i.e., "A" and "D") and minimizes the other elements. A confusion matrix is incorporated in data analytics evaluation methodology to support effectiveness and explain in deep detail the statistical nature of output experimental results observed by the prediction accuracy evaluation metric.

3. Experiments and results

Adopted data model, which is pre-processed according to two class values (i.e., Class 0: PDO Kalamata olive oil and Class 1: Crete olive oil) is incorporated to perform specific experiments and observe derived results. It holds that an experimental setup is necessary to formulate the experimental phase with certain evaluation methods and metrics and observe the results of the current research effort.

Table 2 - Numerical spectra values ranges of Class 0.

Predictive attributes	Min	Max
Tocopherols	232.221	237.223
Phenolic compounds	3.606	8.607
Oxidation products of triglycerides	1.896	5.149
Oxidation products of tocopherols	1.237	3.262
Chlorophylls	45.243	50.246

3.1 Experimental setup

Specific parameters are incorporated to set up the experimental process. Concretely, it is defined as the number of classes (i.e., Class 0 and Class 1), which is assigned to each data sample instance. Intuitively, predictive attributes used to describe a certain class attribute are defined accordingly. Subsequently, a certain statistical data analytics model should be adopted to perform the experiments and observe the results.

3.1.1 Binary classification

Being able to evaluate Since the number of class values is two this classification process is characterized as a binary classification problem. Specifically, two classes are defined as follows: (1) Class 0: PDO Kalamata olive oil, and (2) Class 1: Crete olive oil. Intuitively, the number of predictive attributes is five, which are characterized as follows: (1) 1st predictive attribute: 'tocopherols', (2) 2nd predictive attribute: 'phenolic compounds', (3) 3rd predictive attribute: 'oxidation products of triglycerides', (4) 4th predictive attribute: 'oxidation products of tocopherols', and (5) 5th predictive attribute: 'chlorophylls'. The number of data sample instances is 52, where 29 samples are assigned to Class 0 while the remaining 23 are assigned to Class 1. Specifically, samples distribution for Class 0 is the following: (1) 2 samples from Aris, (2) 2 samples from Thouria, (3) 7 samples from Verga, (4) 15 samples from Arfara, and (5) 3 samples from Meligalas. Subsequently, samples distribution for Class 1 is as follows: (1) 6 samples from Chania, (2) 13 samples from Rethymno, and (3) 4 samples from Heraklion.

Table 3 - Numerical spectra values ranges of Class 1.

Predictive attributes	Min	Max
-----------------------	-----	-----

Tocopherols	131.868	136.866
Phenolic compounds	14.249	18.997
Oxidation products of triglycerides	4.109	8.859
Oxidation products of tocopherols	1.254	3.245
Chlorophylls	52.482	57.481

Numerical spectra values as recorded by a Perkin Elmer LS55 spectrofluorometer using solution 1% w/v olive oil in n-hexane, where $\Delta\lambda$ (i.e., the difference between excitation and emission wavelength) was adjusted to 30 nm for the predictive attributes of Class 0 (i.e., PDO Kalamata olive oil) are observed in the following ranges: (1) tocopherols spectra values are within interval [232.221, 237.223], (2) phenolic compounds values are within range [3.606, 8.607], (3) oxidation products of triglycerides spectra values are within interval [1.896, 5.149], (4) oxidation products of tocopherols values are within range [1.237, 3.262], and (5) chlorophylls spectra values are within interval [45.243, 50.246]. Numerical spectra values of Class 0 (i.e., PDO Kalamata olive oil) are presented in Table 2.

Subsequently, the numerical values of the predictive attributes of Class 1 (i.e., Crete olive oil) take values as described in following ranges: (1) tocopherols spectra values are within interval [131.868, 136.866], (2) phenolic compounds values are within range [14.249, 18,997], (3) oxidation products of triglycerides spectra values are within interval [4.109, 8.859], (4) oxidation products of tocopherols values are within range [1.254, 3.245], and (5) chlorophylls spectra values are within interval [52.482, 57.481]. Numerical values of Class 1 (i.e., Crete olive oil) are presented in Table 3.

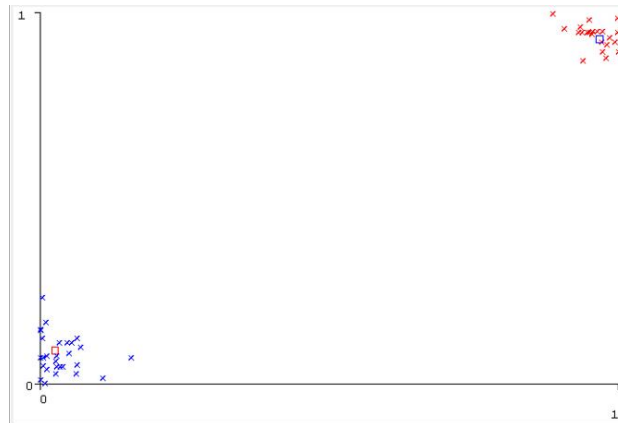


Fig. 2 – Classification visualization of the two classes.

3.1.2 C4.5 statistical data analytics model

There were certain experiments to be able to select the optimum statistical pervasive data analytics model, which is efficient for the examined binary classification problem. Specifically, we performed experiments with several statistical data analytics classifiers available in the Weka machine learning software, (Hall, et al., 2009), (Frank, Hall, & Witten, 2016). Intuitively, the pervasive data analytics model, which has optimal predictive behavior emerged to be the C4.5 tree based statistical data analytics model (i.e., J48 tree-based classification model as implemented in Weka machine learning software) thus it is adopted for further experimentation to observe the derived results of the current research paper.

3.2 Derived results

Evaluation of the experimental requires to define a specific evaluation method (i.e., 10-fold cross-validation) and metrics used to assess the efficiency of the adopted statistical data analytics model, which in this case is the C4.5 tree-based data analytics model. Intuitively, based on certain evaluation parameters specific derived results are observed, which define the effectiveness of the incorporated experimental setup adopted in the current research effort. Concretely, to understand the observed results and be able to explain the research effort's findings it is significant to use the incorporated evaluation method and evaluation metrics. Such knowledge would reveal the inherent complexity that exists in the provided data sources aiming to observe optimal results for the adopted statistical data analytics model.

3.2.1 Observed classification visualization

Applying C4.5 tree-based data analytics model to the supplied data sources for both classes (i.e., Class 0 and Class 1) there are observed promising results denoting a data analytics model with impact during the prediction process. Specifically, in Fig. 2 it can be observed that only a single sample of

Class 0 is falsely classified as sample of Class 1 annotated with a small blue square at the upper right side of the graph. Subsequently, it can also be observed in Fig. 2 that only a single sample of Class 1 is falsely classified as a sample of Class 0 annotated with a small red square at the lower left side of the graph. Intuitively, all the other data samples are correctly classified by the adopted statistical data analytics model

3.2.2 Observed prediction accuracy

Evaluation method, which is incorporated to evaluate the adopted binary classification statistical data analytics model is 10-fold cross-validation. According to the adopted evaluation method observed prediction accuracy is: $a = 0.9615$, which is a high value for prediction accuracy thus proving that the adopted statistical data analytics model is suitable for the examined binary classification problem. Concretely, the high value observed for the prediction accuracy enables the adopted tree-based data analytics model to be incorporated for similar use in new unseen olive oil instances in a further future research that might extend the potentiality of the current research effort to both examined geographical regions of interest (i.e., Kalamata city and Crete island).

3.2.3 Observed correctly classified instances

Considering the adopted evaluation method of 10-fold cross-validation correctly classified instances it occurred to be: $c = 96.15\%$, which indicated that the selected statistical data analytics model is an optimal choice for the examined classification problem.

Table 4 – Binary confusion matrix observed results.

Class 0	Class 1	← Classified as
28	1	Class 0
1	22	Class 1

3.2.4 Observed confusion matrix

Binary confusion matrix results are derived based on a 10-fold validation evaluation method for the examined binary classification problem. Derived results are presented in Table 4. Concretely, it can be observed that most of the classified instances are located in the main diagonal of Table 4. Specifically, the quantity of elements in the main diagonal depicts the significant number of certain instances, which are correctly classified. Concretely, such an optimal prediction behavior indicates a classification model with prediction impact for the examined binary classification problem. Intuitively, such a detailed confusion matrix enables the observation of experimental results in detail thus being able to assess the effectiveness of the adopted data analytics model for differentiating PDO Kalamata olive oil in comparison with Crete olive oil based on provided experimental instances.

4. Discussion

Definition of the problem indicates a binary classification problem of two discrete classes (i.e., Class 0: PDO Kalamata olive oil and Class 1: Crete olive oil), with five separate predictive attributes and a total of 52 data instances. Specifically, 29 sample instances are assigned to Class 0 while the remaining 23 are assigned to Class 1. Subsequently, current research paper has achieved significant values of the observed results based on certain evaluation metrics, which indicate the stability of the examined evaluation parameters. Intuitively, observed prediction accuracy and correctly classified instances are relatively high. Concretely, binary confusion matrix has high values in the main diagonal denoting an accurate predictive model. Subsequently, adoption of C4.5 data analytics model for binary classification required in the current research paper result in robust results able to differentiate PDO Kalamata olive oil in comparison with Crete olive oil based on applied methodological framework.

5. Conclusion and future work

Identifying differences between two regions that produce extra virgin olive oil, such PDO Kalamata olive oil and Crete olive oil is of highly concern in the current research paper. Adopted synchronous photoluminescence spectra of olive oils IoT device can provide initial data sources to compare the origins of PDO Kalamata olive oil and Crete olive oil. In this research paper, we use statistical data analytics models to perform binary classification between Class0 (i.e., PDO Kalamata olive oil) and Class 1 (i.e., Crete olive oil) collected from several geographic origins. Evaluation of the statistical models are based on certain methods and metrics, which have proved to be promising for differentiating PDO Kalamata olive oil compared with Crete olive oil. Note that according to the research outcomes, future work should mainly focus on the incorporation of more detailed input measurement data sources based on improvements in synchronous photoluminescence spectra IoT-enabled technology, thus providing a more stable input to the selected statistical classification model. Intuitively, current research effort could be further used in deep detail to verify authentication and to detect adulteration of PDO Kalamata olive oil thus facing the fraud problem occurring in the olive oil global market.

References

Aparicio, R., Morales, M. T., Ruiz, R. A., Tena, N., & Gonzalez, D. L. G. (2013). Authenticity of olive oil: Mapping and comparing official methods and promising alternatives. *Food Research International*, 54(2), 2025–2038.

- Bagnall, A., Davis, L., Hills, J., & Lines, J. (2012). Transformation Based Ensembles for Time Series Classification. *Proceedings of the 2012 SIAM International Conference on Data Mining (SDM)*, Anaheim, California, USA, April 26 – 28, 307–318.
- Barros, H. A. S., Paixao, L. S., Nascimento, M. H. C., Lacerda, V. J., Figueiras, P. R., & Romao, W. (2021). Use of portable Raman spectroscopy in the quality control of extra virgin olive oil and adulterated compound oils. *Vibrational Spectroscopy*, 116, 1–10.
- Bertran, E., Blance, M., Coello, J., Iturriaga, H., Maspocho, S., & Montoliu, I. (2000). Near infrared spectrometry and pattern recognition as screening methods for the authentication of virgin olive oils of very close geographical origins. *Journal of Near Infrared Spectroscopy*, 8, 2000, 45–52.
- Boer, J. L. D., & Erickson, B. (2019). Setting the Record Straight on Precision, Agriculture Adoption. *Agronomy Journal*, 111(4), 1552–1569.
- Cecchi, L., Migliorini, M., Giambanelli, E., Rossetti, A. Cane, A., Mulinacci, N., & Melani, F. (2020). Authentication of the geographical origin of virgin olive oils from the main worldwide producing countries: A new combination of HS-SPME-GC-MS analysis of volatile compounds and chemometrics applied to 1217 samples. *Food Control*, 112, 1–10.
- Drakopoulou, S. K., Kritikou, A. S., Baessmann, C., & Thomaidis, N. (2024). Untargeted 4D-metabolomics using Trapped Ion Mobility combined with LC-HRMS in extra virgin olive oil adulteration study with lower-quality olive oils. *Food Chemistry*, 434, 1–9.
- Dupuy, N., Dreau, Y. L., Ollivier, D., Artaud, J., Pinatel, C., & Kister, J. (2005). Origin of French Virgin Olive Oil Registered Designation of Origins Predicted by Chemometric Analysis of Synchronous Excitation-Emission Fluorescence Spectra. *Journal of Agricultural and Food Chemistry*, 53(24), 9361–9368.
- Ellis, D. I., Muhamadali, H., Haughey, S. A., Elliott, C. T., & Goodacre, R. (2015). Point-and-shoot: Rapid quantitative detection methods for on-site food fraud analysis—moving out of the laboratory and into the food supply chain. *Analytical Methods*, 7(22), 9375–9716.
- Frank, E., Hall, M. A., & Witten, I. H. (2016). *The WEKA Workbench. Online Appendix for Data Mining: Practical Machine Learning Tools and Techniques* (4th Edition). Morgan Kaufmann.
- Girolamo, F. D., Masotti, A., Lante, I., Scapaticci, M., Calvano, C. D., Zambonin, C., Muraca, M., & Putignani, L. A. (2015). Simple and effective mass spectrometric approach to identify the adulteration of the mediterranean diet component extra-virgin olive oil with corn oil. *International Journal of Molecular Sciences*, 16, 20896–20912.
- Hall, M., Frank, E., Holmes, G., Pfahringer, B., Reutemann, P., & Witten, I. H. (2009). *ACM SIGKDD Explorations Newsletter*, 11(1), 10–18.
- Issa, A., Riachy, M. E., Mitri, C. B., Doumit, J., Skaff, W., & Karam, L. (2023). Influence of geographical origin, harvesting time and processing system on the characteristics of olive-mill wastewater: A step toward reducing the environmental impact of the olive oil sector. *Environmental Technology & Innovation*, 32, 1–12.
- Karoui, R., & Blecker, C. (2011). Fluorescence Spectroscopy Measurement for Quality Assessment of Food Systems—a Review. *Food Bioprocess Technology*, 4, 364–386.
- Khani, S., Ghasemi, & Vanak, Z. P. (2024). Development of computer vision system for classification of olive oil samples with different harvesting years and estimation of chlorophyll and carotenoid contents: A comparison of the proposed method's efficiency with UV-Vis spectroscopy. *Journal of Food Composition and Analysis*, 106078, 1–42.
- Krause, J., Gruger, H., Gebauer, L., Zheng, X., Knobbe, J., Pgnier, T., Kicherer, A., Gruna, R., Langle, T., & Beyerer, J. (2021). Smart Spectrometer-Embedded Optical Spectroscopy for Applications in Agriculture and Industry. *Sensors*, 21, 1–18.
- Li Y., Fang, T. Zhu, S., Huang, F., Chen, Z., & Wang, Y. (2018). Detection of olive oil adulteration with waste cooking oil via Raman spectroscopy combined with iPLS and SiPLS. *Spectrochimica Acta Part A: Molecular Biomolecular Spectroscopy*, 189, 37–43.
- Maggio, R. M., Valli, E., Bendini, A., Caravaca, A. M. G., Toschi, T. G., & Cerretani, L. (2011). A spectroscopic and chemometric study of virgin olive oils subjected to thermal stress. *Food Chemistry*, 127, 216–221.
- Mata, P. D. L., Vidal, A. D., Sendra, J. M. B., Medina, A. R., Rodriguez, L. C., & Canada, M. J. A. (2012). Olive oil assessment in edible oil blends by means of ATR-FTIR and chemometrics. *Food Control*, 23, 449–455.
- Meras, I. D., Manzano, J. D., Rodriguez, D. A., & Pena, A. M. (2018). Detection and quantification of extra virgin olive oil adulteration by means of autofluorescence excitation-emission profiles combined with multi-way classification. *Talanta*, 178, 751–762.
- Miao, X., Ma, J., Miu, X., Zhang, H., Geng, Y., Hu, W., Deng, Y., & Li, N. (2023). Integrated transcriptome and proteome analysis the molecular mechanisms of nutritional quality in 'Chenggu-32' and 'Koroneiki' olives fruits (*Olea europaea* L.). *Journal of Plant Physiology*, 288, 1–12.
- Milanez K. D. T. M., Nobrega, T. C. A., Nascimento, D. S., Insausti, M., Band, B. S. F., & Pontes, M. J. C. (2017). Multivariate modeling for detecting adulteration of extra virgin olive oil with soybean oil using fluorescence and UV-Vis spectroscopies: A preliminary approach. *LWT – Food Science and Technology*, 85, 2017, 9–15.

- Mossoba, M. M., Azizian, H., Kia, A. R. F., Karunathilaka, S. R., & Kramer, J. K. G. (2017). First Application of Newly Developed FT–NIR Spectroscopic Methodology to Predict Authenticity of Extra Virgin Olive Oil Retail Products in the USA. *Lipids*, 52, 443–455.
- Paiho, S., Tuominen, P., Rockman, J., Ylikerala, M., Pajula, J., & Siikavirta, H. (2022). Opportunities of collected city data for smart cities. *IET Smart Cities*, 4(4), 275–291.
- Prieto, A. L., Tena, N., Ruiz, R. A., Gonzalez, D. L. G., & Sirkorska, E. (2020). Monitoring Virgin Olive Oil Shelf-Life by Fluorescence Spectroscopy and Sensory Characteristics: A Multidimensional Study Carried Out under Simulated Market Conditions. *Foods*, 9, 1–20.
- Rajak, P., Ganguly, A., Adhikary, S., & Bhattacharya, S. (2023). Internet of Things and smart sensors in agriculture: Scopes and challenges. *Journal of Agriculture and Food Research*, 14, 1–13.
- Reda, R., Saffaj, T., Bouzida, I., Saidi, O., Belgirir, M., Lakssir, B., & Hadrami, E. M. E. (2023). Optimized variable selection and machine learning models for olive oil quality assessment using portable near infrared spectroscopy. *Spectrochimica Acta Part A: Molecular and Biomolecular Spectroscopy*, 303, 1–11.
- Rotich, V., Riza, D. F. A., Giametta, F., Susuki, T., Ogawa, Y., & Kondo, N. (2020). Thermal oxidation assessment of Italian extra virgin olive oil using an UltraViolet (UV) induced fluorescence imaging system. *Spectrochimica Acta Part A: Molecular and Biomolecular Spectroscopy*, 237, 2020, 1–8.
- Rotondo, A., Mannina, L., & Salvo, A. (2019). Multiple Assignment Recovered Analysis (MARA) NMR for a Direct Food Labeling: The Case Study of Olive Oils. *Food Analytical Methods*, 12, 1238–1245.
- Santos, R. A., Cancilla, J. C., Perez, A. P., Moral, A., & Torrecilla, J. S. (2016). Quantifying binary and ternary mixtures of monovarietal extra virgin olive oils with UV–vis absorption and chemometrics. *Sensors and Actuators B: Chemical*, 234, 115–121.
- Schiano, M. E., Sodano, F., Cassiano, C., Magli, E., Seccia, S., Rimoli, M. G., & Albrizio, S. (2024). Monitoring of seven pesticide residues by LC-MS/MS in extra virgin olive oil samples and risk assessment for consumers. *Food Chemistry*, 442, 1–9.
- Sikorska, E., Gorecki, T., Khmelinskii, I. V., Sikorski, M., & Koziol, J. (2005). Classification of edible oils using synchronous scanning fluorescence spectroscopy. *Food Chemistry*, 89, 217–225.
- Sikorska, E., Swiglo, A. G., Khmelinskii, I., & Sirorski, M. (2005). Synchronous Fluorescence of Edible Vegetable Oils. Quantification of Tocopherols. *Journal of Agriculture and Food Chemistry*, 53(18), 6988–6994.
- Sikorska E., Khmelinskii I., Sikorski M. (2012). Analysis of Olive Oils by Fluorescence Spectroscopy: Methods and Applications. Olive Oil – Constituents, Quality, Health Properties and Bioconversions. *InTech*, 4, 63–88.
- Stefas, D., Gyftokostas, N., Kourelis, P., Nanou, E., Kokkinos, V., Bouras, C., & Couris, S. (2021). Discrimination of olive oils based on the olive cultivar origin by machine learning employing the fusion of emission and absorption spectroscopic data. *Food Control*, 130, 1–8.
- Tapp, H. S., Defernez, M., & Kemsley, E. K. (2003). FTIR Spectroscopy and Multivariate Analysis Can Distinguish the Geographic Origins of Extra Virgin Olive Oils. *Journal of Agricultural and Food Chemistry*, 51, 6110–6115.
- Trabelsi, L., Ncube, B., Hassena, A. B., Zouairi, M., Amar, F. B., & Gargouri, K. (2023). Comparative study of productive performance of two olive oil cultivars Chemlali Sfax and Koroneiki under arid conditions. *South African Journal of Botany*, 154, 356–364.
- Trevor, H., Tibshirani, R., & Friedman, J. H. (2009). *The Elements of Statistical Learning: Data Mining, Inference, and Prediction* (2nd Edition). New York, Springer.
- Vokhmintsev, A. S., Minin, M. G., Henaish, A. M. A., & Weinstein, I. A. (2015). Spectrally resolved thermoluminescence measurements in fluorescence spectrometer. *Measurement*, 66, 90–94.
- Wu, X., Gao, S., Niu, Y., Zhao, Z., Xu, B., Ma, R., Liu, H., & Zhang, Y. (2022). Identification of olive oil in vegetable blend oil by one-dimensional convolutional neural network combined with Raman spectroscopy. *Journal of Food Composition and Analysis*, 108, 1–7.
- Zaroual, H., Chene, C., Hadrami, E. M. E., & Karoui, R. (2022). Comparison of four classification statistical methods for characterizing virgin olive oil quality storage up to 18 months. *Food Chemistry*, 370, 1–16.