



UNIVERSITÀ DI PARMA

ARCHIVIO DELLA RICERCA

University of Parma Research Repository

Detection of soft-refined oils in extra virgin olive oil using data fusion approaches for LC-MS, GC-IMS and FGC-Enose techniques: The winning synergy of GC-IMS and FGC-Enose

This is the peer reviewed version of the following article:

Original

Detection of soft-refined oils in extra virgin olive oil using data fusion approaches for LC-MS, GC-IMS and FGC-Enose techniques: The winning synergy of GC-IMS and FGC-Enose / Tata, Alessandra; Massaro, Andrea; Damiani, Tito; Piro, Roberto; Dall'Asta, Chiara; Suman, Michele. - In: FOOD CONTROL. - ISSN 0956-7135. - 133:(2022). [10.1016/j.foodcont.2021.108645]

Availability:

This version is available at: 11381/2922208 since: 2022-04-30T09:17:08Z

Publisher:

ELSEVIER SCI LTD

Published

DOI:10.1016/j.foodcont.2021.108645

Terms of use:

Anyone can freely access the full text of works made available as "Open Access". Works made available

Publisher copyright

note finali coverpage

(Article begins on next page)

Food Control

Detection of soft-refined oils in extra virgin olive oil using data fusion approaches for LC-MS, GC-IMS and FGC-Enose techniques: the winning synergy of GC-IMS and FGC-Enose --Manuscript Draft--

Manuscript Number:	FOODCONT-D-21-00923R2
Article Type:	Research Paper
Keywords:	Soft-deodorization; deacidification; FGC-Enose; adulteration; GC-IMS; LC-MS; soft-refined olive oil (SROO); data fusion
Corresponding Author:	Michele Suman Parma, ITALY
First Author:	Michele Suman
Order of Authors:	Michele Suman Alessandra Tata Andrea Massaro Tito Damiani Roberto Piro Chiara Dall'Asta
Abstract:	<p>Extra virgin olive oil (EVOO) is frequently adulterated by mixing it with soft refined oils (SROO). The differentiation of EVOO from its blends with SROO is not possible with the most common approaches, and, for this reason, the discriminating power of liquid chromatography-high resolution mass spectrometry (LC-MS), gas-chromatography ion mobility spectrometry (GC-IMS) and flash gas-chromatography electronic nose (FGC-Enose) was examined previously. Here, the combination of the above-mentioned techniques for an improvement in classification power of the methods is explored. A total of 43 commercial EVOOs and 18 illegal mixtures of SROO with EVOO were previously analysed by LC-(+/-)MS, GC-IMS and FGC-Enose. Low-level and mid-level data fusion of the four datasets were performed. The merged unique fingerprints were submitted to partial least squared discriminant analysis (PLS-DA), and the extrapolated most informative variables were used to build support vector machine (SVM) classifiers. Statistical indicators were calculated and compared to find out the best classifier. The results of PLS-DA-SVM strategies on the combination of datasets demonstrated that, after low-level data fusion, the discriminatory capability of the two merged GC-based techniques was remarkably improved as compared to the individual techniques. This indicates that merging the datasets before PLS-DA better retrieves the most informative variables and, thus, enhances group separation and classification of unknowns. The combination of LC(+/-)MS datasets, both by mid- and low-level data fusion, did not show significant enhancement in terms of discrimination of EVOO from SROO as compared to the individual LC(+/-)MS matrix. The low-level combination of the four datasets (LC(+/-)MS, GC-IMS, FGC-Enose) was successful, although this laborious option is not a viable path in industry quality assurance. This study primarily provides new paths for the authentication of EVOO, taking advantage of merging multimodal LC-(+/-)MS, GC-IMS and FGC-Enose data, with consequent improvement in the performances of the classification models. The most promising results were achieved by the low-level data fusion of GC-IMS and FGC-Enose data.</p>
Suggested Reviewers:	Philipp Weller p.weller@hs-mannheim.de Expertise in Analytical Data Fusion Approaches Marta Ferreira marta.ferreiro@uca.es Expertise in GC-IMS Analysis

	<p>Tullia Gallina Toschi Tultullia.gallinatoschi@unibo.it Expertise in Olive Oil topics and fraud issues</p>
--	--



To the kind attention of:
Editor

Dr. Q. Rao, PhD

(Food Science; Food Chemistry; Food Quality; Food Safety)
Florida State University College of Human Sciences Nutrition,
Food & Exercise Sciences, Tallahassee, Florida, United States of America
Food Control - Elsevier

Authors:

*Alessandra Tata¹, Andrea Massaro¹, Tito Damiani², Roberto Piro¹, Chiara Dall'Asta², Michele Suman^{3,4} **

¹ Istituto Zooprofilattico Sperimentale delle Venezie, Laboratorio di Chimica Sperimentale, Vicenza, Italy

² Department of Food and Drug, University of Parma, Parma, Italy

³ Department of Analytical Food Science, Barilla G. e R. Fratelli S.p.A., Parma, Italy

⁴ Department for Sustainable Food Process, Catholic University Sacred Heart, Piacenza, Italy

Title:

“Advantages and disadvantages of data fusion of LC-MS, GC-IMS and FGC-Enose techniques in the authentication of extra virgin olive oil”

Dear Editor,

the present paper describes an original data-fusion exercise devoted to face recent fraud issues within Extra Virgin Olive Oil food chain. In particular taking into account of LC-(+/-)MS, GC-IMS and FGC-Enose analytical data, low-level data fusion of GC-IMS and FGC-ENose datasets demonstrated to be effective in order to generate an optimal model within a new framework for the authentication of EVOO.

It was a positive synergic effort among a control authority (Istituto Zooprofilattico), an academic (University of Parma) and an industrial (Barilla Advanced Research Labs) research labs.

The present manuscript has not been previously submitted/published and is not currently in press, under review or being considered for publication by another journal. Therefore, we would like you to evaluate it for publication and we would be honored in case it will be taken into consideration.

On behalf of all the authors

Yours sincerely.

Michele Suman

Parma, 6th April 2021

Dr. Michele Suman, PhD

Barilla G.R. F.lli SpA

Research, Development & Quality

Food Safety & Authenticity Research Manager

Food Safety Fellow Technical Ladder

Adjunct Professor of AgriFood Authenticity at Catholic University Sacred Heart – Milan/Piacenza

Chair ILSI Process Related Compounds & Natural Toxins Task Force

Chair Italian National Normative Organization (UNI) - Food Authenticity Commission

Scientific Board Member Italian Chemistry Society-Food Chemistry Inter-divisional Group

Via Mantova 166 - 43100 Parma (Italy)

☎ phone +39 0521 262332

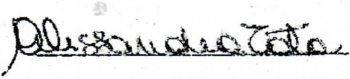

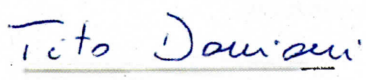
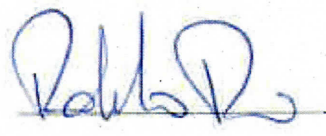

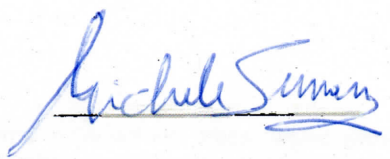
☎ mobile +39 3386938349

✉ mail michele.suman@barilla.com

🌐 web www.barillagroup.com

This statement is signed by all the authors to indicate agreement that the above information is true and correct (a photocopy of this form may be used if there are more than 10 authors):

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

Author's name (typed)	Author's signature	Date
ALESSANDRA TATA		04/03/2021
ANDREA MASSARO		04/03/2021
TITO DANIANI		04/03/2021
ROBERTO PIRO		04/04/2021
CHIARA DALL'ASTA		04/03/2021
MICHELE SUMAN		10/03/2021
_____	_____	_____
_____	_____	_____
_____	_____	_____
_____	_____	_____

Conflicts of Interest Statement

Manuscript title: ADVANTAGES AND DISADVANTAGES OF DATAFUSION OF
LC-MS, GC-IMS AND FGC-ENOXE
~~MULTIPLE MASS SPECTROMETRIC~~ TECHNIQUES IN THE AUTHENTICATION
OF EXTRA VIRGIN OLIVE OIL

The authors whose names are listed immediately below certify that they have NO affiliations with or involvement in any organization or entity with any financial interest (such as honoraria; educational grants; participation in speakers' bureaus; membership, employment, consultancies, stock ownership, or other equity interest; and expert testimony or patent-licensing arrangements), or non-financial interest (such as personal or professional relationships, affiliations, knowledge or beliefs) in the subject matter or materials discussed in this manuscript.

Author names:

ALESSANDRA TATA
ANDREA TASSARO
TITO DANIANI
ROBERTO PIRO
CHIARA DALL'ASTA
MICHELE SODAN

The authors whose names are listed immediately below report the following details of affiliation or involvement in an organization or entity with a financial or non-financial interest in the subject matter or materials discussed in this manuscript. Please specify the nature of the conflict on a separate sheet of paper if the space below is inadequate.

Author names:



To the kind attention of:
Editor

Dr. Q. Rao, PhD

(Food Science; Food Chemistry; Food Quality; Food Safety)
Florida State University College of Human Sciences Nutrition,
Food & Exercise Sciences, Tallahassee, Florida, United States of America
Food Control - Elsevier

Parma, 14th October 2021

Dear Editors,

with reference to the manuscript entitled "Detection of soft-refined oils in extra virgin olive oil using data fusion approaches for LC-MS, GC-IMS and FGC-Enose techniques: the winning synergy of GC-IMS and FGC-Enose", which we have submitted to your attention, we would like, as requested, to indicate the correspondent suggested highlights:

Highlights

- Extra virgin olive oil (EVOO) can be adulterated by mixing it with soft refined oils (SROO)
- LC-MS, GC-IMS and FGC-ENose were evaluated for their fraud detection potentialities
- Low-level and mid-level data fusion of those analytical dataset were performed
- The discriminatory capability of the two merged GC-based techniques was significantly improved
- Combining GC-based techniques, data fusion and a PLS-DA-SVM strategy provides a new framework for effective authentication of EVOO

Please do not hesitate to contact me for any other needs.

Best Regards,
Yours sincerely,
Michele Suman

Dr. Michele Suman, PhD

Barilla G.R. F.lli SpA
Research, Development & Quality
Food Safety & Authenticity Research Manager
Food Safety Fellow Technical Ladder
Adjunct Professor of AgriFood Authenticity at Catholic University Sacred Heart – Milan/Piacenza
Chair ILSI Process Related Compounds & Natural Toxins Task Force
Chair Italian National Normative Organization (UNI) - Food Authenticity Commission
Scientific Board Member Italian Chemistry Society-Food Chemistry Inter-divisional Group
Via Mantova 166 - 43100 Parma (Italy)
phone +39 0521 262332
mobile +39 3386938349

✉ mail michele.suman@barilla.com
🌐 web www.barilagroup.com



To the kind attention of:
 Prof. Andrea Armani, DVM, PhD,
 Dipartimento di Scienze Veterinarie
 Viale delle Piagge, 2, 56124 Pisa (PI)
 e-mail: andrea.armani@unipi.it
 Editor Food Control / Elsevier

Ms Ichiko Charis Howells
 On Behalf of the Editorial Board - Food Control

Authors:

Alessandra Tata^a, Andrea Massaro^a, Tito Damiani^b, Roberto Piro^a, Chiara Dall'Asta^b, Michele Suman^{c,d}*

^a Istituto Zooprofilattico Sperimentale delle Venezie, Laboratorio di Chimica Sperimentale, Vicenza, Italy

^b Department of Food and Drug, University of Parma, Parma, Italy

^c Department of Analytical Food Science, Barilla G. e R. Fratelli S.p.A., Parma, Italy

^d Department for Sustainable Food Process, Catholic University Sacred Heart, Piacenza, Italy

Dear Editor,

with reference to the manuscript entitled "Detection of soft-refined oils in extra virgin olive oil using data fusion approaches for LC-MS, GC-IMS and FGC-Enose techniques: the winning synergy of GC-IMS and FGC-Enose", which we have submitted to your attention, we would like, as requested, to make the following CRediT Statements:

CRediT author statement

Terms, Definition, Conceptualization: *Tata, Massaro,*

Ideas, formulation or evolution of overarching research goals and aims: *Tata, Piro, Dall'Asta, Suman*

Methodology, Development or design of methodology; creation of models: *Tata, Massaro*

Validation, Verification, whether as a part of the activity or separate, of the overall replication/ reproducibility of results/experiments and other research outputs: *Tata, Massaro, Damiani*

Application of statistical, mathematical, computational, or other formal techniques to analyze or synthesize study data: *Tata, Massaro*

Investigation, Conducting a research and investigation process, specifically performing the experiments, or data/evidence collection: *Tata, Massaro, Damiani*

Writing - Original Draft, Preparation, creation and/or presentation of the published work, specifically writing the initial draft: *Tata, Massaro*

Writing - Review & Editing: *Damiani, Piro, Dall'Asta, Suman*

Preparation, creation and/or presentation of the published work, specifically visualization/ data presentation: *Tata, Massaro*

Supervision, Oversight and leadership responsibility for the research activity planning and execution, including mentorship external to the core team: *Piro, Dall'Asta, Suman*

Project administration, Management and coordination responsibility for the research activity planning and execution: *Piro, Dall'Asta, Suman*

Please do not hesitate to contact me for any other needs.

Yours sincerely,

On behalf of all the authors

Michele Suman

Parma, 14th October 2021

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

Dr. Michele Suman, PhD
Barilla G.R. F.Ili SpA
Research, Development & Quality
Food Safety & Authenticity Research Manager
Adjunct Professor of AgriFood Authenticity at Catholic University Sacred Heart – Milan/Piacenza
Chair ILSI Process Related Compounds & Natural Toxins Task Force
Chair Italian National Normative Organization (UNI) - Food Authenticity Commission
Via Mantova 166 - 43100 Parma (Italy)
mobile +39 3386938349
mail michele.suman@barilla.com
web www.barillagroup.com



To the kind attention of:
Prof. Andrea Armani, DVM, PhD,
Dipartimento di Scienze Veterinarie
Viale delle Piagge, 2, 56124 Pisa (PI)
e-mail: andrea.armani@unipi.it
Editor Food Control

Ms Ichiko Charis Howells
On Behalf of the Editorial Board - Food Control
Food Control - Elsevier

Authors:

Alessandra Tata^a, Andrea Massaro^a, Tito Damiani^b, Roberto Piro^a, Chiara Dall'Asta^b, Michele Suman^{c,d}*

^a Istituto Zooprofilattico Sperimentale delle Venezie, Laboratorio di Chimica Sperimentale, Vicenza, Italy

^b Department of Food and Drug, University of Parma, Parma, Italy

^c Department of Analytical Food Science, Barilla G. e R. Fratelli S.p.A., Parma, Italy

^d Department for Sustainable Food Process, Catholic University Sacred Heart, Piacenza, Italy

Title:

“Detection of soft-refined oils in extra virgin olive oil using data fusion approaches for LC-MS, GC-IMS and FGC-Enose techniques: the winning synergy of GC-IMS and FGC-Enose”

Answers to the comments and suggestions from the reviewers – FOODCONT-D-21-00923R1

Dear Editor,

with reference to your opinion on publication of the present work, we would like to thank again for the valuable final review we received. We are honored that this submitted paper can be accepted for publication based on last fine tunings accordingly to the reviewer(s)' comments.

Therefore, the manuscript has been modified according to these reviewers' requests. The detailed responses to the comments and suggestions are reported here below.

Reviewer 1 Comments:

The manuscript can be accepted after this revision, but the title should be changed since, it seems that confirm the utility of data fusion of data obtained from LC-MS, GC-IMS and FGC-Enose techniques for the detection of soft refined oils in extra virgin olive oil. The authors should revise the manuscript because the way of presenting the results can be confusing since till the end of the manuscript it cannot be found that the combination of GC-IMS and FGC-Enose fingerprints using a low-level data fusion approach is the most powerful classification tool.

Response: The reviewer has raised an interesting point, therefore, we modified the title accordingly and we made clear in the abstract that the low-level combination of GC-IMS and FGC-Enose is the most powerful.

Resultant changes to the title: “Detection of soft-refined oils in extra virgin olive oil using data fusion approaches for LC-MS, GC-IMS and FGC-Enose techniques: the winning synergy of GC-IMS and FGC-Enose”

Resultant changes to the abstract: “The most promising results were achieved by the low-level data fusion of GC-IMS and FGC-Enose data.”

Reviewer 3 Comments:

I really appreciate that all my comments addressed in the first review have been properly explained by the authors.

Yet, I regret to say that I do not agree with the authors comment: "Having a small dataset (60 samples) due to reasons explained above, therefore, the proportion of our data split was based on the concept that the more training data we have, the better our model will be. In other words, big training data maximizes the performance of the model and provides higher confidence in the resulting accuracy".

The key point for both the training and test set is to be representative of the case under study.

Regarding the training set, it should contain as many samples as required to proper cover the data variability. Let's us say (just as an example) that with 20 samples all the variability is considered, therefore 20 samples are enough. There are several papers/algorithms that deal with training sample selection, such as Kennard-Stone, PCA score distribution, etc. The final number of training samples is strong depending on the sample/data distribution, whether it is homogeneous or heterogeneous. I am aware that it is not a simple decision, but models build with lower number of samples (in order to increase the test set) might be checked. Test set is used to check the performance of the model, if not enough test samples are used, the performance values based on the test set are not reliable. If that the case, (as it happens in that paper with only 6 test samples) then the best option is to used cross-validation instead of an independent test set.

Response: We thank the reviewer for raising this interesting point. Actually, we tested the models with the same independent samples used in the previous studies from Damiani et al 2020 and Cavanna et al. 2020. Indeed, the three authentic EVOO samples of the test set were previously selected with a Kennard-Stone algorithm, while the other three "NOT EVOO" samples (DEO3, DEO_DEA2, and Mix D) chosen with the aim to predict both pure adulterated samples and mixtures. In order to clarify this point, we added this info to the manuscript. We also removed from the manuscript the comment related to the "concept that big training data maximizes the performance of the model and provides higher confidence in the resulting accuracy on test set"

Resultant changes to material and methods: "The test set was comprised of three authentic EVOO (CP-30, CP-31, CP-32) and three SROO (DEO3, DEO_DEA2, MIX_D) as previously done by Cavanna et al 2020 and Damiani et al 2020. The three authentic EVOO samples of the test set were selected with a Kennard-Stone algorithm, while the other three "NOT EVOO" samples (DEO3, DEO_DEA2, and Mix D) chosen with the aim to predict both pure adulterated samples and mixtures (Cavanna et al 2020)."

Parma, 14th October 2021

On behalf of all the authors. Best regards.

Dr. Michele Suman, PhD

Barilla G.R. F.Ili SpA

Research, Development & Quality

Food Safety & Authenticity Research Manager

Adjunct Professor of AgriFood Authenticity at Catholic University Sacred Heart – Milan/Piacenza

Chair ILSI Process Related Compounds & Natural Toxins Task Force

Chair Italian National Normative Organization (UNI) - Food Authenticity Commission

Via Mantova 166 - 43100 Parma (Italy)

☎ mobile +39 3386938349

✉ mail michele.suman@barilla.com

🌐 web www.barillagroup.com

1 **Detection of soft-refined oils in extra virgin olive oil using data fusion approaches for**
2 **LC-MS, GC-IMS and FGC-Enose techniques: the winning synergy of GC-IMS and**
3 **FGC-Enose**

4 Alessandra Tata ^a, Andrea Massaro^a, Tito Damiani^b, Roberto Piro^a, Chiara Dall'Asta^b, Michele
5 Suman^{c,d*}

6 ^aIstituto Zooprofilattico Sperimentale delle Venezie, Laboratorio di Chimica Sperimentale,
7 Vicenza, Italy

8 ^b Department of Food and Drug, University of Parma, Parma, Italy

9 ^c Department of Analytical Food Science, Barilla G. e R. Fratelli S.p.A., Parma, Italy

10 ^d Department for Sustainable Food Process, Catholic University Sacred Heart, Piacenza, Italy

11

12

13 *Corresponding author: Dr. Michele Suman Advanced Research Laboratory, Barilla G. e R.

14 Fratelli S.p.A., Parma, Italy

15 email: michele.suman@barilla.com; michele.suman@unicatt.it

16

17

18 **Abbreviations:** LC-MS: liquid chromatography-mass spectrometry; GC-IMS: gas-
19 chromatography ion mobility spectrometry; FGC-Enose: flash gas-chromatography electronic
20 nose; EVOO Extra Virgin Olive Oil; SROO: soft-refined olive oil; PLS-DA, Partial Least
21 Squared Discriminant Analysis; SVM, support vector machine; ROC curve, Receiver
22 Operating Characteristic curve; AUC area under the curve.

23

24 **Keywords:** Soft-deodorization, deacidification, FGC-Enose, adulteration, GC-IMS, LC-MS,
25 soft-refined olive oil (SROO), data fusion

Abstract

Extra virgin olive oil (EVOO) is frequently adulterated by mixing it with soft refined oils (SROO). The differentiation of EVOO from its blends with SROO is not possible with the most common approaches, and, for this reason, the discriminating power of liquid chromatography-high resolution mass spectrometry (LC-MS), gas-chromatography ion mobility spectrometry (GC-IMS) and flash gas-chromatography electronic nose (FGC-Enose) was examined previously. Here, the combination of the above-mentioned techniques for an improvement in classification power of the methods is explored.

A total of 43 commercial EVOOs and 18 illegal mixtures of SROO with EVOO were previously analysed by LC-(+/-)MS, GC-IMS and FGC-Enose. Low-level and mid-level data fusion of the four datasets were performed. The merged unique fingerprints were submitted to partial least squared discriminant analysis (PLS-DA), and the extrapolated most informative variables were used to build support vector machine (SVM) classifiers. Statistical indicators were calculated and compared to find out the best classifier. The results of PLS-DA-SVM strategies on the combination of datasets demonstrated that, after low-level data fusion, the discriminatory capability of the two merged GC-based techniques was remarkably improved as compared to the individual techniques. This indicates that merging the datasets before PLS-DA better retrieves the most informative variables and, thus, enhances group separation and classification of unknowns. The combination of LC(+/-)MS datasets, both by mid- and low-level data fusion, did not show significant enhancement in terms of discrimination of EVOO from SROO as compared to the individual LC(+)MS matrix. The low-level combination of the four datasets (LC(+/-)MS, GC-IMS, FGC-Enose) was successful, although this laborious option is not a viable path in industry quality assurance.

This study primarily provides new paths for the authentication of EVOO, taking advantage of merging multimodal LC-(+/-)MS, GC-IMS and FGC-Enose data, with consequent

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51 improvement in the performances of the classification models. The most promising results were
52 achieved by the low-level data fusion of GC-IMS and FGC-Enose data.

53

54 **1. Introduction**

55 Due to its high economic value and unique sensorial and nutritional characteristics, extra virgin
56 olive oil (EVOO) is considered at high risk of fraud(Casadei *et al.*, 2021). Recently, more
57 sophisticated adulterations have been developed. The mixtures of EVOO with soft deacidified
58 and soft deodorized olive oils are considered the most critical frauds because they are not easily
59 detectable by regular methods(Conte *et al.*, 2020).

60 The detection of soft refined products in EVOO has been recently attempted by near infrared
61 (NIR) spectroscopy (Gertz, Matthäus and Willenberg, 2020) and diacylglycerol
62 determination(Gómez-Coca *et al.*, 2020). Recently, the adulteration of EVOO with soft-refined
63 olive oil (SROO) has raised the interest of our research group, as four non-targeted methods
64 capable of detecting this fraud were developed and validated separately; these were liquid
65 chromatography-mass spectrometry (LC-MS) in positive and negative ion mode, gas-
66 chromatography ion mobility spectrometry (GC-IMS) and flash gas-chromatography electronic
67 nose (FGC-Enose) (Damiani *et al.*, 2020; Cavanna *et al.*, 2020).

68 Data fusion is a chemometric technique that merges the outcomes of multiple analytical sources.
69 It has recently emerged as an attractive means to enhance the prediction power of a model for
70 food authentication (Callao and Ruisánchez, 2018; Hu *et al.*, 2019; Márquez *et al.*, 2016). Low-
71 level data fusion is a valuable chemometric strategy capable of concatenating multiple datasets
72 and improving the classification performances by retrieving the discriminative variables from
73 different techniques (Andrade *et al.*, 2021). Mid-level data fusion aims at merging datasets by
74 reducing their high dimensionality and teasing out solely the most informative variables capable
75 of codifying each group in the study (Jandric *et al.*, 2021; Tata *et al.*, 2021; Riuzzi *et al.*, 2021).

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

76 Data fusion models were applied to EVOO for the detection of its adulteration with vegetable
77 oils (Schwolow *et al.*, 2019; Li, Xiong and Min, 2019), the assessment of its geographical origin
78 (Casale *et al.*, 2010a; Casale *et al.*, 2012; Pizarro *et al.*, 2013; Nescatelli *et al.*, 2014; Bajoub *et*
79 *al.*, 2017) and the reveal of sensory defects (Borràs *et al.*, 2016).

80 Most of the common data fusion models applied to EVOO have merged data from analytical
81 techniques that provide similar information, such as Raman, near infrared and medium infrared
82 spectroscopies (Casale *et al.*, 2010b; Li, Xiong and Min, 2019; Pizarro *et al.*, 2013; Bevilacqua
83 *et al.*, 2013; Jiménez-Carvelo, Lozano and Olivieri, 2019; Casale *et al.*, 2012; Bragolusi *et al.*,
84 2021) or chromatographic profiles recorded at three different wavelengths (Nescatelli *et al.*,
85 2014). On the other hand, data fusion could be very useful when complementary information is
86 fused and included in one unique model (Schwolow *et al.*, 2019; Assis *et al.*, 2019; Borràs *et*
87 *al.*, 2016; Casale *et al.*, 2010a; Casale *et al.*, 2007).

88 In the present study, data from the three complementary techniques, each of them characterized
89 by distinct information (volatile and non-volatile chemical profiles) were merged by low and
90 mid-level data fusion for the discrimination of authentic EVOO and fraudulent SROO blends.

91 Although promising results have been achieved in food authentication assessment (Damiani *et*
92 *al.*, 2020; Cavanna *et al.*, 2020), reports on the combination of data from different mass
93 spectrometric techniques for the improvement of detection of the SROO blends are still limited.

94 The present study aimed to evaluate the enhanced prediction power obtained by low-level and
95 mid-level data fusion and outline any possible disadvantages.

96 The comparison was carried out through the estimation of statistical indicators, i.e., accuracy,
97 sensitivity, specificity, for a training set and probability of predictions for a set of validation
98 samples. To the best of our knowledge, this is the first study exploring data fusion strategies for
99 the detection of SROO blends in EVOO.

100

101

102 **2. Materials and methods**

103 *2.1 Dataset collection and analysis*

104 The datasets used for this study were acquired in our previous studies (Damiani *et al.*, 2020;
105 Cavanna *et al.*, 2020). Therefore, all the details about sample collection and analyses are
106 reported in detail in our previous publications.

107 Briefly, a total of 43 commercial Italian EVOOs, obtained over three harvest seasons (i.e.,
108 2015/2016, n = 18; 2016/2017, n = 8; 2017/2018, n = 17), were considered as authentic samples.
109 In addition, soft-deodorization and deacidification were carried out on commercial virgin and
110 lampante olive oils to create counterfeit soft-refined samples (SROO).

111 In order to create counterfeited samples potentially compliant with the legislation, the official
112 EVOO physic-chemical quality parameters(Regulation, 2016) were analysed in these refined
113 oils.

114 Based on the obtained results, 18 illegal blends were prepared at different percentages by
115 mixing the so-obtained SROO with authentic EVOOs randomly chosen from the sample set.

116 Authentic and counterfeit olive oil samples were analysed using three different techniques,
117 namely GC-IMS, FGC-Enose, and LC-(+/-)MS.

118 Partially satisfactory classification models were obtained from the separate volatile profiles
119 (Damiani *et al.* 2020) and from the LC-MS profiles (Cavanna *et al.* 2020).

120

121 *2.2 Data fusion strategies and multivariate statistical analysis*

122 In order to improve the prediction of authentic and adulterated EVOO, LC-(+/-)MS, GC-IMS
123 and FGC-Enose data were merged via both low level and mid-level data fusion strategies using
124 RStudio 3.6.2 and Metabonalyt 5.0 web platform.

125

126 2.2.1 Low-level data fusion

1
2 127 Each dataset was pre-processed by removing the C¹³ isotopes and the *m/z* ions with more than
3
4 128 75% of non-acquired intensities (missing values) across all the samples. Each dataset was
5
6
7 129 normalized by sum and scaled by Pareto. Each pre-processed dataset was split into training set
8
9
10 130 (55 samples) and test set (6 samples). The test set was comprised of three authentic EVOO (CP-
11
12 131 30, CP-31, CP-32) and three SROO (DEO3, DEO_DEA2, MIX_D) as previously done by
13
14 132 Cavanna *et al* 2020 and Damiani *et al* 2020. The three authentic EVOO samples of the test set
15
16
17 133 were selected with a Kennard-Stone algorithm, while the other three “NOT EVOO” samples
18
19 134 (DEO3, DEO_DEA2, and Mix D) chosen with the aim to predict both pure adulterated samples
20
21
22 135 and mixtures (Cavanna *et al* 2020).

23
24 136 Low-level data fusions of: i) two LC-MS instrumental ion modes; ii) GC-IMS and FGC-Enose,
25
26
27 137 and; iii) multimodal LC-MS and FGC-Enose and GC-IMS were performed.

28
29 138 The pre-processed signals of each training set were simply concatenated, mean-centered and
30
31
32 139 processed as a unique fingerprint of the samples.

33
34 140 The merged training sets were submitted to the supervised partial least squared discriminant
35
36
37 141 analysis (PLS-DA) with the aim of extrapolating the most informative variables.

38
39 142 The PLS-DA variables with coefficients >55 were retained and used to construct the linear
40
41
42 143 SVM classification models which was validated on the merged test set.(Massaro *et al.*, 2021)

43
44 144 The criterion used to extrapolate the most significant features was based on the inspection of
45
46
47 145 PLS-DA coefficient plot (not shown) reporting the informative variables in a descending order
48
49
50 146 (from the one with highest coefficient to that with the lowest).

51
52 147 The "elbow" of the graph, where the coefficient of the informative variables leveled off, was
53
54
55 148 considered as limit point.

56
57 149 The variables placed to the right of this point, corresponding to coefficient equal to 55, were
58
59
60 150 retained as significant.

151

1
2 152 *2.2.2 Mid-level data fusion*
3

4 153 Briefly, each pre-processed dataset (split into training and test sets) was submitted to supervised
5
6
7 154 PLS-DA. We selected the first five components of the PLS-DA of each dataset and we retrieved
8
9
10 155 from them the most significant variables. As recommended by Hair et al (Hair *et al.*, 2006) only
11
12 156 the ions with absolute values for PLS-DA loadings >0.3 were retained and used to build the
13
14 157 SVM classification models. Further details of the mid-level data fusion strategy adopted can be
15
16
17 158 found elsewhere (Massaro *et al.*, 2021)
18

19 159

20
21
22 160 *2.2.3 Validation of the classification model*
23

24 161 Support vector machine (SVM) classification models were built with the extrapolated
25
26 162 molecular features using the Biomarker Analysis section of Metaboanalyst 5.0 after low-level
27
28
29 163 and mid-level data fusions. Each SVM model was cross-validated by Monte Carlo cross
30
31 164 validation (MCCV) using a repeated, balanced sub-sampling procedure. In details, the MCCV
32
33
34 165 split training data in $2/3$ for training the model and $1/3$ for testing it.
35
36 166 For each iteration, the training/test split was different. In the first iteration, the model was tested
37
38
39 167 on training data and test errors were calculated. After 100 iterations, the average of the test
40
41 168 errors was determined and sensitivity (true positive rate), specificity (true negative rate) and
42
43
44 169 accuracy were calculated.

45
46 170 The overall prediction power of the SVM models was estimated based on the area under the
47
48 171 curve (AUC) of the receiver operating characteristic (ROC) curve. Finally, the SVM models
49
50
51 172 were tested for their ability to classify six samples from the merged test set that was withheld
52
53 173 previously.

54
55
56 174

57
58 175
59
60
61
62
63
64
65

176 3. Results

177 The combination of analytical methods (LC(+/-)-MS, GC-IMS and FGC-Enose) was assessed
178 to evaluate possible improvements in discriminating of EVOO from their blends with SROO.
179 First, a low-level data fusion of LC(+/-)-MS datasets was conducted. The resultant global data
180 matrix was split into training and test sets. The training set was submitted to multivariate
181 statistical analysis by means of PLS-DA (**Figure 1A**). A good trend of separation was
182 observed, with Component1 and Component2 capable of explaining 35.7 % and 10.2% of the
183 data variance, respectively (**Figure 1B**). The m/z values and associated retention times with a
184 higher discriminatory capacity (coefficient >55) were retained and used to build a SVM
185 classifier. The SVM model was cross-validated by MCCV on the training set (**Figure 1A**, right
186 side) with accuracy, sensitivity and specificity reaching 0.94, 0.93 and 0.95 respectively (**Table**
187 **1**). The ROC curve, a graph plotting true positive and false positive rates of the SVM
188 classification model at all classification thresholds, showed an AUC equal to 0.97 (**Figure 1C**).
189 These excellent accuracy, sensitivity and specificity parameters increased in the blinded
190 verification which was able to correctly classify 6/6 samples. The results of the predictions on
191 the test set, and the correlated probabilities, can be visualised in **Table S1** of the supplementary
192 material. The averaged probability of all samples is above 96%.
193 Subsequently the combination of GC-IMS and FGC-Enose approaches by low-level data fusion
194 was evaluated. To this aim, GC-IMS and FGC-Enose datasets were both split into training and
195 test sets, concatenated, and PLS-DA was performed on the fused data (**Figure 2A**). The PLS-
196 DA score plot is reported in **Figure 2B** with the EVOO samples grouped decently by
197 Component1 and Component2. The SVM model, built with the selected variables, showed an
198 accuracy, sensitivity and specificity on the training set of 0.96, 0.93 and 0.97, respectively, and
199 an AUC of the ROC curve equal to 0.99 (**Table 1** and **Figure 2C**). The SVM correctly classified

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

200 6/6 samples in the test set with an averaged probability above 93%, although with a low
201 probability of predicting one sample (MIX_D) (**Table S2**).

202 Finally, the LC(+/-)MS, GC-IMS and FGC-Enose datasets were merged by a low-level data
203 fusion approach. Compared to the previous two techniques, the score plot showed improved
204 clustering of the two groups (authentic and non-authentic EVOO) in the study, with the first
205 and second components, C1 and C2, explaining 26.9% and 11.0 % of the total variance of the
206 model, respectively (**Figure 3**).

207 The results of the cross-validation of SVM, built with the variables with coefficient >55
208 retrieved from fused-PLS-DA, are shown in **Table 1**. In this case, the SVM model built with
209 the combination of the most informative variables of the three techniques reached an accuracy,
210 sensitivity and specificity on the training set of 0.96, 0.93 and 0.97 respectively and an AUC of
211 the ROC curve equal to 0.98 (**Table 1** and **Figure 3C**). The SVM correctly classified 6/6
212 samples in the test set with an averaged probability above 93%, although with a low probability
213 of predicting the sample MIX_D (**Table S3**).

214 Mid-level data fusion was also attempted for the alternated combination of all four datasets
215 (**Figures S1, S2 and S3** of the supplementary material), with less satisfactory results, especially
216 in terms of the ROC curve in cross-validation and the probability of predictions for the test set,
217 as compared to the low-level data fusion.

218 For this reason, a summary of mid-level data fusion results of the cross-validation of the SVM
219 models and their validation on the merged test set are only shown in the supplementary material
220 (**Tables S4-S7**).

221 Note that the best classification performances in this case were achieved by the mid-level data
222 fusion of the two LC matrices. (**Table S4**). With the mid-level data fusion of the four datasets
223 less trustable classifier was obtained (**Table S4**).

224
225

226 4. Discussion

227
228 In previous studies, GC-IMS, FGC-Enose and LC(+/-)MS datasets were statistically
229 analysed separately. Headspace-based techniques (i.e., GC-IMS, FGC-Enose) showed great
230 potential as rapid screening platforms and exhibited remarkable reproducibility over the time;
231 yet, the EVOO's volatile fingerprint seemed to be heavily affected by chemical changes
232 occurring in ordinary shelf-life conditions. On the other hand, LC-MS enabled the identification
233 of fraud-specific markers; however, it suffers of limited sensitivity (i.e., fraud detected at >40%
234 SROO addition). In this study, we want to explore the possibility of merging the data and
235 evaluate possible improvements in the discrimination of genuine EVOO from SROO. In
236 particular, the main aim was to provide a robust data fusion approach to be coupled with quick
237 fingerprint analysis that could be applied in an industrial environment for rapid EVOO
238 authentication. Low-level fusion was first used to pick up correlations between variables of
239 different blocks of data. Low-level fusion is based on the simple concatenation of data to which
240 a single model is applied to pick up correlations between variables belonging to different
241 datasets (Biancolillo *et al.*, 2014; Borràs *et al.*, 2015). It has the limitations of high volume of
242 features, which is difficult to handle, and the possible predominance of one data source over
243 the others. In order to exclude this possible issue, we checked the number of variables of each
244 dataset. We had a thousand variables in each LC-MS dataset and a total of one-hundred thirty
245 variables in the GC matrices. Besides the predominance of the LC-MS source, the difference

246 in block sizes did not affect the PLS-DA weighting of the GC variables that appear as the most
247 significant features in low-level data fusion of the four datasets.

248 On the other hand, mid-level data fusion is characterized by an initial high dimensional data
249 reduction, by means of supervised or unsupervised tools capable of extracting the most
250 informative variables from each separate dataset (Pirro *et al.*, 2014; Borràs *et al.*, 2015).

251 After both low-level and mid-level data fusion, PLS-DA–SVM strategies were applied to
252 concatenated datasets to obtain classification rates for cross-validation and validation on the
253 test set. The SVM models that followed the mid-level data fusion provided less powerful
254 classification, and for this reason, results were included in SI only, and are not discussed further.

255 In the individual techniques, the LC(+)-MS profiles showed high accuracy, R^2 and Q^2 (Cavanna
256 *et al.*, 2020). The accuracy is the capability of the model to correctly classify the samples, the
257 R^2 parameter indicates the goodness of fit of the PLS-DA model (how well it explains the
258 dataset) and Q^2 provides a measure of exactness between the predicted and actual data (Triba
259 *et al.*, 2015; Worley and Powers, 2013). Further details are reported elsewhere (Anderssen *et al.*,
260 2006; Westerhuis *et al.*, 2008). It is worth noting that LC-MS is a highly informative
261 technique that can be used for the identification of chemical markers to be further used in target
262 analysis. While being extremely powerful, this approach is costly and requires high-level
263 laboratory skills. Its application in an industrial environment is, therefore, suggested only for
264 explorative analysis or for confirmatory purposes, whereas it cannot be applied for routine
265 controls. Although in the present study we used linear SVM as the classifier instead of PLS-
266 DA (PLS-DA was employed just to extrapolate the most informative variables used to build the
267 classification model), it does not seem that either the mid or the low-level combination of the
268 two datasets resulted in improvements to the classification figures of merit. However, the
269 performance obtained from the fusion of the LC-(+/-)-MS can be regarded as a benchmark for
270 evaluating the discrimination potential shown by the data fusion applied to volatile fingerprints.

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

271 In the individual techniques, the soft independent modelling by class analogy (SIMCA) models
272 developed on the GC-IMS and FGC-Enose fingerprint datasets were able to correctly recognize
273 the SROO blends as non-authentic products, even at the lowest adulteration percentage (i.e.
274 10%) (Damiani *et al.*, 2020). Only one EVOO sample was wrongly recognized as not EVOO,
275 confirming the high potential of the two separately employed techniques (Damiani *et al.*, 2020).
276 After the application of low-level data fusion, the SVM model developed herein achieved
277 extremely high sensitivity, specificity and accuracy with fully correct predictions for the test
278 set. In contrast to our previous study, we were able to include EVOO 15/16 (CP_1-CP_12), oils
279 that negatively altered the performance of our previous SIMCA model (Damiani *et al.*, 2020).
280 Therefore, the chemometric approach followed in the present work, and based on the fusion of
281 both volatile fingerprint datasets, showed an improvement in the discrimination potential of the
282 model compared to each technique alone. This fused dataset approach is able to overcome the
283 difficulties related to partial overlap of EVOO's chemical features in the volatile fraction
284 characteristics, thereby differentiating oil resulting from fraudulent practice from naturally aged
285 oil subjected to long storage conditions.
286 When compared to the SVM model obtained by fusing LC-(+)MS and LC-(-)MS datasets, the
287 quality parameters on the training set were slightly higher for the GC-fused model, while the
288 probability of correct prediction in the validation test set was lower (0.93 versus 0.96 for GC-
289 fused and LC-fused model, respectively), even though the same outcomes for sample
290 classification were seen.
291 Overall, it can be concluded these the two models are almost comparable in terms of
292 classification performances, although the GC-fused model showed undeniable advantages in
293 terms of cost-effectiveness and ease of handling in an industrial quality control routine
294 approach.

1
2 295 On the other hand, it must be underlined that MS offers the opportunity to identify the chemical
3
4 296 markers responsible for classification, and to monitor them over time. Therefore, its superior
5
6 297 use for explorative and confirmatory purposes is without question.

7 298 To gain a comprehensive overview of the potential of data fusion in EVOO classification, all
8
9 299 four datasets were fused, and the resultant model was compared to the previous one in terms of
10
11 300 performance.

12
13 301 In this case, the statistical indicators obtained in both mid- and low-level data fusion were still
14
15 302 satisfactory, but lower than those obtained from the combination of the two GC-based
16
17 303 approaches. We observed a low AUC when running the mid-level data fusion of the four
18
19 304 datasets (**Table S4**).

20
21 305 On the other hand, considering the analytical and chemometric efforts required to collect and
22
23 306 fuse datasets from four different techniques, with little to no improvement obtained in the
24
25 307 overall model, this approach is far from offering a useful solution currently applicable within
26
27 308 industrial production monitoring.

28
29 309 In conclusion, the combination of GC-IMS and FGC-Enose fingerprints using a low-level data
30
31 310 fusion approach is the most powerful classification tool we know of to date that could be used
32
33 311 for identifying soft refinement of EVOO in an industrial quality assurance setting. Notably, this
34
35 312 approach is based on datasets obtained using cost-effective and easy-to-handle techniques.

36
37 313

38
39 314

40 41 315 **5. Conclusion**

42
43 316 Data fusion strategies to authenticate EVOO were tested taking into account of LC-(+/-)MS,
44
45 317 GC-IMS and FGC-Enose analytical data. Specifically, low-level data fusion of GC-IMS and
46
47 318 FGC-ENose datasets produces an optimal model for classifying SROO and EVOO with an
48
49 319 overall accuracy of 0.96 and the advantage of the rapid acquisition of the spectra.

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

320 This approach, combining GC-based techniques, data fusion and a PLS-DA-SVM strategy,
321 likely provides a new framework for the authentication of EVOO in a possible industrial quality
322 assurance setting.
323
324
325

326 References

- 327 Andrade, D. F., de Almeida, E., de Carvalho, H. W. P., Pereira-Filho, E. R.
328 and Amarasiriwardena, D. (2021) 'Chemical inspection and elemental analysis
329 of electronic waste using data fusion - Application of complementary
330 spectroanalytical techniques', (1873-3573 (Electronic)).
- 331 Assis, C., Pereira, H. V., Amador, V. S., Augusti, R., de Oliveira, L. S.
332 and Sena, M. M. (2019) 'Combining mid infrared spectroscopy and paper spray
333 mass spectrometry in a data fusion model to predict the composition of
334 coffee blends', *Food Chemistry*, 281, pp. 71-77.
- 335 Bajoub, A., Medina-Rodríguez, S., Gómez-Romero, M., Ajal, E. A., Bagur-
336 González, M. G., Fernández-Gutiérrez, A. and Carrasco-Pancorbo, A. (2017)
337 'Assessing the varietal origin of extra-virgin olive oil using liquid
338 chromatography fingerprints of phenolic compound, data fusion and
339 chemometrics', *Food Chemistry*, 215, pp. 245-255.
- 340 Bevilacqua, M., Bucci, R., Magrì, A. D., Magrì, A. L. and Marini, F. (2013)
341 'Data Fusion for Food Authentication. Combining near and Mid Infrared to
342 Trace the Origin of Extra Virgin Olive Oils', *NIR news*, 24(2), pp. 12-15.
- 343 Biancolillo, A., Bucci, R., Magrì, A. L., Magrì, A. D. and Marini, F.
344 (2014) 'Data-fusion for multiplatform characterization of an italian craft
345 beer aimed at its authentication', *Analytica Chimica Acta*, 820, pp. 23-31.
- 346 Borràs, E., Ferré, J., Boqué, R., Mestres, M., Aceña, L. and Busto, O.
347 (2015) 'Data fusion methodologies for food and beverage authentication and
348 quality assessment - A review', *Analytica Chimica Acta*, 891, pp. 1-14.
- 349 Borràs, E., Ferré, J., Boqué, R., Mestres, M., Aceña, L., Calvo, A. and
350 Busto, O. (2016) 'Olive oil sensory defects classification with data fusion
351 of instrumental techniques and multivariate analysis (PLS-DA)', *Food
352 chemistry*, 203, pp. 314-322.
- 353 Bragolusi, M., Massaro, A., Tata, A. and Piro, R. (2021) 'A data fusion
354 model of NIR and RAMAN techniques for the geographical screening of Italian
355 extra virgin olive oil', *NIRItalia online 2021*.
- 356 Callao, M. P. and Ruisánchez, I. (2018) 'An overview of multivariate
357 qualitative methods for food fraud detection', *Food Control*, 86, pp. 283-
358 293.
- 359 Casadei, E., Valli, E., Panni, F., Donarski, J., Farrús Gubern, J., Lucci,
360 P., Conte, L., Lacoste, F., Maquet, A., Brereton, P., Bendini, A. and
361 Gallina Toschi, T. (2021) 'Emerging trends in olive oil fraud and possible
362 countermeasures', *Food Control*, 124, pp. 107902.
- 363 Casale, M., Armanino, C., Casolino, C. and Forina, M. (2007) 'Combining
364 information from headspace mass spectrometry and visible spectroscopy in
365 the classification of the Ligurian olive oils', *Analytica Chimica Acta*,
366 589(1), pp. 89-95.
- 367 Casale, M., Casolino, C., Oliveri, P. and Forina, M. (2010a) 'The potential
368 of coupling information using three analytical techniques for identifying
369 the geographical origin of Liguria extra virgin olive oil', *Food Chemistry*,
370 118(1), pp. 163-170.
- 371 Casale, M., Oliveri, P., Casolino, C., Sinelli, N., Zunin, P., Armanino,
372 C., Forina, M. and Lanteri, S. (2012) 'Characterisation of PDO olive oil
373 Chianti Classico by non-selective (UV-visible, NIR and MIR spectroscopy)
374 and selective (fatty acid composition) analytical techniques', *Analytica
375 Chimica Acta*, 712, pp. 56-63.
- 376 Casale, M., Sinelli, N., Oliveri, P., Di Egidio, V. and Lanteri, S. (2010b)
377 'Chemometrical strategies for feature selection and data compression
378 applied to NIR and MIR spectra of extra virgin olive oils for cultivar
379 identification', *Talanta*, 80(5), pp. 1832-1837.
- 380 Cavanna, D., Hurkova, K., Džuman, Z., Serani, A., Serani, M., Dall'Asta,
381 C., Tomaniova, M., Hajslova, J. and Suman, M. (2020) 'A Non-Targeted High-
382 Resolution Mass Spectrometry Study for Extra Virgin Olive Oil Adulteration
383 with Soft Refined Oils: Preliminary Findings from Two Different
384 Laboratories', *ACS Omega*, 5(38), pp. 24169-24178.

385 Conte, L., Bendini, A., Valli, E., Lucci, P., Moret, S., Maquet, A.,
386 Lacoste, F., Brereton, P., García-González, D. L., Moreda, W. and Gallina
387 Toschi, T. (2020) 'Olive oil quality and authenticity: A review of current
388 EU legislation, standards, relevant methods of analyses, their drawbacks
389 and recommendations for the future', *Trends in Food Science & Technology*,
390 105, pp. 483-493.
391 Damiani, T., Cavanna, D., Serani, A., Dall'Asta, C. and Suman, M. (2020)
392 'GC-IMS and FGC-Enose fingerprint as screening tools for revealing extra
393 virgin olive oil blending with soft-refined olive oils: A feasibility
394 study', *Microchemical Journal*, 159, pp. 105374.
395 Gertz, C., Matthäus, B. and Willenberg, I. (2020) 'Detection of Soft-
396 Deodorized Olive Oil and Refined Vegetable Oils in Virgin Olive Oil Using
397 Near Infrared Spectroscopy and Traditional Analytical Parameters', *European
398 Journal of Lipid Science and Technology*, 122(6), pp. 1900355.
399 Gómez-Coca, R. B., Pérez-Camino, M. d. C., Bendini, A., Gallina Toschi, T.
400 and Moreda, W. (2020) 'Olive oil mixtures. Part two: Detection of soft
401 deodorized oil in extra virgin olive oil through diacylglycerol
402 determination. Relationship with free acidity', *Food Chemistry*, 330, pp.
403 127226.
404 Hu, O., Chen, J., Gao, P., Li, G., Du, S., Fu, H., Shi, Q. and Xu, L.
405 (2019) 'Fusion of near-infrared and fluorescence spectroscopy for
406 untargeted fraud detection of Chinese tea seed oil using chemometric
407 methods', *Journal of the Science of Food and Agriculture*, 99(5), pp. 2285-
408 2291.
409 Jandric, Z., Tchaikovsky, A., Zitek, A., Causon, T., Stursa, V., Prohaska,
410 T. and Hann, S. (2021) 'Multivariate modelling techniques applied to
411 metabolomic, elemental and isotopic fingerprints for the verification of
412 regional geographical origin of Austrian carrots', *Food Chemistry*, 338, pp.
413 127924.
414 Jiménez-Carvelo, A. M., Lozano, V. A. and Olivieri, A. C. (2019)
415 'Comparative chemometric analysis of fluorescence and near infrared
416 spectroscopies for authenticity confirmation and geographical origin of
417 Argentinean extra virgin olive oils', *Food Control*, 96, pp. 22-28.
418 Li, Y., Xiong, Y. and Min, S. (2019) 'Data fusion strategy in quantitative
419 analysis of spectroscopy relevant to olive oil adulteration', *Vibrational
420 Spectroscopy*, 101, pp. 20-27.
421 Massaro, A., Stella, R., Negro, A., Bragolusi, M., Miano, B., Arcangeli,
422 G., Biancotto, G., Piro, R. and Tata, A. (2021) 'New strategies for the
423 differentiation of fresh and frozen/thawed fish: A rapid and accurate non-
424 targeted method by ambient mass spectrometry and data fusion (part A)',
425 *Food Control*, 130, pp. 108364.
426 Márquez, C., López, M. I., Ruisánchez, I. and Callao, M. P. (2016) 'FT-
427 Raman and NIR spectroscopy data fusion strategy for multivariate
428 qualitative analysis of food fraud', (1873-3573 (Electronic)).
429 Nescatelli, R., Bonanni, R. C., Bucci, R., Magri, A. L., Magri, A. D. and
430 Marini, F. (2014) 'Geographical traceability of extra virgin olive oils
431 from Sabina PDO by chromatographic fingerprinting of the phenolic fraction
432 coupled to chemometrics', *Chemometrics and Intelligent Laboratory Systems*,
433 139, pp. 175-180.
434 Pirro, V., Oliveri, P., Ferreira, C. R., González-Serrano, A. F., Machaty,
435 Z. and Cooks, R. G. (2014) 'Lipid characterization of individual porcine
436 oocytes by dual mode DESI-MS and data fusion', *Analytica chimica acta*, 848,
437 pp. 51-60.
438 Pizarro, C., Rodríguez-Tecedor, S., Pérez-del-Notario, N., Esteban-Díez, I.
439 and González-Sáiz, J. M. (2013) 'Classification of Spanish extra virgin
440 olive oils by data fusion of visible spectroscopic fingerprints and
441 chemical descriptors', *Food Chemistry*, 138(2), pp. 915-922.
442 Regulation, E. (2016) '2095 (2016). Amending Regulation of EEC No:
443 2568/91', *Official Journal of the European Communities*, 326, pp. 1-6.
444 Riuzzi, G., Tata, A., Massaro, A., Bisutti, V., Lanza, I., Contiero, B.,
445 Bragolusi, M., Miano, B., Negro, A., Gottardo, F., Piro, R. and Segato, S.

446 (2021) 'Authentication of forage-based milk by mid-level data fusion of
1 447 (+/-) DART-HRMS signatures', *International Dairy Journal*, 112, pp. 104859.
2 448 Schwolow, S., Gerhardt, N., Rohn, S. and Weller, P. (2019) 'Data fusion of
3 449 GC-IMS data and FT-MIR spectra for the authentication of olive oils and
4 450 honeys-is it worth to go the extra mile?', *Anal Bioanal Chem*, 411(23), pp.
5 451 6005-6019.
6 452 Tata, A., Pallante, I., Massaro, A., Miano, B., Bottazzari, M., Fiorini,
7 453 P., Dal Prà, M., Paganini, L., Stefani, A., De Buck, J., Piro, R. and
8 454 Pozzato, N. (2021) 'Serum Metabolomic Profiles of Paratuberculosis Infected
9 455 and Infectious Dairy Cattle by Ambient Mass Spectrometry', *Frontiers in*
10 456 *Veterinary Science*, 7(1214).
11 457

12
13 458
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

Detection of soft-refined oils in extra virgin olive oil using data fusion approaches for LC-MS, GC-IMS and FGC-Enose techniques: the winning synergy of GC-IMS and FGC-Enose.

Alessandra Tata, Andrea Massaro, Tito Damiani, Roberto Piro, Chiara Dall'Asta, Michele Suman

Supplementary material

Table S1. Classification independent adulterated (OTHER) and authentic extra virgin olive oil (EVOO) by the support vector machine (SVM) model, built using informative variables from low-level data fusion of multimodal high performance liquid chromatography-high resolution mass spectrometry (HPLC-HRMS) datasets.

SAMPLES	ACTUAL	PREDICTED	AVAREGED PROBABILITY
CP_30	AUTHENTIC	EVOO	0.96514
CP_31	AUTHENTIC	EVOO	0.92409
CP_32	AUTHENTIC	EVOO	0.99294
DEO3	ADULTERATED	OTHER	0.98983
DEO_DEA_2	ADULTERATED	OTHER	0.9958
MIX_D	ADULTERATED	OTHER	0.91528

Table S2. Classification independent adulterated (OTHER) and authentic extra virgin olive oil (EVOO) by the support vector machine (SVM) model, built using informative variables from low-level data fusion of gas chromatography coupled with ion mobility spectrometry (GC-IMS) and flash gas chromatography electronic nose (FGC-Enose) datasets.

SAMPLES	ACTUAL	PREDICTED	AVAREGED PROBABILITY
CP_30	AUTHENTIC	EVOO	0.99903
CP_31	AUTHENTIC	EVOO	0.95723
CP_32	AUTHENTIC	EVOO	0.98761
DEO3	ADULTERATED	OTHER	0.99879
DEO_DEA_2	ADULTERATED	OTHER	0.99938
MIX_D	ADULTERATED	OTHER	0.67833

Table S3. Classification independent adulterated (OTHER) and authentic extra virgin olive oil (EVOO) by the support vector machine (SVM) model, built using informative variables from low-level data fusion of gas chromatography coupled with ion mobility spectrometry (GC-IMS), flash gas chromatography electronic nose (FGC-Enose) and multimodal high performance liquid chromatography-high resolution mass spectrometry (HPLC-HRMS) datasets.

SAMPLES	ACTUAL	PREDICTED	AVAREGED PROBABILITY
CP_30	AUTHENTIC	EVOO	0.99753
CP_31	AUTHENTIC	EVOO	0.97704
CP_32	AUTHENTIC	EVOO	0.99371
DEO3	ADULTERATED	OTHER	0.99936
DEO_DEA_2	ADULTERATED	OTHER	0.99969
MIX_D	ADULTERATED	OTHER	0.65853

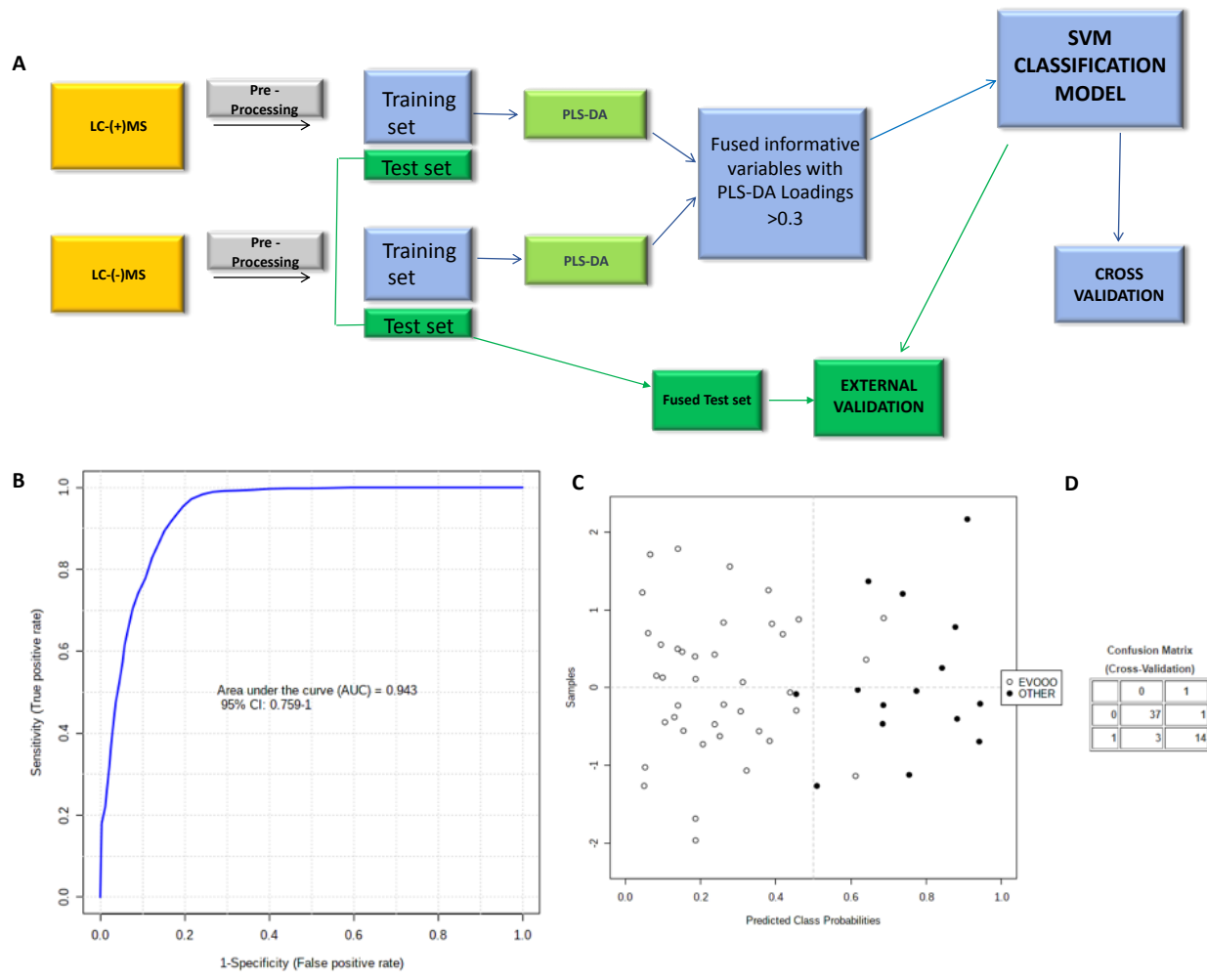


Figure S1. Mid-level data fusion of multimodal high performance liquid chromatography-high resolution mass spectrometry (HPLC-HRMS) datasets and multivariate statistical analysis, aimed at the classification of extra virgin olive oil (EVOO) A) Flow-chart of the mid level data fusion of multimodal HPLC-HRMS with extraction of the most informative variables by Partial least squares-discriminant analysis (PLS-DA) of the single datasets and built-in support vector machine (SVM). B) Receiver operating characteristic (ROC) the performance of a classification model in cross-validation on training set. C) The predictions of SVM model in the cross-validation with D) the resulting confusion matrix.

Table S4. Accuracy Sensitivity specificity obtained by mid-level data fusion of GC-IMS and FGC-Enose and (+/-)HPLC-HRMS datasets and built-in SVM model.

Merged technique	Sensitivity on training set	Specificity on training set	Accuracy on training set	AUC of the ROC	Samples correctly classified in validation on test set	Probability of predictions in validation on test set
HPLC-(+/-)HRMS	0.93	0.93	0.93	0.94	6/6	0.93
GC-IMS FGC-Enose	0.86	0.98	0.95	0.90	6/6	0.88
GC-IMS FGC-Enose HPLC-(+/-)HRMS	0.93	0.98	0.96	0.88	6/6	0.94

Table S5. Classification independent adulterated (OTHER) and authentic extra virgin olive oil (EVOO) by the support vector machine (SVM) model, built using informative variables from mid-level data fusion of multimodal high performance liquid chromatography-high resolution mass spectrometry (HPLC-HRMS) datasets.

SAMPLES	ACTUAL	PREDICTED	AVAREGED PROBABILITY
CP_30	AUTHENTIC	EVOO	0.99749
CP_31	AUTHENTIC	EVOO	0.97675
CP_32	AUTHENTIC	EVOO	0.99858
DEO3	ADULTERATED	OTHER	0.73324
DEO_DEA_2	ADULTERATED	OTHER	0.96799
MIX_D	ADULTERATED	OTHER	0.94743

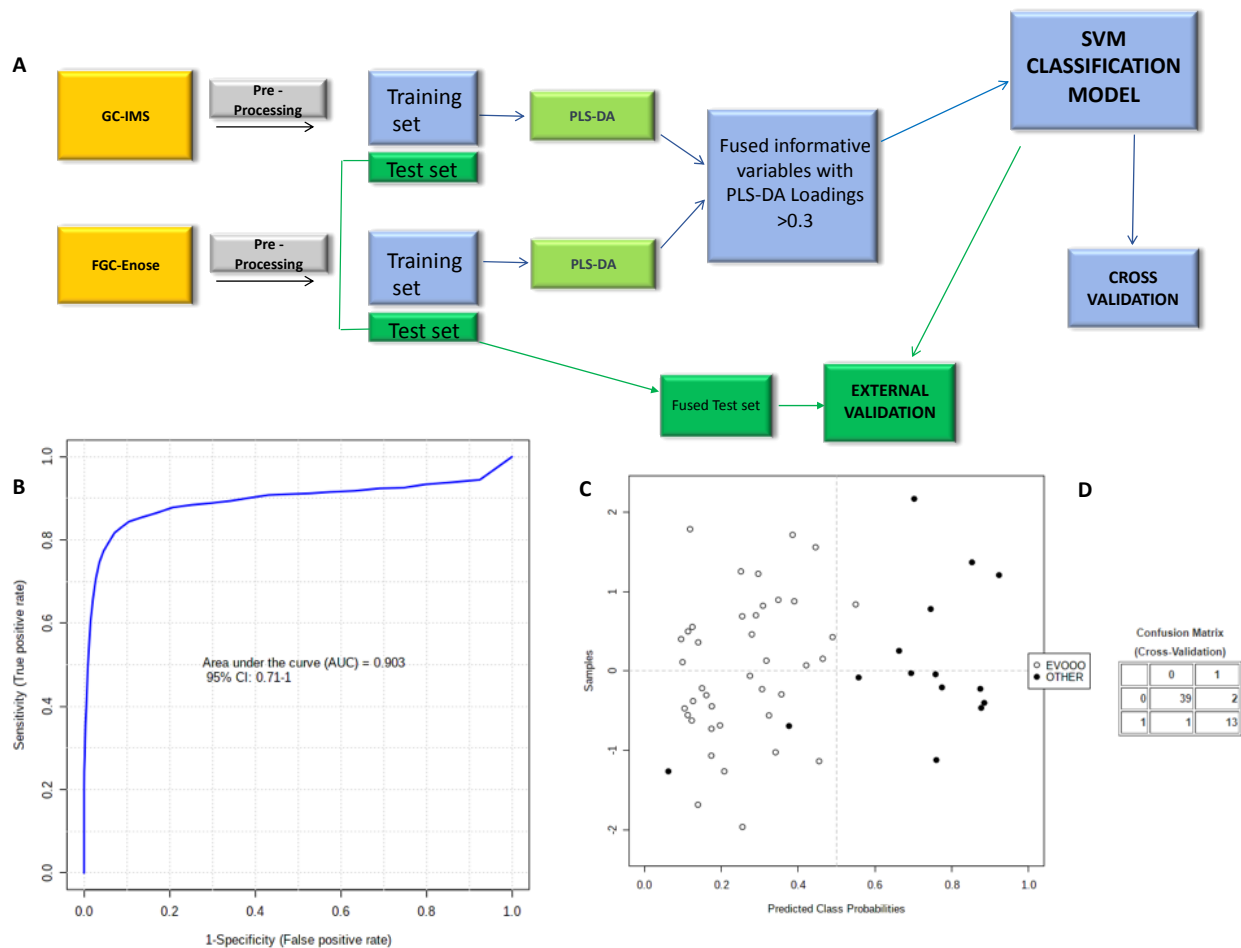


Figure S2. Mid-level data fusion of gas chromatography coupled with ion mobility spectrometry (GC-IMS) and flash gas chromatography electronic nose (FGC-Enose) datasets and multivariate statistical analysis aimed at the classification of extra virgin olive oil (EVOO). A) Flow-chart of the mid-level data fusion of GC-IMS and FGC-Enose with extraction of the most informative variables by Partial least squares-discriminant analysis (PLS-DA) of the single datasets and built-in support vector machine (SVM). B) Receiver operating characteristic (ROC) the performance of a classification model in cross-validation on training set. C) The predictions of SVM model in the cross-validation with D) the resulting confusion matrix.

Table S6. Classification independent adulterated (OTHER) and authentic extra virgin olive oil (EVOO) by the support vector machine (SVM) model, built using informative variables from mid-level data fusion of gas chromatography coupled with ion mobility spectrometry (GC-IMS) and flash gas chromatography electronic nose (FGC-Enose) datasets.

SAMPLES	ACTUAL	PREDICTED	AVAREGED PROBABILITY
CP_30	AUTHENTIC	EVOO	0.9968
CP_31	AUTHENTIC	EVOO	0.87557
CP_32	AUTHENTIC	EVOO	0.89127
DEO3	ADULTERATED	OTHER	0.99869
DEO_DEA_2	ADULTERATED	OTHER	0.95555
MIX_D	ADULTERATED	OTHER	0.60246

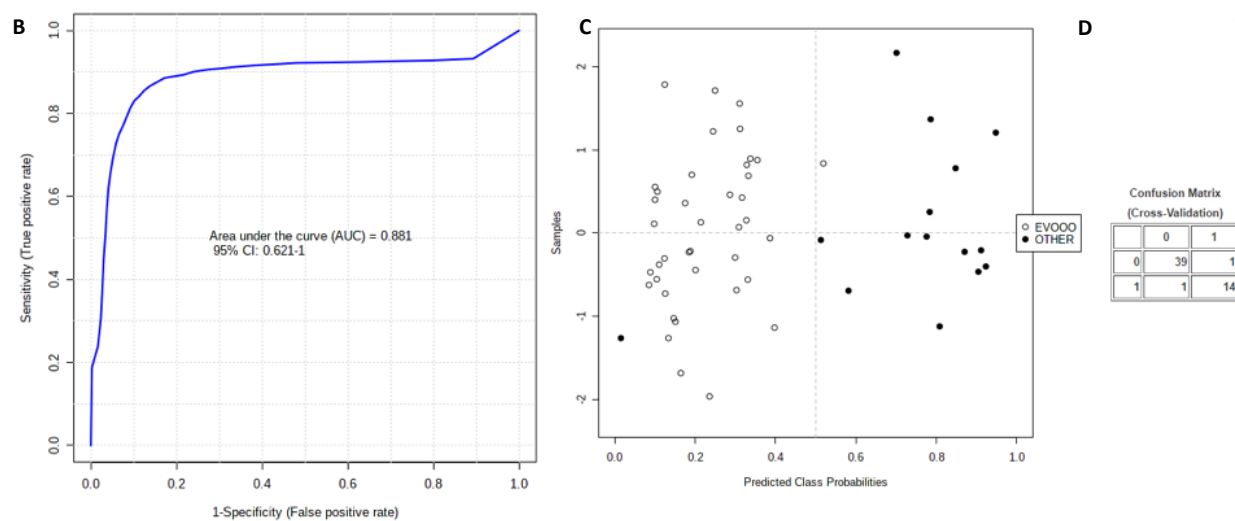
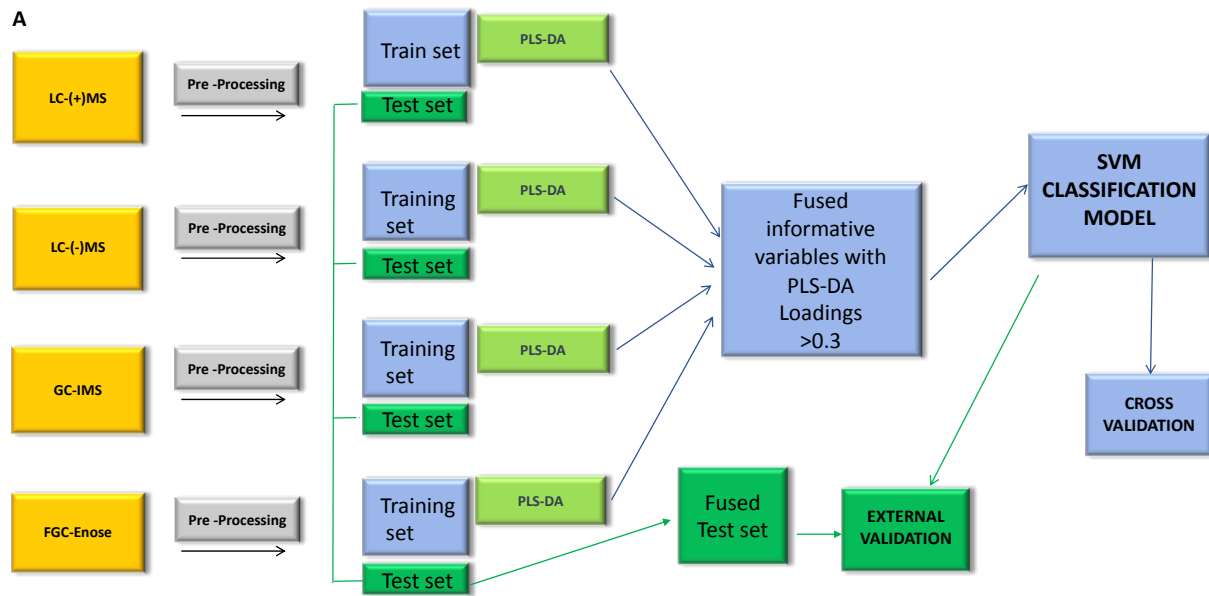


Figure S3. Mid-level data fusion of multimodal high performance liquid chromatography-high resolution mass spectrometry (HPLC-HRMS), gas chromatography coupled with ion mobility spectrometry (GC-IMS) and flash gas chromatography electronic nose (FGC-Enose) datasets coupled to multivariate statistical analysis aimed at the classification of extra virgin olive oil (EVOO). A) Flow-chart of the mid-level data fusion of the four datasets with extraction of the most informative variables by Partial least squares-discriminant analysis (PLS-DA) from each datasets and built-in support vector machine (SVM). B) Receiver operating characteristic (ROC) the performance of a classification model in cross-validation on training set. C) The predictions of SVM model in the cross-validation with D) the resulting confusion matrix.

Table S7. Classification independent adulterated (OTHER) and authentic extra virgin olive oil (EVOO) by the support vector machine (SVM) model, built using informative variables from low-level data fusion of gas chromatography coupled with ion mobility spectrometry (GC-IMS), flash gas chromatography electronic nose (FGC-Enose) and multimodal high performance liquid chromatography-high resolution mass spectrometry (HPLC-HRMS) datasets.

SAMPLE	ACTUAL	PREDICTED	AVAREGED PROBABILITY
CP_30	AUTHENTIC	EVOO	0.985
CP_31	AUTHENTIC	EVOO	0.9778
CP_32	AUTHENTIC	EVOO	0.96445
DEO3	ADULTERATED	OTHER	0.99454
DEO_DEA_2	ADULTERATED	OTHER	0.98428
MIX_D	ADULTERATED	OTHER	0.74301

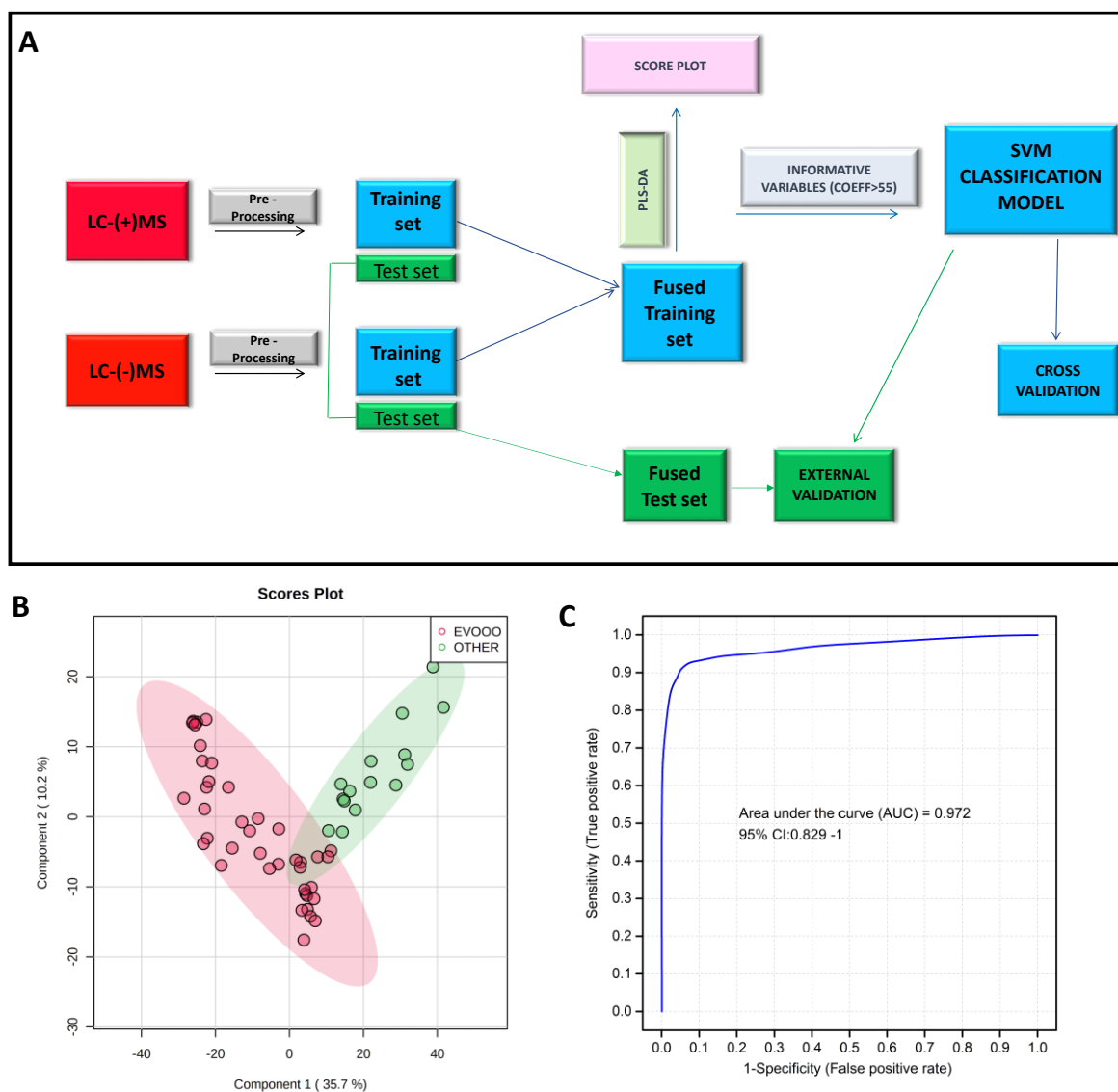


Figure 1. Flow chart of the low-level data fusion and multivariate statistical analysis of multimodality high pressure liquid chromatography-high resolution mass spectrometry (LC(+/-) MS) datasets. A) The flow chart showing the combination of LC(+/-)MS datasets after low-level data fusion. B) PLS-DA score plot that allowed visualization of the discrimination of the two groups in the study. C) The prediction power of the SVM model was estimated based on the area under the curve (AUC) of the receiver operating characteristic (ROC) curve.

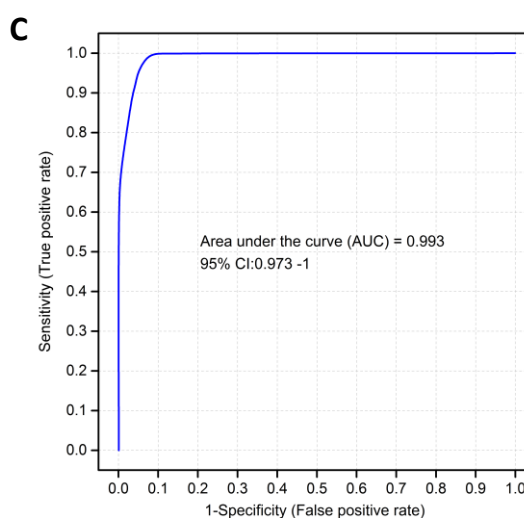
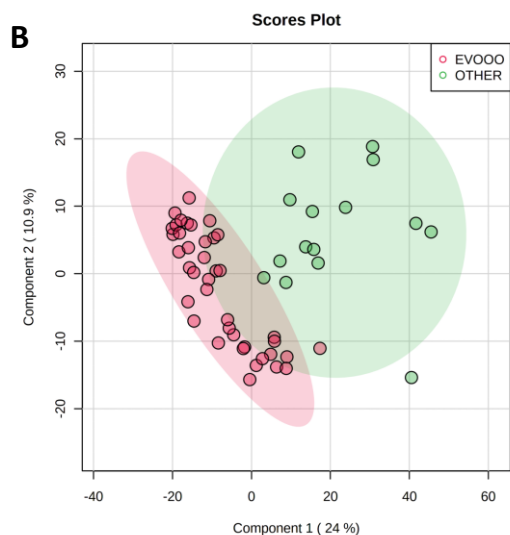
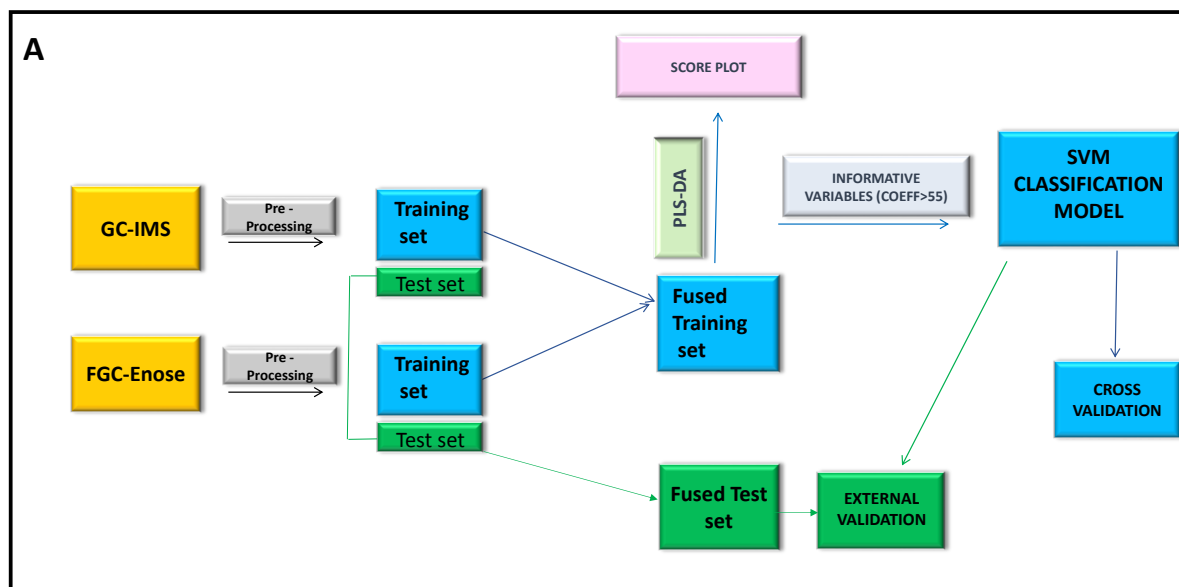


Figure 2. Flow chart of the low-level data fusion and multivariate statistical analysis of gas-chromatography ion mobility spectrometry (GC-IMS) and flash gas-chromatography electronic nose (FGC-Enose) datasets. **A**) The flow chart showing the combination of GC-IMS and FGC-Enose datasets after low-level data fusion. **B**) PLS-DA score plot that allowed visualization of the discrimination of the two groups in the study. **C**) The prediction power of the SVM model was estimated based on the area under the curve (AUC) of the receiver operating characteristic (ROC) curve.

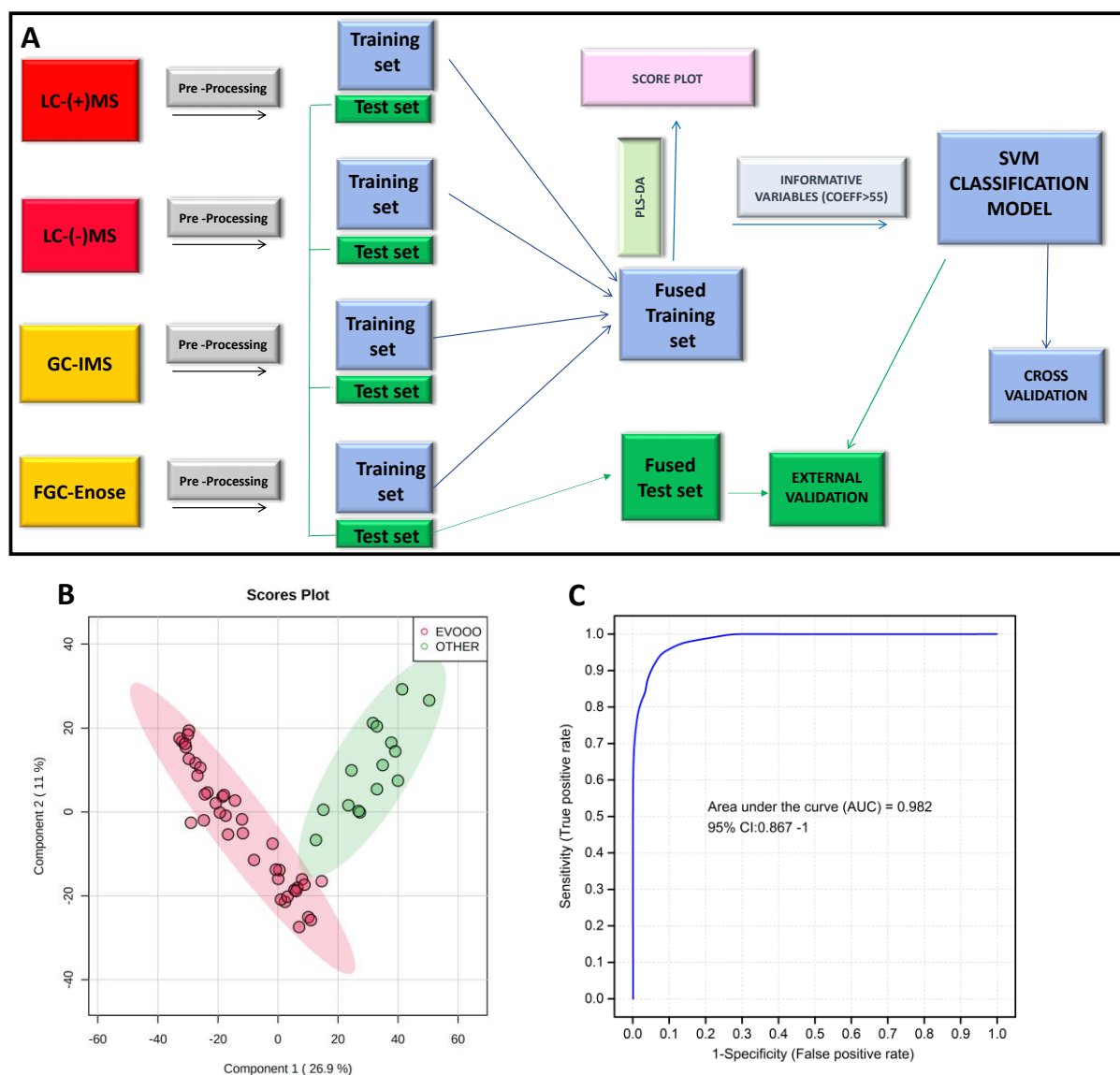


Figure 3. Flow chart of the low-level data fusion and multivariate statistical analysis of multimodality high pressure liquid chromatography-high resolution mass spectrometry (LC-(+/-) MS), gas-chromatography ion mobility spectrometry (GC-IMS) and flash gas-chromatography electronic nose (FGC-Enose) datasets. A) The flow chart showing the combination of the four datasets after low-level data fusion. B) PLS-DA score plot that allowed visualization of the discrimination of the two groups in the study. C) The prediction power of the SVM model was estimated based on the area under the curve (AUC) of the receiver operating characteristic (ROC) curve.

Table 1. Statistical figures of merit of Support Vector Machine (SVM) models obtained in cross-validation on training set (sensitivity, specificity and accuracy) after combining the three analytical approaches by low-level data fusion. Number of samples correctly classified and probability of predictions during in validation on test set are also reported.

Merged technique	Sensitivity on training set	Specificity on training set	Accuracy on training set	AUC of the ROC	Samples correctly classified in validation of the test set	Average probability of the predictions on test set
LC-(+/-)MS	0.93	0.95	0.94	0.97	6/6	0.96
GC-IMS FGC-Enose	0.93	0.97	0.96	0.99	6/6	0.93
GC-IMS FGC-Enose LC-(+/-)MS	0.93	0.96	0.96	0.98	6/6	0.94