

Accepted Manuscript

Cultivar classification of Apulian olive oils: use of artificial neural networks for comparing NMR, NIR and merceological data

Giulio Binetti, Laura Del Coco, Rosa Ragone, Samanta Zelasco, Enzo Perri, Cinzia Montemurro, Raffaele Valentini, David Naso, Francesco Paolo Fanizzi, Francesco Paolo Schena

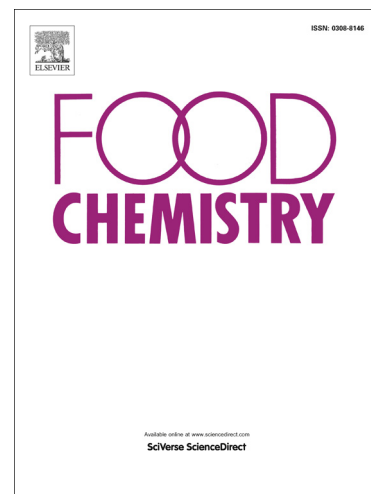
PII: S0308-8146(16)31424-8
DOI: <http://dx.doi.org/10.1016/j.foodchem.2016.09.041>
Reference: FOCH 19821

To appear in: *Food Chemistry*

Received Date: 30 May 2016
Revised Date: 2 September 2016
Accepted Date: 6 September 2016

Please cite this article as: Binetti, G., Coco, L.D., Ragone, R., Zelasco, S., Perri, E., Montemurro, C., Valentini, R., Naso, D., Fanizzi, F.P., Schena, F.P., Cultivar classification of Apulian olive oils: use of artificial neural networks for comparing NMR, NIR and merceological data, *Food Chemistry* (2016), doi: <http://dx.doi.org/10.1016/j.foodchem.2016.09.041>

This is a PDF file of an unedited manuscript that has been accepted for publication. As a service to our customers we are providing this early version of the manuscript. The manuscript will undergo copyediting, typesetting, and review of the resulting proof before it is published in its final form. Please note that during the production process errors may be discovered which could affect the content, and all legal disclaimers that apply to the journal pertain.



1 **Cultivar classification of Apulian olive oils: use of artificial neural networks for**
2 **comparing NMR, NIR and merceological data**

3 **Giulio Binetti^{a,‡}, Laura Del Coco^{b,‡}, Rosa Ragone^{c,‡}, Samanta Zelasco^d, Enzo Perri^d, Cinzia**
4 **Montemurro^e, Raffaele Valentini^f, David Naso^a, Francesco Paolo Fanizzi^{b*}, Francesco Paolo**
5 **Schena^{c*}**

6 ^a Dipartimento di Ingegneria Elettrica e dell'Informazione, Politecnico di Bari, Via E. Orabona, 4,
7 70125 Bari, Italia; e-mail: giulio.binetti@poliba.it, david.naso@poliba.it

8 ^b Dipartimento di Tecnologie Biologiche ed Ambientali, Università del Salento, Prov.le Lecce-
9 Monteroni, 73100 Lecce, Italia; e-mail: laura.delcoco@unisalento.it, fp.fanizzi@unisalento.it

10 ^c Consorzio C.A.R.S.O., Università di Bari, Strada Provinciale Casamassima Km 3, 70010
11 Valenzano (Bari), Italia; e-mail: rosaragone@yahoo.it, paolo.schena@uniba.it

12 ^d Consiglio per la Ricerca in Agricoltura e l'Analisi dell'Economia Agraria-Centro di ricerca per
13 l'Olivicoltura e l'Industria Olearia, Contrada Li Rocchi, 87036 Rende (Cosenza), Italia; e-mail:
14 samanta.zelasco@crea.gov.it, enzo.perri@crea.gov.it

15 ^e Dipartimento di Biologia e Chimica Agro-Forestale ed Ambientale, Sezione di Genetica e
16 Miglioramento, Università di Bari, via Amendola 165/a, 70126 Bari, Italia; e-mail:
17 cinzia.montemurro@uniba.it

18 ^f Oliveti Terra di Bari O.P. Olivicoli Soc. Coop. Agricola, 6/A, Via Brigata 6/A, 70124 Bari, Italia;
19 e-mail: agrivalent@libero.it

20
21 [‡] These authors have contributed equally to this work.

22 * Corresponding authors: G. Binetti, F.P. Fanizzi, F.P. Schena

23

24 **ABSTRACT**

25 The development of an efficient and accurate method for extra-virgin olive oils cultivar and
26 origin authentication is complicated by the broad range of variables (e.g., multiplicity of varieties,
27 pedo-climatic aspects, production and storage conditions) influencing their properties. In this study,
28 artificial neural networks (ANNs) were applied on several analytical datasets, namely standard
29 merceological parameters, near-infra red data and ^1H nuclear magnetic resonance (NMR)
30 fingerprints, obtained on mono-cultivar olive oils of four representative Apulian varieties (Coratina,
31 Ogliarola, Cima di Mola, Peranzana). We analysed 888 samples produced at a laboratory-scale
32 during two crop years from 444 plants, whose variety was genetically ascertained, and on 17
33 industrially produced samples. ANN models based on NMR data showed the highest capability to
34 classify cultivars (in some cases, accuracy > 99%), independently on the olive oil production
35 process and year; hence, the NMR data resulted to be the most informative variables about the
36 cultivars.

37

38 **Keywords:**

39 Artificial neural networks; olive oil; cultivar classification; merceological analysis; near-infra red
40 spectroscopy; nuclear magnetic resonance spectroscopy.

41

42 1. INTRODUCTION

43 Extra-virgin olive oil (EVOO) is considered as one of the best sources of “good” fatty acids and
44 antioxidants, with positive effects on human health [1]. It is a complex food matrix, difficult to be
45 analyzed. Monitoring quality (from harvesting through transformation to storage) and
46 authentication (detection of adulterations, identification of geographical origin and variety) are
47 nowadays the main challenges for olive oil industry and food control laboratories. In an attempt to
48 protect customers and producers against false declarations, international organizations have
49 established the guidelines for olive oil certification, indicating the methods to determine several
50 chemical and physical parameters and reference limits. However, the official procedures often result
51 to be inadequate to screen a large number of samples, time consuming, and insufficient for a quick
52 and detailed examination [2]. As regards the authentication of the cultivar of EVOOs, the
53 development of an efficient and accurate method is complicated by the broad range of variables
54 influencing the olive oil properties: pedo-climatic aspects and process conditions interact with
55 genetic characteristics.

56 The general strategy followed by researchers is to get the metabolic fingerprints of a large
57 amount of oils obtained from many varieties, and then to build up models by means of multivariate
58 statistical analyses (MVA) to predict unknown samples. Among the modern analytical methods,
59 nuclear magnetic resonance (NMR) spectroscopy seems very attractive, because it requires an easy
60 sample preparation and rapidly provides a complete metabolic profile of olive oils, giving
61 information about either the lipid fraction and several minor compounds (sterols, tocopherols,
62 polyphenols, oxidized products, etc.) [3-6]. Near infra-red (NIR) spectroscopy [7,8], gas and high-
63 performance liquid chromatography with mass spectrometry (MS) [9,10], and electronic sensors
64 have been exploited too [11,12].

65 MVA programs handle the large amounts of metabolic data produced by innovative techniques,
66 and extract the main variables that discriminate between the categories under examination. Beside
67 the exploratory methods, such as Principal Component, Hierarchical Cluster and Tree Clustering
68 Analysis (unsupervised), and Linear and Quadratic Discriminant Analysis, Partial Least Square
69 discriminant analysis (PLS-DA) regression (supervised) are some of the most extensively used
70 classification approaches. However, sometimes their usage as a reliable approach to classify
71 cultivars and geographic origins and to unravel adulterations has been negatively criticized [2].

72 Recently, artificial neural networks (ANNs) have been introduced in food analysis [13]. ANNs
73 are a set of mathematical methods, which attempt to mimic the functioning of the human brain [14].
74 They consist of sophisticated non-linear computational tools that are capable of modelling
75 extremely complex functions learning by example: the data structure is automatically learnt from
76 representative data by means of opportunely designed training algorithms. There are some examples
77 of usage of ANNs for olive oil classification according to geographical origin, year of production,
78 merceological category, adulteration, processing, and blending. Generally, in those works ANNs
79 have been built using only one kind of analytical data, such as data obtained through mass
80 spectrometry [15], NIR [16], electronic sensors [11,12], NMR [17,18], or traditional standardized
81 methods [19,20]. The use of ANNs for cultivar classification of olive oils has not been completely
82 explored. Bucci R. et al. have shown that, using chemometrics and ANNs and choosing the
83 chemical indices routinely determined for the oil quality control as descriptors, it was possible to
84 accurately attribute the cultivar of 153 Italian EVOOs obtained from five different varieties [21].
85 Peres A.M. et al. measured ten biometrical parameters of oil samples belonging to six Portuguese
86 cultivars, collected in different groves during four crop years, and created an artificial neural
87 network able of predicting the variety of unknown samples more accurately than using linear
88 discriminant analysis [22].

89 In our study, ANN models were set up using multiple types of information, namely standard
90 merceological parameters, NIR data, and NMR fingerprints, in order to find the most accurate ANN
91 model for cultivar classification.

92 **2. MATERIALS AND METHODS**

93 **2.1 Selection of plants and olive oil extraction**

94 We focused on four representative cultivars of Apulia, namely Coratina, Ogliarola Barese, Cima di
95 Mola, and Peranzana. Forty hundred-fifth plants were selected in 15 areas (30 plants per area)
96 across the Foggia and Bari provinces, in order to cover different pedo-climatic regions. We chose
97 240 plants for Coratina, 60 for Ogliarola Barese, 90 for Cima di Mola, and 60 for Peranzana, based
98 on phenotype characteristics. Every tree was marked with an identification code. Harvest was
99 performed at the optimal olive ripening stage, in different periods depending on cultivar and
100 growing conditions, in two subsequent crop years (2013-2014 and 2014-2015). Considering the
101 very high number of samples, analyses were conducted in the same period for the same cultivar,
102 following each other, compatibly with the time required for each test. The drupes harvested from
103 each plant were milled within 24 hours, by using a mini olive press (Spremoliva C30 milling
104 machine, Toscana Enologica Mori, Tavarnelle Val di Pesa, FI, Italy). About 25-30 kg of olives were
105 processed in each working cycle, lasting approximately 2-3 hours: the machine was cleaned
106 carefully after each cycle. The produced monocultivar olive oils were filtered and stored in sealed
107 dark glass bottles at room temperature prior to analysis.

108 **2.2 Genetic characterization**

109 Molecular characterization was conducted on the 450 olive accessions in order to verify genetic
110 correspondence with reference cultivars (Coratina, Ogliarola barese, Cima di Mola and Peranzana).
111 Reference cultivars were collected from the CREA-OLI olive collection located in Mirto Crosia
112 (CS, Italy). A set of highly discriminant microsatellite molecular markers was used (DCA3, DC5,

113 DCA8, DCA15, DCA18, GAPU71b, GAPU101) for this purpose [23]. Genomic DNA was
114 extracted from young leaves, dried in silica gel for 2-3 days. An amount of 50 mg of dried leaf
115 tissue was ground using TissueLyser II (Qiagen, Hilden, Germany); DNA was extracted using the
116 GenElute™ Plant Genomic DNA Miniprep Kit (Sigma-Aldrich, St. Louis MO, USA), and
117 quantified at NanoDrop 2000 UV-Vis Spectrophotometer (Thermo Scientific, Waltham, MA,
118 USA). PCR amplification was carried out using KAPA3G Plant PCR Kit (Kapa Biosystems,
119 Wilmington, MA, USA). Amplicons were analyzed by using a 16-capillary DNA sequencer (3130
120 Applied Biosystems, Foster City, CA, USA), equipped with the GeneMapper 3.7v software. Two
121 laboratories conducted the molecular analysis (CREA-OLI and DISSPA, University of Bari). To
122 verify the reproducibility of the results, a ring test was conducted on 15 samples randomly chosen.

123 **2.3 Merceological analysis**

124 Twenty-three merceological parameters were measured and subsequently used as inputs for the
125 ANNs: total tocopherols, total phenols, free acidity, peroxide value, UV spectrophotometric indices
126 (K_{232} = UV absorbance at 232 nm, K_{270} = UV absorbance at 270 nm, $\Delta K = K_{232} - K_{270}$), NI R.T.6.9
127 and NI R.T. 7.4 (unidentified peaks), C14:0 (myristic acid), C16:0 (palmitic acid), C16:1is and
128 C16:1c (isomers of palmitoleic acid), C17:0 (eptadecanoic acid), C17:1 (eptadecenoic acid), C18:0
129 (stearic acid), C18:1 (oleic acid), C18:2 (linoleic acid), C18:3 (linolenic acid), C20:0 (arachidic
130 acid), C20:1 (eicosenoic acid), C22:0 (behenic acid), C24:0 (lignoceric acid). All the measurements
131 were carried out according to the official methods of the European Regulation/Commission
132 Regulation EEC no. 2568/91 and its subsequent modifications (EC Reg. 2568/1991) [24].
133 Moreover, the quantification of phenols and tocopherols was performed according to the method
134 described by COI [25].

135 **2.4 Near Infra-Red (NIR) Spectrometry**

136 The Near Infra-Red (NIR) analysis was conducted using the XDS analyser instrument (Foss,
137 Analytical A/S, Denmark) equipped with an infrared reading system. The XDS NIR instrument is
138 supported by the RINA (Remote Internet Analysis) software suite. A dispersive grating
139 monochromator permits a highest signal/noise ratio, thus it guarantees an efficient analysis of
140 complex matrices and dispersions. For each olive oil sample, an aliquot was weighted in a 1.5 mL
141 quartz cuvettes for NIR analysis, without any preliminary treatment. Spectra were acquired at
142 constant temperature of 40°C, according to standard instrument procedures for olive oil analysis, in
143 the wavelength range from 700 nm to 2500 nm. Three measurements were performed for each
144 sample, in order to minimize errors due to instrumental fluctuations. Twenty parameters were
145 evaluated: acidity, peroxides, K_{232} , K_{270} , ΔK , palmitic acid, palmitoleic acid, eptadecanoic acid,
146 eptadecenoic acid, stearic acid, oleic acid, linoleic acid, linolenic acid, arachidic acid, eicosenoic
147 acid, polyphenols, tocopherols, methyl esters, ethyl esters, methyl ethyl esters.

148 **2.5 Nuclear Magnetic Resonance (NMR) Spectrometry**

149 The mono-cultivar EVOO samples, belonging to the four most representative varieties of Apulia,
150 collected during two crop years, were analyzed by ^1H NMR spectroscopy. Briefly, each NMR
151 sample was prepared dissolving ~140 mg of olive oil in CDCl_3 and adjusting the mass ratio of olive
152 oil: CDCl_3 to 13.5%:86.5%. Next, 600 μL of the prepared mixture was transferred into a 5-mm
153 NMR tube. ^1H NMR spectra were recorded on a Bruker Avance spectrometer (Bruker, Karlsruhe,
154 Germany), operating at 400.13 MHz, $T = 300\text{ K}$, equipped with a PABBI 5-mm inverse detection
155 probe incorporating a z axis gradient coil. NMR experiments were performed under full automation
156 for the entire process after loading individual samples on a Bruker Automatic Sample Changer
157 (BACS-60), interfaced with the software IconNMR (Bruker). Automated tuning and matching,
158 locking and shimming, and calibration of the 90° hard pulse $P(90^\circ)$ were done for each sample
159 using standard Bruker routines, ATMA, LOCK, TOPSHIM and PULSECAL, to optimize NMR
160 conditions. For each sample, after a 5-min waiting period for temperature equilibration, standard

161 one-dimensional (^1H ZG) NMR experiments were performed. The relaxation delay (RD) and
162 acquisition time (AQ) were set to 4 s and ~ 3.98 s, respectively, resulting in a total recycle time of
163 ~ 7.98 s. FIDs were collected into time domain (TD) = 65536 (64 k) complex data points by setting:
164 spectral width (SW) = 20.5524 ppm (8223.685 Hz), receiver gain (RG) = 4, number of scans (NS) =
165 16. The accumulation of 16 scans was preferred because of some metabolites present in high
166 concentrations [26].

167 The NMR raw data set was pre-processed using Topspin 2.1 and AMIX 3.9.15 (Bruker BioSpin
168 GmbH, Germany). The FIDs were multiplied by an exponential line broadening function (0.3 Hz)
169 before Fourier transformation and automatically phased. Spectra were referenced to the TMS signal
170 at 0.00 ppm, used as an internal standard and obtaining good peak alignment.

171 We focused on the spectral region within 10.00-0.5 ppm, excluding the signal of the residual non-
172 deuterated chloroform and its carbon satellites (7.6-6.9 ppm). This region was reduced in small
173 intervals (buckets) of equal size (0.04 ppm) by applying the rectangular bucketing procedure,
174 obtaining 221 buckets. Total sum normalization was applied to minimize small differences due to
175 total olive oil concentration and/or acquisition conditions among samples. The Pareto scaling
176 method, which is performed by dividing the mean-centered data by the square root of the standard
177 deviation, was then applied to the variables. Each bucket (NMR variable) was ~~labeled~~ labelled with
178 its average value of chemical shift (ppm), and subsequently used as ANN input. Multivariate
179 analysis (Principal Component Analysis, PCA) of NMR bucket-reduced data was carried out using
180 SIMCA 14 software (Umetrics, Umea, Sweden).

181 All chemical reagents for analyses were of analytical grade. Deuterated chloroform (CDCl_3
182 99.8%-d) containing tetramethylsilane TMS (0.03% v/v) was purchased from Armar Chemicals
183 (Döttingen, Switzerland).

184 2.6 Artificial Neural Networks

185 The most common artificial neural network is the Multi-Layer Perceptron (MLP), successfully
 186 used also for classification and pattern recognition [27-29]. It is a universal function approximator
 187 which can solve non-linearly separable problems and learn any arbitrarily complex linear function
 188 with an arbitrary accuracy level [30,31]. In general, a MLP is composed by one input layer with p
 189 inputs, one or more hidden layers with n hidden neurons, and one output layer with q outputs
 190 (Figure 1). The output of the j -th hidden neuron is computed as

$$191 \quad y_j^h = f^h \left(\sum_{k=1}^p w_{k,j}^h x_k + b_j^h \right), \quad j = 1, \dots, n,$$

192 where f^h is the activation function of the hidden neuron, $w_{k,j}^h$ is the weight between the input x_k
 193 and the hidden neuron j , and b_j^h is a bias term for the hidden neurons. Then, considering a MLP
 194 with one hidden layer, the i -th output is computed as

$$195 \quad y_i = f^o \left(\sum_{j=1}^n w_{j,i}^o y_j^h + b_i^o \right), \quad i = 1, \dots, q,$$

196 where f^o is the activation function of the output neuron, $w_{j,i}^o$ is the weight between the hidden
 197 neuron j and the output neuron i , and b_i^o is a bias term for the output neuron.

198 In this study, the Multi-Layer Perceptron model was used for olive oil cultivar classification, and
 199 was configured as a pattern recognition network with sigmoidal and softmax activation functions
 200 for hidden and output neurons, respectively [28]. The training process was performed using the
 201 Scaled Conjugate Gradient algorithm to minimize the cross-entropy cost function. The dataset was
 202 normalized with a zero mean and unit standard deviation transformation. Then, the dataset was
 203 divided into training, validation, and test sets: the training set was used for the learning process, the
 204 validation set was used during the learning process to avoid overfitting by adopting the early stop
 205 strategy [29], and the test set was used to properly evaluate the classifier performance on an

206 independent data set. In addition, a 5-fold cross validation was performed to limit the bias of the
207 model performance associated with a random sampling of the training data. This means that, 5 folds
208 were randomly created with roughly the same cultivar proportion as in the initial dataset; then, 5
209 different MLP models were trained by using different folds for training, validation and test. We did
210 many trials, varying the number of hidden layers from 1 to 2, the number of hidden neurons from 3
211 to 20, and re-training each network 1000 times with different random initial weights. Finally, we
212 chose the best model based on the accuracy, i.e. a global criterion defined as the percentage of
213 correct prediction in the dataset. Considering the confusion matrix (i.e., an $m \times m$ table that
214 compares the classifier outputs with the actual values in the dataset), the accuracy was defined as:

$$215 \quad a = \frac{TP + TN}{(TP + TN + FP + FN)},$$

216 where TP, TN, FP and FN are the true positive, true negative, false positive and false negative
217 values, respectively. Then, since 5 MLP models were trained on different folds using the k-folds
218 cross validation, the classifier performance was defined as the average value of the accuracies for
219 the 5 trained MLP models.

220 The LASSO (*Least Absolute Shrinkage and Selection Operator*) algorithm [32,33], was used to
221 identify the most informative predictors for the olive oil cultivars among all the predictors. Indeed,
222 there were cases of data correlation/redundancy; for example, the NMR signals at 1.30 ppm, 2.82
223 ppm, and 0.9 ppm refer, respectively, to the methylenic, bis-allylic and methyl protons of the same
224 molecule (linolenic acid). Thus, the LASSO algorithm was used to reduce the number of predictors
225 in the regression models by selecting the most informative predictors, and to produce shrinkage
226 estimates with potentially lower predictive errors than ordinary least squares.

227 In addition to the LASSO algorithm, two heuristic approaches based on the standard deviation of
228 each predictor were also considered, to further investigate the importance of the predictors for the

229 olive oil cultivar. The rationale was that, an input with a smaller range of variability may be
230 expected to provide a lower influence on the ANN output compared to an input with a greater range
231 of variability. Thus, in the first heuristic approach we chose the predictors with greater standard
232 deviations (I), while in the second one we chose the predictors with greater standard deviations
233 normalized by their respective average values (II).

234 Different ANNs were trained using different combinations of variables: the full dataset, the
235 subsets of variables provided by LASSO, and the two subsets of variables selected according to
236 standard deviation and normalized standard deviation criteria (Table 1), in order to find the best
237 model for cultivar classification performance and to compare the prediction capability of the
238 different predictors among merceological, NIR, and NMR data. The MATLAB software was used
239 for LASSO analysis and neural networks training and validation.

240 **3. RESULTS**

241 **3.1 Genetic data**

242 The set of microsatellites used for the molecular characterization (DCA3, DC5, DCA8, DCA15,
243 DCA18, GAPU71b, GAPU101) discriminated efficiently all the analyzed accessions, showing a
244 unique molecular profile corresponding to the four reference cultivars (Table S1). Molecular data
245 were highly comparable between CREA-OLI and DISSPA laboratories, except at the DCA9
246 (182/194 vs 172/186) and DCA18 (175/177 vs 177/179) loci for ‘Peranzana’ and ‘Ogliarola’
247 accessions, respectively, where allele assignments were different. Ring test conducted on 15
248 accessions randomly chosen confirmed the reproducibility of the analysis. Genetic analysis showed
249 a discrepancy for only 6 out of 450 accessions, which were excluded from the study. Among them,
250 four accessions were riconducibile to Apulian varieties (Simona, Pasola and Cima di Melfi), one was
251 different from the reference profile at two loci (DCA9, DCA18), and the last one did not correspond

252 to any known cultivar (Table S1). In conclusion, 444 olive oil samples of each crop year were used
253 for further analyses.

254 3.2 Merceological data

255 Average and standard deviation of all merceological parameters are reported in Table S2.
256 Coratina samples presented the highest content of phenols, higher than the minimum level
257 ($>250\text{mg}\cdot\text{Kg}^{-1}$) necessary to boast healthy claims, as established by Reg. CEE 2568/91 and its
258 subsequent modifications [24]. Another interesting aspect of Coratina oils concerned the content of
259 eicosenoic acid, that was higher than the maximum value fixed at 40% by regulations in 22% of
260 samples.

261 3.3 NIR data

262 Results obtained through NIR spectroscopy are summarized in Table S3. NIR profiles, recorded
263 in both crop years, showed that Coratina samples were enough different from samples belonging to
264 the other three cultivars. Coratina oils presented the highest content of oleic and eicosenoic acids
265 and the lowest content of all the other fatty acids and total tocopherols. The average values of
266 acidity, peroxides, K_{232} and K_{270} , and ΔK of all samples, were compatible with the “extra-virgin”
267 definition, based on the limits fixed by the Commission Regulation EC. No 1989/2003 [34].

268 3.4 NMR data

269 ^1H NMR profiling of olive oils dissolved into deuterated chloroform is a well-established
270 technique in metabolomics and for classification of olive oil cultivars [6,35]. NMR spectroscopy
271 combines targeted and non-targeted analysis within one single measurement and provides a
272 remarkable level of reproducibility of the data. Moreover, due to a highly reproducible and very
273 detailed fingerprinting, it is possible to differentiate samples even if only small changes occur.
274 Determination of both fatty acid profile and unsaponifiable fraction is usually obtained from proton

275 ^1H NMR spectrum according to literature data [26]. The olefinic protons $-\text{CH}=\text{CH}-$ of all
276 unsaturated fatty acids were assigned at 5.4-5.3 ppm, the proton signal at 5.14 ppm was assigned to
277 $>\text{CHOCOR}$ of sn 1,2 DGs; bis-allylic protons ($=\text{CHCH}_2\text{CH}=\text{}$) of linolenic and linoleic acids were
278 assigned at 2.85-2.70 ppm, the methylene (CH_2) protons of at 1.2 ppm, and the terminal methyl
279 group protons of all saturated and unsaturated chains at 1.0-0.8 ppm. Signals in the range between
280 4.75 and 4.55 ppm referred to different terpenes, while protons of aldehydes and phenolic
281 compounds resonate in the range 9.7-9.1 and 7.0-5.6 ppm, respectively.

282 Interestingly, among the 221 buckets constituting the reduced NMR spectrum of each sample, only
283 24 buckets were selected on the basis of their ability to discriminate between cultivars, and then
284 used for training ANNs (see Section 3.5). A preliminary work on a MVA analysis related to the
285 complete NMR data set (221 buckets for 900 samples) has been already reported [36]. On the first
286 attempt, in order to reveal a general data grouping of all the samples, an unsupervised PCA analysis
287 was applied to the whole data (^1H NMR-bucket-reduced spectra). Visual inspection of three
288 dimensional PCA scoreplot, reported in Figure 2, showed a certain degree of separation in particular
289 for the Coratina samples, while a certain degree of overlap was observed among the three remaining
290 classes, Cima Di Mola, Ogliarola and Peranzana.

291

292 3.5 ANN data

293 We exploited the ANN methodology for cultivar classification of 888 mono-cultivar olive oil
294 samples obtained from the two crop years. Different sets of data, chosen according to the four
295 criteria described in Section 2.6, were used to train ANNs, in order to find the most accurate ANN
296 model, and consequently the most informative analytical technique (Table 1).

297 We recorded globally 43 “traditional” variables for merceological and NIR analyses and 221
298 “innovative” variables for NMR analysis. ANN models were trained using data gathered from the

299 two crop years, and were validated on an independent test set composed by 176 samples not used
300 during the training process. They showed similar optimal performances independently on the data
301 type. In fact, the values of classification accuracy (i.e., the prediction rate on the independent test
302 set) of the ANNs trained with all 43 merceological and NIR variables and the ANNs trained with all
303 221 NMR variables were 98.9% and 99.0%, respectively (Table 2). Moreover, ANN models
304 showed similar complexity: the ANNs trained with “traditional” variables and ANNs trained with
305 “innovative” variables were composed by 2 layers with 20 and 17 hidden neurons, respectively
306 (Table 2).

307 The trained ANN models were also used to independently classify the complete dataset of the
308 two crop years. The classification accuracy remained high independently on the crop year. In fact,
309 the values of classification accuracy of the ANNs trained with “traditional” variables and ANNs
310 trained with “innovative” variables were 99.5% and 99.4% for the first crop (2013-2014) and 99.6%
311 and 99.7% for the second crop (2014-2015), respectively (Table 3).

312 The application of the LASSO algorithm allowed us to select a subset of 29 variables for
313 merceological and NIR analyses (Table 1). The ANNs trained with this subset performed as well as
314 the ANNs trained with all 43 traditional variables, having values of accuracy about 99% (Table 2,
315 and 3). This means that, the findings removed by the LASSO algorithm were not informative for the
316 purpose of cultivar classification of olive oils. Analogously, ANNs created using the subset of 24
317 NMR variables selected by LASSO presented a classification capability similar to that of ANNs
318 created using all 221 NMR variables (Table 2, and 3). Hence, only few regions of the $^1\text{H-NMR}$
319 spectrum obtained from an olive oil sample were very informative about its cultivar. Moreover,
320 reducing the number of predictors, the complexity of the ANN models was also reduced, from 2
321 layers with 20 hidden neurons to 2 layers with 13 hidden neurons for traditional variables and from
322 2 layers with 17 hidden neurons to 1 layer with 16 hidden neurons (Table 2).

323 Since inputs with different variability may have different influence on the ANN outputs, we
324 chose other two subsets of variables to be used for ANN training based on two criteria: (I) greater
325 standard deviations or (II) greater standard deviations normalized by their respective average
326 values. According to (I), we selected 3 traditional and 5 innovative variables, whereas according to
327 (II) we selected 4 traditional and 10 innovative variables (Table 1). In both cases, ANN models
328 trained with the subsets of merceological and NIR data presented significantly lower accuracy in
329 cultivar classification than models trained with NMR data, i.e. 80.5% and 91.6% using (I), and
330 60.8% and 97.5% using (II), respectively (Table 2). Similar trends were also observed when using
331 the full dataset of the two independent crop years: performances ranged from 57.5% to 81.6% for
332 “traditional” data and from 92.4% to 99.5% for “innovative” variables (Table 3).

333 3.6 Validation on industrially produced olive oils

334 The ANNs described above, that were obtained using data from monocultivar olive oils
335 produced at a laboratory-scale, were tested on a group of 17 monocultivar olive oils industrially
336 produced, belonging to the four cultivar considered. For the 17 testing samples, we recorded all the
337 same merceological, NIR, and NMR variables, as done with the previous collection of 888 olive
338 oils. The aim of this step was to verify if the trained ANNs could work with commercialized (large-
339 scale produced) olive oils as well as with oils obtained by mini olive press, despite of differences in
340 operating conditions (volume, instrumentation, storage, etc...). Table 4 shows the values of
341 accuracy in assigning the cultivar of the testing samples relative to the trained ANN models.

342 The ANN model built using the subset of 24 NMR variables, found by application of LASSO
343 algorithm, correctly classified 15 out of 17 testing samples, with the highest accuracy (88.2%)
344 among the compared ANN models.

345 4. DISCUSSION

346 In our study, for the first time, ANN models were set up using multiple types of information,
347 namely standard merceological parameters, NIR data profiles, and NMR fingerprints, in order to
348 find the most accurate ANN model for cultivar classification. For training and validation of the
349 ANNs, we used data from 888 mono-cultivar olive oil samples produced at a laboratory scale,
350 belonging to four varieties (ascertained by genetic analysis), and collected during two crop years.
351 The ANNs were also tested on a smaller dataset composed by 17 samples industrially produced, to
352 verify if the trained ANNs could work with these as well, despite of differences in operating
353 conditions (volume, instrumentation, storage, etc...).

354 Overall, the ANNs seems to be an excellent approach to classify olive oils according to cultivar;
355 in particular, the ANN models based on NMR data shows the best performance, when considering
356 the samples produced either at a laboratory scale and at a large scale.

357 In fact, the high values of accuracy reported for ANN models built with NMR data (in some
358 cases > 99%) suggest that, NMR spectroscopy supplied a quantity of information, useful for
359 cultivar classification, similar or higher than merceological analysis and NIR spectroscopy together.
360 Moreover, the method of olive oil milling has a weak influence on performances of ANNs trained
361 with NMR data (Table 4). Most information are contained in narrow regions of the entire NMR
362 spectrum, namely the 24 variables extracted by LASSO. Consequently, they are the most suitable
363 attributes for classification of samples according to the cultivars examined here. In details, these
364 spectral regions comprise the peaks of phenolic compounds (NMR resonances at 6.86 6.74, 6.58,
365 6.30 ppm), aldehydes (9.06, 7.98, 7.94 ppm), acyl groups of all TGs (triglycerides, 5.34 ppm), acyl
366 groups of *sn* 1,2 DGs (diglycerides, 5.14, 4.94, 3.70 ppm), fatty acids such as linoleic and linolenic
367 acids (2.82, 2.78, 1.30 ppm), and cycloartenol (0.58 ppm). Our results confirm previous literature
368 findings. Indeed, the content of different polyphenols and aldehydes has been largely associated
369 with the organoleptic properties typical of a cultivar [2]. Differences in the level of saturated and
370 unsaturated fatty acids have been observed among cultivars [37]. Cycloartenol, that is a triterpenoid
371 of the sterol class, has already been found to be affected by cultivar, as well as triterpene alcohol

372 composition and total triterpene alcohol content [38]. On the other hand, the level of *sn* 1,2 DGs or,
373 more correctly, the *sn* 1,2 DGs / *sn* 1,3 DGs ratio, has been correlated with the degree of lipid
374 degradation due to the activity of hydrolytic enzymes, that prevalently depends on oil production
375 and storage processes [39].

376 In conclusion, we assess that the combination of NMR fingerprinting with ANN modelling
377 could provide an effective, robust, and rapid method for classifying olive oil cultivars.

378 CONCLUSIONS

379 The application of ANNs evidence that NMR data showed the highest capability to classify
380 cultivars ~~are the most informative variables~~. ANN modelling of NMR data provide an accurate and
381 efficient approach for cultivar classification of olive oils and prediction of unknown samples,
382 independently from methods of milling and year of production. In addition, it is known that NMR
383 technique is advantageous in terms of required time and costs (especially compared to traditional
384 merceological methods). As well, the ANN approach presents its own advantages compared to other
385 chemometric techniques, e.g. it is a non-linear method, fitting better to the data; no particular
386 manipulation of raw data is needed; the MLP design is relatively simple, presenting connections in
387 parallel and sequence between neurons; it learns and does not need reprogramming, thus
388 implementation is not difficult; it can handles data of different origins more easily than using other
389 approaches.

390 However, the limit of the ANN models is that they are applicable exclusively to monovarietal
391 samples belonging to the four varieties (Coratina, Ogliarola, Cima di Mola, Peranzana) used for
392 their training.

393 Acknowledgments

394 This study was supported by a grant from the Italian Ministry of University and Research-MIUR
395 (PON01_01958 PIVOLIO). We gratefully acknowledge Oliveti Terra di Bari (Bari, Italy) and
396 Olearia Basile (Andria, BT, Italy) farms for sample supply.

397 **References**

- 398 [1] G. Sindona, A Marker of Quality of Olive Oils: The Expression of Oleuropein, Olives and
399 Olive Oil in Health and Disease Prevention (2010) 95–100.
- 400 [2] P. Dais, E. Hatzakis, Quality assessment and authentication of virgin olive oil by NMR
401 spectroscopy: A critical review, *Anal. Chim. Acta.* 765 (2013) 1–27.
- 402 [3] L. Del Coco, S.A. De Pascali, V. Iacovelli, G. Cesari, F.P. Schena, F.P. Fanizzi, Following
403 the olive oil production chain: 1D and 2D NMR study of olive paste, pomace, and oil, *Eur. J. Lipid*
404 *Sci. Technol.* 116 (2014) 1513–1521.
- 405 [4] C. Lucas-Torres, A. Pérez, B. Cabañas, A. Moreno, Study by ^{31}P NMR spectroscopy of the
406 triacylglycerol degradation processes in olive oil with different heat-transfer mechanisms, *Food*
407 *Chem.* 165 (2014) 21–8.
- 408 [5] R.I.M. Almoselhy, M.H. Allam, M.H. El-Kalyoubi, A.A. El-Sharkawy, ^1H NMR spectral
409 analysis as a new aspect to evaluate the stability of some edible oils, *Ann. Agric. Sci.* 59 (2014)
410 201–206.
- 411 [6] F. Longobardi, A. Ventrella, C. Napoli, E. Humpfer, B. Schütz, H. Schäfer, et al.,
412 Classification of olive oils according to geographical origin by using ^1H NMR fingerprinting
413 combined with multivariate analysis, *Food Chem.* 130 (2012) 177–183.
- 414 [7] M. Casale, N. Sinelli, P. Oliveri, V. Di Egidio, S. Lanteri, Chemometrical strategies for
415 feature selection and data compression applied to NIR and MIR spectra of extra virgin olive oils for

- 416 cultivar identification., *Talanta*. 80 (2010) 1832–7.
- 417 [8] M. Casale, R. Simonetti, Review: Near infrared spectroscopy for analysing olive oils, *J.*
418 *Near Infrared Spectrosc.* 22 (2014) 59–80.
- 419 [9] A. Bajoub, A. Carrasco-Pancorbo, E.A. Ajal, N. Ouazzani, A. Fernández-Gutiérrez,
420 Potential of LC-MS phenolic profiling combined with multivariate analysis as an approach for the
421 determination of the geographical origin of north Moroccan virgin olive oils, *Food Chem.* 166
422 (2015) 292–300.
- 423 [10] P. Agozzino, G. Avellone, D. Bongiorno, L. Ceraulo, S. Indelicato, S. Indelicato, et al.,
424 Determination of the cultivar and aging of Sicilian olive oils using HPLC-MS and linear
425 discriminant analysis, *J. Mass Spectrom.* 45 (2010) 989–995.
- 426 [11] M.S. Cosio, D. Ballabio, S. Benedetti, C. Gigliotti, Geographical origin and authentication
427 of extra virgin olive oils by an electronic nose in combination with artificial neural networks, *Anal.*
428 *Chim. Acta.* 567 (2006) 202–210.
- 429 [12] L.G. Dias, A. Fernandes, A.C. a Veloso, A. a S.C. Machado, J. a. Pereira, A.M. Peres,
430 Single-cultivar extra virgin olive oil classification using a potentiometric electronic tongue, *Food*
431 *Chem.* 160 (2014) 321–329.
- 432 [13] E. Funes, Y. Allouche, G. Beltrán, A. Jiménez, A Review: Artificial Neural Networks as
433 Tool for Control Food Industry Process, *J. Sens. Technol.* 5 (2015) 28–43.
- 434 [14] J. Zupan, J. Gasteiger, Neural networks: A new method for solving chemical problems or
435 just a passing phase?, *Anal. Chim. Acta.* 248 (1991) 1–30.
- 436 [15] R. Goodacre, D.B. Kell, Pyrolysis mass spectrometry and its applications in biotechnology,
437 *Curr. Opin. Biotechnol.* 7 (1996) 20–8.

- 438 [16] F. Marini, A.L. Magrì, R. Bucci, A.D. Magrì, Use of different artificial neural networks to
439 resolve binary blends of monocultivar Italian olive oils, *Anal. Chim. Acta.* 599 (2007) 232–40.
- 440 [17] D.L. García-González, L. Mannina, M. D’Imperio, A.L. Segre, R. Aparicio, Using ^1H and
441 ^{13}C NMR techniques and artificial neural networks to detect the adulteration of olive oil with
442 hazelnut oil, *Eur. Food Res. Technol.* 219 (2004) 545–548.
- 443 [18] S. Rezzi, D.E. Axelson, K. Héberger, F. Reniero, C. Mariani, C. Guillou, Classification of
444 olive oils using high throughput flow ^1H NMR fingerprinting with principal component analysis,
445 linear discriminant analysis and probabilistic neural networks, *Anal. Chim. Acta.* 552 (2005) 13–24.
- 446 [19] F. Marini, Artificial neural networks in foodstuff analyses: Trends and perspectives A
447 review, *Anal. Chim. Acta.* 635 (2009) 121–131.
- 448 [20] S.F. Silva, C.A.R. Anjos, R.N. Cavalcanti, R.M.D.S. Celeghini, Evaluation of extra virgin
449 olive oil stability by artificial neural network, *Food Chem.* 179 (2015) 35–43.
- 450 [21] R. Bucci, A.D. Magrì, A.L. Magrì, D. Marini, F. Marini, Chemical authentication of extra
451 virgin olive oil varieties by supervised chemometric procedures, *J. Agric. Food Chem.* 50 (2002)
452 413–418.
- 453 [22] A.M. Peres, P. Baptista, R. Malheiro, L.G. Dias, A. Bento, J.A. Pereira, Chemometric
454 classification of several olive cultivars from Trás-os-Montes region (northeast of Portugal) using
455 artificial neural networks, *Chemom. Intell. Lab. Syst.* 105 (2011) 65–73.
- 456 [23] V. Alba, C. Montemurro, W. Sabetta, A. Pasqualone, A. Blanco, SSR-based identification
457 key of cultivars of *Olea europaea* L. diffused in Southern-Italy, *Sci. Hortic. (Amsterdam)* 123
458 (2009) 11–16.
- 459 [24] European Commission, Regulation No 2568/91 on the characteristics of olive oil and olive-

- 460 residue oil and on the relevant methods of analysis (1991).
- 461 [25] International Olive Council, COI/T.20/Doc No 29 Determination of biophenols in olive oils
462 by HPLC (2009) 1–8.
- 463 [26] A. Barison, C.W.P. da Silva, F.R. Campos, F. Simonelli, C.A. Lenz, A.G. Ferreira, A simple
464 methodology for the determination of fatty acid composition in edible oils through ^1H NMR
465 spectroscopy, *Magn. Reson. Chem.* 48 (2010) 642–650.
- 466 [27] G.P. Zhang, Neural networks for classification: a survey, *IEEE Trans. Syst. Man Cybern.*
467 Part C Applications, 30 (2000) 451–462.
- 468 [28] C.M. Bishop, *Neural networks for pattern recognition*, 2005.
- 469 [29] C.M. Bishop, *Pattern Recognition and Machine Learning*, 2006.
- 470 [30] K. Hornik, M. Stinchcombe, H. White, Universal approximation of an unknown mapping
471 and its derivatives using multilayer feedforward networks, *Neural Networks*. 3 (1990) 551–560.
- 472 [31] K. Hornik, Approximation capabilities of multilayer feedforward networks, *Neural*
473 *Networks* 4 (1991) 251–257.
- 474 [32] R. Tibshirani, R.S. Society, Regression and shrinkage via the LASSO, *J R Stat Soc, Ser B.*
475 58 (1996) 267–288.
- 476 [33] R. Tibshirani, Regression shrinkage and selection via the LASSO: A retrospective, *J. R.*
477 *Stat. Soc. Ser. B Stat. Methodol.* 73 (2011) 273–282.
- 478 [34] European Commission, Regulation No 1989/2003 amending Regulation No 2568/91 on the
479 characteristics of olive oil and olive-pomace oil and on the relevant methods of analysis (2003).
- 480 [35] S. Piccinonna, R. Ragone, M. Stocchero, L. Del Coco, S.A. De Pascali, F.P. Schena, et al.,

481 Robustness of NMR-based metabolomics to generate comparable data sets for olive oil cultivar
 482 classification. An inter-laboratory study on Apulian olive oils, Food Chem. 199 (2016) 675–683.

483 [36] Fanizzi, F.P., Del Coco, L., Girelli, C.R., De Pascali, S.A. Harvest year effects on Apulian
 484 EVOOs evaluated by ¹H NMR based metabolomics. IV International Conference on Foodomics,
 485 Cesena (FC), Book of abstract (2015) 31–32.

486 [37] R.M. Alonso-Salces, K. Héberger, M. V. Holland, J.M. Moreno-Rojas, C. Mariani, G.
 487 Bellan, et al., Multivariate analysis of NMR fingerprint of the unsaponifiable fraction of virgin olive
 488 oils for authentication purposes, Food Chem. 118 (2010) 956–965.

489 [38] R. Aparicio, G. Luna, Characterisation of monovarietal virgin olive oils, Eur. J. Lipid Sci.
 490 Technol. 104 (2002) 614–627.

491 [39] B. Nieva-echevarría, E. Goicoechea, M.J. Manzanos, M.D. Guillén, A method based on ¹H
 492 NMR spectral data useful to evaluate the hydrolysis level in complex lipid mixtures, FRIN. 66
 493 (2014) 379–387.

494

495 **Table 1**

496 Feature selection results for merceological, NIR and NMR data (L acid: linoleic acid; Ln acid:
 497 linolenic acid; DGs: diglycerides; TGs: triglycerides).

Data	Merceological + NIR	NMR
Full dataset	43 variables (23 + 20 variables)	221 variables
LASSO selection	29 variables merc_tocopherols_tot NIR_acidity merc_phenols_tot NIR_peroxide merc_peroxide NIR_K232 merc_K232 NIR_K270 merc_K270 NIR_palmitoleic_acid merc_C14 NIR_eptadecanoic_acid merc_C16_1IS NIR_eptadecenoic_acid	24 variables NMR_906, 798, 794 (aldehydes) NMR_686, 674, 658, 630 (phenols) NMR_534 (TGs <i>CH=CH</i>) NMR_514, 494 (<i>sn</i> 1,2 DGs > <i>CHOCOR</i>) NMR_466 (terpenes) NMR_390 NMR_370 (<i>sn</i> 1,2 DGs <i>CH₂OH</i>)

	merc_C16_1C merc_C17_0 merc_C17_1 merc_C18_1 merc_C20_0 merc_C24_0	NIR_stearic_acid NIR_oleic_acid NIR_linoleic_acid NIR_linolenic_acid NIR_eicosenoic_acid NIR_tocopherols NIR_methyl_esters NIR_ethyl_esters NIR_methylethyl_esters	NMR_326 NMR_282 (Ln acid =CHCH ₂ CH=) NMR_278 (L acid =CHCH ₂ CH=) NMR_254 NMR_218 NMR_198 (Acyl groups CH ₂ CH=CH) NMR_190 (Acyl groups OCOCH ₂ CH ₂) NMR_158 (Acyl groups OCOCH ₂ CH ₂) NMR_154 (Acyl groups OCOCH ₂ CH ₂) NMR_130 (CH ₂ Ln and L acids) NMR_058 (CH ₂ cycloartenol)
Selection with standard deviation criterion	3 variables: merc_tocopherols_tot merc_phenols_tot NIR_tocopherols		5 variables: NMR_130 (CH ₂ Ln and L acids) NMR_198 (Acyl groups CH ₂ CH=CH) NMR_534 (TGs CH=CH) NMR_278 (L acid =CHCH ₂ CH=) NMR_158 (Acyl groups OCOCH ₂ CH ₂)
Selection with normalized standard deviation criterion	4 variables: merc_C24_0 merc_C17_0 merc_C17_1 NIR_ethyl_esters		10 variables: NMR_058 (CH ₂ cycloartenol) NMR_278 (L acid =CHCH ₂ CH=) NMR_514 (<i>sn</i> DGs 1,2 >CHOCOR) NMR_190 (Acyl groups OCOCH ₂ CH ₂) NMR_370 (<i>sn</i> 1,2 DGs CH ₂ OH) NMR_154 (Acyl groups OCOCH ₂ CH ₂) NMR_674 (phenolic compounds -Ph-H) NMR_658 (phenolic compounds -Ph-H) NMR_282 (Ln acid =CHCH ₂ CH=) NMR_198 (Acyl groups CH ₂ CH=CH)

498

499

500

501 **Table 2**

502 Architecture of the best ANNs and accuracy computed on the independent test set composed by
503 samples of both the two crop years.

Data	Merceological + NIR	NMR
Full dataset	98.9% 2 layers with 20 hidden neurons	99.0% 2 layers with 17 hidden neurons
LASSO selection	98.9% 2 layers with 13 hidden neurons	98.3% 1 layer with 16 hidden neurons
Selection with standard deviation criterion	80.5% 1 layer with 16 hidden neurons	91.6% 1 layer with 12 hidden neurons

Selection with normalized standard deviation criterion	60.8% 2 layers with 12 hidden neurons	97.5% 1 layer with 15 hidden neurons
---	--	---

504

505

506

507 **Table 3**

508 Accuracy for the best ANNs evaluated on the full dataset of the two crop years.

Data	Merceological + NIR		NMR	
Crop year	2013/2014	2014/2015	2013/2014	2014/2015
Full dataset	99.5%	99.6%	99.4%	99.7%
LASSO selection	99.7%	99.2%	98.8%	99.5%
Selection with standard deviation criterion	81.6%	79.6%	92.4%	92.6%
Selection with normalized standard deviation criterion	57.5%	64.0%	97.8%	99.5%

509

510

511

512 **Table 4**

513 Accuracy for the best ANNs evaluated on the independent test of industrially produced samples.

Data	Merceological + NIR	NMR
Full dataset	47.0%	82.3%
LASSO selection	52.9%	88.2%
Selection with standard deviation criterion	47.0%	58.8%
Selection with normalized standard deviation criterion	23.5%	76.5%

514

515 **Figure Captions**

516 **Figure 1.** An example of MLP with 2 input neurons, 3 hidden neurons in the hidden layer and 1
517 output neuron

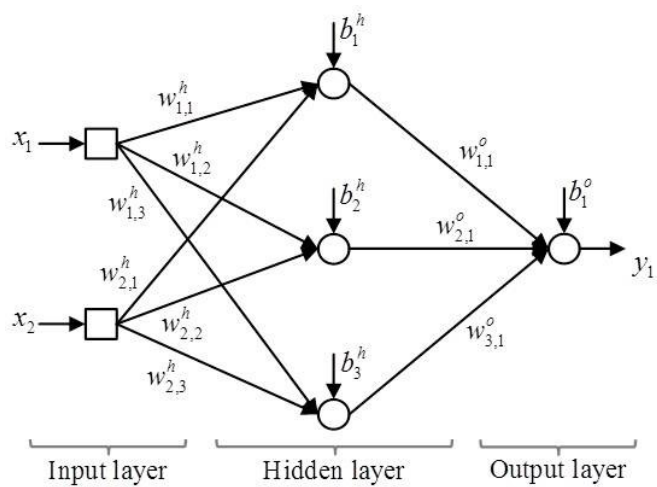
518 **Figure 2.** 3D PCA scoreplot for ^1H NMR-bucket-reduced spectra of monovarietal EVOO
519 samples, showing the general data grouping of all the samples.

520

521

522

ACCEPTED MANUSCRIPT

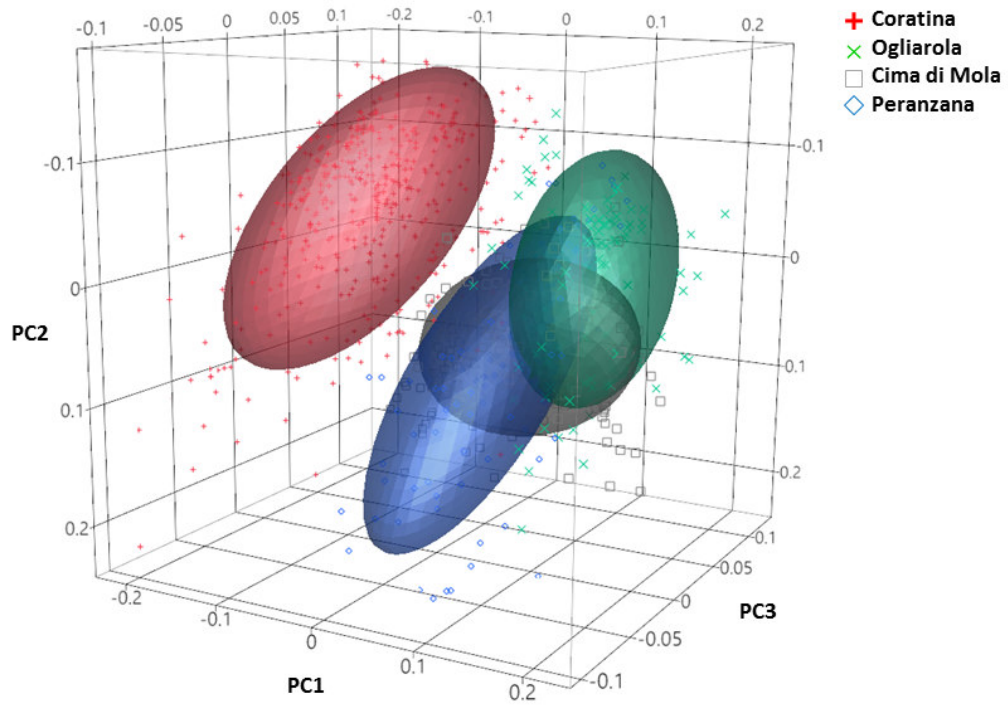


523

524

525

ACCEPTED MANUSCRIPT



526

527

528

ACCEPTED

- 529 1. Cultivar discriminating ability of ANNs on Apulian monocultivar EVOOs was studied.
- 530 2. Merceological, NIR and ^1H NMR data were used as ANNs training sets.
- 531 3. ANN models based on NMR data showed the highest accuracy in classifying cultivars.
- 532 4. The most information about cultivars was contained in very few NMR peaks.
- 533 5. Performance was not influence by the milling method nor the crop year.
- 534
- 535
- 536

ACCEPTED MANUSCRIPT