

1 **Discrimination of geographical origin of oranges (*Citrus sinensis***  
2 ***L. Osbeck*) by mass spectrometry-based electronic nose and**  
3 **characterization of volatile compounds**

4  
5 Valentina Centonze<sup>a</sup>, Vincenzo Lippolis<sup>b\*</sup>, Salvatore Cervellieri<sup>b</sup>, Anna Damascelli<sup>b</sup>, Grazia  
6 Casiello<sup>a</sup>, Michelangelo Pascale<sup>b</sup>, Antonio Francesco Logrieco<sup>b</sup>, Francesco Longobardi<sup>a,b</sup>

7  
8 <sup>a</sup>*Dipartimento di Chimica, Università di Bari "Aldo Moro", Via Orabona 4, 70126 Bari, Italy.*

9 <sup>b</sup>*Institute of Sciences of Food Production (ISPA), National Research Council of Italy (CNR),*  
10 *Via G. Amendola 122/O, 70126 Bari, Italy*

11  
12  
13  
14  
15  
16 \* Corresponding author: Tel: +39-080-5929457; fax +39-080-5929374;

17 e-mail address: [vincenzo.lippolis@ispa.cnr.it](mailto:vincenzo.lippolis@ispa.cnr.it)

18

19

20 **Abstract**

21 An untargeted method using headspace solid-phase microextraction coupled to electronic  
22 nose based on mass spectrometry (HS-SPME/MS-eNose) in combination with chemometrics  
23 was developed for the discrimination of oranges of three geographical origins (Italy, South  
24 Africa and Spain). Three multivariate statistical models, i.e. PCA/LDA, SELECT/LDA and  
25 PLS-DA, were built and relevant performances were compared. Among the tested models,  
26 SELECT/LDA provided the highest prediction abilities in cross-validation and external  
27 validation with mean values of 97.8% and 95.7%, respectively. Moreover, HS-SPME/GC-MS  
28 analysis was used to identify potential markers to distinguish the geographical origin of  
29 oranges. Although 28 out of 65 identified VOCs showed a different content in samples  
30 belonging to different classes, a pattern of analytes able to discriminate simultaneously  
31 samples of three origins was not found. These results indicate that the proposed MS-eNose  
32 method in combination with multivariate statistical analysis provided an effective and rapid  
33 tool for authentication of the orange's geographical origin.

34

35 **Key-words:** oranges, MS-based electronic nose, geographical origin, volatile compounds,  
36 chemometrics.

37

38

39        **1. Introduction**

40        Sweet orange (*Citrus sinensis* L. Osbeck) is one of the most popular fruits all over the world  
41        because it is very well-accepted by consumers for its nutritional, nutraceutical and sensorial  
42        attributes. Sweet oranges are usually classified into three main groups (i.e. common, navel,  
43        blood), with a diversification in terms of agronomical features within each group. Sweet  
44        orange accounting for 70% of worldwide citrus production is widely consumed both as fresh  
45        fruit as well as fruit juice. An annual global production was estimated at 73 million tons for  
46        oranges in 2016, and the main producers were Brazil, China, India and the United States  
47        (FAO, 2016).

48        Orange production at EU level was higher than 6 million tons for the 2016/2017 harvest and it  
49        is mainly concentrated in the Mediterranean basin with Spain and Italy representing about  
50        80% of the total yield, followed by Greece and Portugal (Citrus Annual, 2017).

51        For both the high productive vocation of its geographical area and the quality of products, Italy  
52        occupies a prominent position in the production of oranges that amounted to 1,2 million tons in  
53        the 2016/2017 season with geographic origin brands recognized by the European  
54        Commission (DOOR, 2018). Sicily and Calabria are the predominant production regions  
55        covering 80% of total Italian production (ISTAT, 2016). Although Italian oranges are  
56        considered of a premium quality, in the last few years the country has lost its leading role in  
57        the Mediterranean basin, due to the high costs of production and considerable loss of  
58        production due to the recent epidemic of the Citrus tristeza virus (CTV). In this contest, Italy  
59        commonly imports oranges from Spain and South-Africa which are gradually increasing their  
60        productions (Citrus Annual, 2017). Imports are mainly requested to cover the lack of  
61        availability of Italian products in some periods of the year (i.e. summer months), providing  
62        year-round availability for consumers. However, there are overlapping periods between Italian

63 and foreign productions, with an increased possibility for the consumer to buy mislabelled  
64 products with a lower product quality. Indeed, the geographical origin identification becomes  
65 more important for food with an origin label (e.g. "Made in Italy") since these products having  
66 acquired an "added-value" are more likely to become a target for frauds.

67 Considering that the illegal food trade is increasing around the world, traceability certifying the  
68 food authenticity including a correct labelling of origin is of great importance for traders,  
69 producers and consumers. Determination of food origin is commonly applied to control  
70 products with labelled or undeclared geographical provenances, for customs control and for  
71 self-control programs in the food industry.

72 For this reason, the development of rapid and reliable analytical methods to assess the  
73 geographical origin of food is highly demanded. For food authentication, non-targeted  
74 analysis (fingerprint) in combination with multivariate statistical analysis is a promising  
75 approach that allowing the detection of many metabolites as possible and permits to classify  
76 samples based on pattern of metabolites. A variety of non-targeted analytical techniques,  
77 mainly based on vibrational spectroscopy, mass spectrometry and NMR, have been applied  
78 to discriminate geographical origin of several food matrices and recently reviewed (Cubero-  
79 Leon et al. 2014, Essingler et al. 2014, Danezis et al. 2016), including fresh oranges (Diaz et  
80 al. 2014, Jandric and Cannavan 2017). Among these techniques, non-chromatographic mass  
81 spectrometry (MS) is an emerging approach for food authentication studies due to its several  
82 advantages in terms of rapidity, sensitivity, selectivity and high-throughput analysis (Danezis  
83 et al. 2016). In particular, the mass spectrometry-based electronic nose (MS-eNose)  
84 technique based on the use of headspace solid-phase microextraction (HS-SPME) directly  
85 coupled to MS is one of the most innovative approach that can be used to analyse volatile  
86 organic compounds (VOCs) of complex matrices. This technique provides a global mass

87 spectrometric fingerprint of VOCs of a sample, analyzed without chromatographic separation,  
88 in which each m/z ratio acts as a “sensor” whose intensity derives from the contribution of  
89 each compound producing that fragment. The main advantages of the proposed methodology  
90 are the minimum sample preparation required and the speed of analysis. The MS-eNose  
91 technique has been successfully applied to characterize several food matrices (i.e. coffee,  
92 raw spirits, milk and honey) for different purposes including the discrimination of geographical  
93 origin (Perez Pavòn et al. 2006, Jeleń et al. 2010, Liberto et al. 2013, Smyth and Cozzolino  
94 2013).

95 Indeed, the aroma is one of the most important factors to discriminate food products and  
96 determine their quality. In particular, the aroma of fresh squeezed orange juice is a complex  
97 mixture of VOCs that is related to several factors, including orange cultivar, environment,  
98 geographical origin, degree of ripeness and storage conditions (Perez-Cacho and Rouseff  
99 2008a, Cuevas et al. 2017). The VOCs of orange consist of esters, alcohols, aldehydes,  
100 ketones, terpenes and furans (Perez-Cacho and Rouseff 2008a, Cuevas et al. 2017). VOCs  
101 of fresh oranges and orange juices have been successfully characterised for different  
102 purposes by several studies using GC-MS analysis (Cuevas et al. 2017, Cerdà-Calero et al.  
103 2013, Reinhard et al. 2008, Cerdà-Calero et al. 2012, Reid 2003, Zierler et al. 2004).  
104 However, to date the MS-eNose technique has not been applied for the discrimination of  
105 geographical origin of oranges.

106 For this reason, the aim of this study was to demonstrate the feasibility of MS-eNose  
107 technique applied to the VOCs analysis for the discrimination of the geographic origin of  
108 oranges. In particular, a robust and suitable non-targeted MS-based electronic nose method  
109 in combination with multivariate statistical analysis was developed and validated for the  
110 discrimination of oranges of three different geographical origins, i.e. Italy, South Africa and

111 Spain. Moreover, an HS-SPME/GC-MS method was used to characterize the possible pattern  
112 of volatile compounds having a role in this discrimination.

113

114

## 115 **2. Materials and Methods**

### 116 *2.1 Chemicals and reagents*

117 Methanol (HPLC grade) and (E)-3-hexen-1-ol ( $\geq 98\%$ ) was purchased from Sigma Aldrich  
118 (Milan, Italy). Ten milliliters headspace vials with magnetic screw cap containing a pierceable  
119 PTFE/silicon septa were purchased from Agilent Technologies (Palo Alto, CA, USA). Helium  
120 at a purity of 99.9995% was obtained by Sapio s.r.l. (Bari, Italy). The automatic solid-phase  
121 microextraction (SPME) fiber holder was obtained from Gerstel (Mulheim an der Ruhr,  
122 Germany).

123 SPME-Fast Fit Fiber Assembly (FFA) divinylbenzene/carboxen/polydimethylsiloxane  
124 (DVB/CAR/PDMS, 50/30  $\mu\text{m}$ , 1 cm fiber length), SPME Fiber Assembly for manual use  
125 divinylbenzene/carboxen/polydimethylsiloxane (DVB/CAR/PDMS, 50/30  $\mu\text{m}$ , 1 cm fiber  
126 length), polydimethylsiloxane/divinylbenzene (PDMS/DVB, 65  $\mu\text{m}$ , 1 cm fiber length),  
127 carboxen/polydimethylsiloxane (CAR/PDMS, 85  $\mu\text{m}$ , 1 cm fiber length) and the manual SPME  
128 holder were purchased from Supelco (Bellafonte, PA, USA).

129

### 130 *2.2 Samples collection and sample preparation*

131 Orange samples of the 2014/2015 crop season were collected from producers. A total of 137  
132 samples of different cultivars of three different geographical origins, i.e. Italy, South Africa and  
133 Spain, was collected. Table 1 reports the number and the cultivars of samples collected for  
134 each geographic area.

135 The collected samples (five oranges for each sample) were squeezed and the juice was  
136 frozen at -20 ° C until the analyses. The stored juice after thawing was centrifuged for 20 min  
137 at 13000 rpm. Aliquots (2g) of the supernatant were placed in 10 mL headspace vials, adding  
138 as internal standard (E)-3-hexen-1-ol in methanol to obtain a concentration of 2 µg/g. Then,  
139 vials were sealed for the analysis with both analytical methods (i.e. HS-SPME/MS-eNose and  
140 HS-SPME/GC-MS). The extraction, desorbition and sample introduction of the samples were  
141 performed automatically in HS-SPME/MS-eNose and manually in HS-SPME/GC-MS analysis.

142

### 143 2.3 HS-SPME/MS-eNose analysis

144 The squeezed orange juice samples were analysed by the mass spectrometry-based  
145 electronic nose (MS-eNose) GERSTEL Headspace ChemSensor System (GERSTEL,  
146 Mülheim, Germany) consisted of a headspace multi-purpose sampler MPS 2 (Gerstel,  
147 Mulheim an der Ruhr, Germany) and the Agilent 7890A GC System (Agilent Technologies,  
148 Palo Alto, CA, USA), modified for non-separative analysis with a deactivated fused-silica  
149 tubing (transfer column, 10 m x 0.18 mm i.d., 0 µm film thickness, Agilent Technologies),  
150 coupled to the Agilent 5975C inert MSD mass spectrometer. Moreover, the MS-eNose was  
151 online integrated with a multi-purpose sampler MPS 2 (Gerstel, Mulheim an der Ruhr,  
152 Germany), which was equipped with headspace incubation chamber and SPME sampling  
153 unit. An HS-SPME/MS-eNose protocol of analysis was *in-house* optimized according the  
154 procedure reported by Cefola et al. 2018. In particular, the headspace vial was kept at  
155 temperature of 40 °C for 10 min in the incubator-agitator of the MPS 2 autosampler to  
156 generate the headspace. The extraction from the headspace was performed by exposing a  
157 divinylbenzene/carboxen/polydimethylsiloxane (DVB/CAR/PDMS) fiber at 40 °C for 30 min.  
158 After extraction, compounds were thermally desorbed exposing the fiber in the CIS-4

159 programmed temperature vaporization (PTV) injector (Gerstel) of the MS-eNose at 250 °C for  
160 5 min. Then, the MS-eNose analyses were carried out for 5 min using the following  
161 experimental conditions: the injection port fitted with a 1 mm i.d. liner was maintained at 250  
162 °C in splitless mode; the oven, transfer line, ion source and quadrupole temperatures were  
163 180, 280, 230 and 150 °C, respectively; the helium flow rate was held constant at 1 mL/min;  
164 Electron impact Ionization (EI+) mode with an electron energy of 70 eV was used; the mass  
165 spectrometer acquired data in full scan mode (scan range: 40–300 amu).  
166 For each analysis mass spectral fingerprint was obtained by the software Chemsensor 6.912  
167 (Gerstel, Mülheim and der Ruhr, Germany) corresponding to the sum of mass spectra  
168 obtained in the time range 0.22-2.0 min. Mass intensities of mass spectral fingerprint were  
169 estimated as relative abundances by comparing the mass intensity of each ion with the  
170 intensity of ion at 43 amu of the internal standard (i.e. (E)-3-hexen-1-ol).

171

#### 172 *2.4 HS-SPME/GC-MS analysis*

173 A subset of 27 orange samples of different geographical origin (9 samples for each origin)  
174 were randomly selected from the entire set of 137 samples and analysed in duplicate by HS-  
175 SPME/GC–MS. The extraction and desorption steps of volatile compounds were performed  
176 following the same experimental parameters optimized for the HS-SPME/MS-eNose method,  
177 while GC-MS analysis was carried out according the procedure reported by Cefola et al. 2018  
178 with some modifications. In particular, the GC-MS analyses were carried out by an Agilent  
179 6890 Series GC system (Agilent Technologies, Palo Alto, CA, USA) equipped with a VF-  
180 WAXms (60 m x 0.25 mm i.d., 0.25 µm film thickness, Agilent Technologies) fused-silica  
181 capillary column and coupled to an Agilent 5973 Network Mass Selective Detector mass  
182 spectrometer. The injection port fitted with a 0.75 mm i.d. liner was maintained at 250 °C in



183 splitless mode. The analyses were performed with programmed temperature: initial  
184 temperature 40 °C maintained for 6 min, from 40 to 120 °C at 2 °C/min, 120 to 230 °C at 10  
185 °C/min, the final temperature being maintained for 10 min. The helium flow rate was held  
186 constant at 1 mL/min. The transfer line, ion source and quadrupole temperatures were 280,  
187 290 and 150 °C, respectively. Electron impact Ionization (EI+) mode with an electron energy  
188 of 70 eV was used. The mass spectrometer acquired data in full scan mode (scan range: 40–  
189 300 amu). The compounds were identified by comparison of experimental mass spectra with  
190 ones present in the NIST v2.0 and Wiley 138 libraries using a match quality higher than 70.  
191 The identification of volatile compounds was also verified by comparison of their linear  
192 retention indices (LRI) determined in relation to the retention times of C5–C14 and C8-C40 n-  
193 alkanes series, with those reported in literature [Zellner et al., 2008]. Quantification of  
194 compounds was performed by the same method of internal standardization used for HS-  
195 SPME/MS-eNose analysis. The amount of each identified compound was estimated by  
196 comparing the total ion current (TIC) peak area with (E)-3-hexen-1-ol peak area and  
197 expressed as area ratio. All mass spectrum fingerprints were combined to obtain a data  
198 matrix containing 137 objects and 260 variables that was submitted to statistical analyses.

199

## 200 *2.5 Statistical data analysis*

201 Before chemometric analysis, data obtained by the Chemsensor 6.912 software were pre-  
202 treated by baseline correction, through noise subtraction, and by internal normalization of the  
203 signal from each sample (Perez Pavon et al. 2006). The internal normalization allowed to  
204 remove the efficiency loss during the extraction process of volatile components due to the  
205 variability of the SPME fiber performance and MS signal instability. Subsequently, data were  
206 pre-processed by Pareto scaling and then submitted to multivariate statistical analyses.

207 The presence of outliers was evaluated observing the influence plot obtained by applying  
208 PCA (Principal Component Analysis) for each single class of different geographical origin.  
209 Samples identified as extreme outliers will be excluded. For this reason, NIPALS (Non-linear  
210 Iterative Partial Least Squares) algorithm was applied, considering V-fold equal to 10 in the  
211 cross validation process (CV=10), establish the exact number of PCs to use to build PCA  
212 models. PCA was also applied as exploratory technique with the aim to visualize if sample  
213 clustering was present as a function of the geographical origin of the samples (Jolliffe 2002).  
214 Then, supervised pattern recognition techniques, i.e. Linear Discriminant Analysis (LDA) and  
215 Partial Least Squares Discriminant Analysis (PLS-DA) (Oliveri and Downey 2012), were used  
216 in order to classify orange samples on the basis of their geographical origin. For this purpose,  
217 the data matrix was divided in two subsets: a modeling set (containing 90 samples) and a test  
218 set (containing 47 samples). In particular, the modeling set, built using only Navel orange  
219 samples, was represented by 30 samples for each class (different geographical origin) and  
220 was used to build the statistical models, while the test set, consisting of 19 Italian, 22 African  
221 and 6 Spanish samples, was used to their validation.

222 In the case of LDA, to prevent model overfitting two different strategies PCA (unsupervised  
223 approach) and SELECT (a supervised feature selection algorithm) were used to reduce the  
224 number of variables that exceeded the number of objects (Berrueta et al. 2007, Casale et al.  
225 2010, Vandeginste et al 1998). In particular, the number of variables should not exceed  $(n -$   
226  $g)/3$ , where n is the number of objects and g is the number of categories, i.e. 29, considering  
227 90 objects (number of samples) and 3 categories (number of geographical origins).

228 On the other hand, PLS-DA was used as an alternative approach to avoid variables reduction  
229 being it frequently used in the case of large number of variables (Massart et al. 1997, Oliveri ,  
230 2017).

231 The PCA/LDA and PLS-DA models were built evaluating the proper number of principal  
232 components and latent variables, respectively, which returned the lowest root mean square  
233 error of cross validation (RMSECV). This parameter can guarantee that feature variables are  
234 collected as much as possible and they are not overfitted. Therefore, performances of the  
235 PCA/LDA, SELECT/LDA and PLS-DA models were compared in terms of recognition ability,  
236 i.e. its ability to correctly classify the samples used for the building of the model, prediction  
237 ability in cross-validation (CV), i.e. its ability to correctly classify samples of a test set  
238 generated in a V-fold cross validation (with V equal to 10) and prediction ability in external  
239 validation calculated using the test set.

240 Univariate statistical analysis, i.e. one-way analysis of variance (ANOVA) followed by a post  
241 hoc Tukey's honestly significant difference (HSD) test ( $p < 0.05$ ), was performed to assess  
242 the differences between mean peak area ratios of identified volatile molecules of orange  
243 samples of three different geographical origins obtained by HS-SPME/GC-MS analysis.

244 Data analyses were performed by using Pirouette software ver. 4.0 (Infometrix, Inc., Bothell,  
245 WA, USA), V-Parvus release 2010 (<http://www.parvus.unige.it>, Genova, Italy), Classification  
246 Toolbox in Matlab (Mathworks Inc., Natick, Massachusetts, USA) and Statistica 6.0 (StatSoft,  
247 Tulsa).

248

249

### 250 **3. Results and Discussion**

#### 251 *3.1 Geographical origin discrimination using HS-SPME/MS-eNose*

252 In order to find anomalous samples (outliers), data were processed in specific PCA models  
253 for each geographical origin (class) showing that 9, 11, and 9 PCs explained 97.0, 97.0, and  
254 95.0% of the total variance, for the Italian, South African, and Spanish origin, respectively.

255 Influence plots obtained by plotting the Mahalanobis distance versus sample residual showed  
256 that all the samples coming from a specific class fit in the respective model then excluding the  
257 presence of outliers.

258 Subsequently, to get a general overview of the data distributions an explorative PCA was  
259 applied on all data and by plotting the PC1 vs. PC2 sample scores (Figure 1) a poor visual  
260 clustering of the objects based on their geographical origin was showed (PC1 and PC2  
261 explained respectively 42.5% and 14.5% of the total variance). Moreover, no significant  
262 separation was evidenced when the score plots of the remaining PCs were observed. This  
263 aspect was also confirmed evaluating the PC Fisher weights (FW) values, i.e. the measure of  
264 the between-class variance/within-class variance ratio, that resulted to be considerably lower  
265 than 1 (data not shown), meaning that no single PC was sufficiently suitable to distinguish  
266 samples for their geographical origin (Harper et al. 1977). Therefore, these results highlighted  
267 the necessity to use supervised techniques, i.e. discriminant techniques such as LDA and  
268 PLS-DA. These classification techniques were applied to data matrix divided into the two data  
269 subsets: modeling set and test set. The overall results of these classification models are  
270 reported and compared in Table 2.

271 In the case of LDA, two variable reduction strategies, i.e. PCA and SELECT were adopted to  
272 avoid model overfitting. In particular, PCA was applied to compress the information and the  
273 number of PCs chosen to get the lowest error in prediction cross validation and then used to  
274 build the PCA/LDA model was of 13 (CV procedure,  $V=10$ ). The PCA/LDA model provided  
275 mean values of the recognition (classification) ability and CV prediction ability of 82.2% and  
276 78.9%, respectively (Table 2). In particular, the model correctly predicted 21/30 Italian  
277 samples, 22/30 South African samples and 28/30 Spanish samples. Despite these low  
278 performances, the applicability of the model was also evaluated by an external test obtaining

279 a similar mean prediction ability of 80.9%. In the case of SELECT, 29 variables out of 260  
280 were selected and then used to build the SELECT/LDA model. Mean percentages of  
281 recognition ability and CV prediction ability obtained using SELECT/LDA model, were 100.0%  
282 and 97.8%, respectively (Table 2). In particular, this model permitted to correctly classify all  
283 South African and Spanish samples and 28 samples out of 30 Italian samples giving a value  
284 of specific prediction percentage of 93.3% for this class.

285 Therefore, the supervised selection approach SELECT allowed significantly improvement of  
286 results in terms of recognition and prediction abilities than those obtained using the  
287 unsupervised PCA compression method. These results can be justified by considering that  
288 the direction of maximum variability of data, used in the PCA approach, of the data could not  
289 correspond to the direction of maximum discrimination among defined classes. Indeed, if the  
290 variability associated with geographical origin is small with respect to the total variability, the  
291 use of PCA variable reduction method can partially hide its specific feature contribution. On  
292 the other hand, SELECT, by choosing the variables that contain the best information for the  
293 under study classification, provides decorrelated variables avoiding redundant information. In  
294 order to confirm these results, the SELECT/LDA model was also validated using the external  
295 set. An external prediction ability of 95.7% was obtained with all South African and Spanish  
296 samples correctly recognized (specific prediction rates of 100.0%) while only 2 samples out of  
297 19 Italian samples were not correctly assigned, with a specific prediction rates of 89.5%.

298 Furthermore, PLS-DA was applied to test an alternative multivariate statistical approach of  
299 classification and avoiding the process of variables reduction. By implementing a 10-fold  
300 cross-validation, 12 latent variables guaranteed the optimal model complexity, leading to a  
301 97.8% average recognition rate. In particular, the totality of the Italian and Spanish samples  
302 were correctly classified, and only two out of 30 African samples was not correctly assigned.

303 The average CV prediction rate was 85.6% with CV prediction abilities for the Italian, African  
304 and Spanish categories of 83.3%, 76.7% and 96.7%, respectively. Moreover, the external  
305 validation procedure provided prediction abilities of 84.2% for Italy, 81.8% for South Africa  
306 and 100.0% for Spain orange samples, corresponding to an average prediction rate of 85.1%.  
307 These results showed that although PLS-DA model permitted acceptable prediction abilities,  
308 they were significantly lower than those obtained by the SELECT/LDA model.  
309 These results demonstrated that HS-SPME/MS-eNose experimental data contain enough  
310 information to allow the construction of appropriate models for the discrimination of orange  
311 samples on the basis of their geographical origin.

312

### 313 *3.2 Characterization of the pattern of volatile compounds by HS-SPME/GC-MS*

314 In order to identify the most important volatile organic compounds (VOCs) to be used as  
315 markers in distinguishing oranges according to the country of origin, 27 orange samples of the  
316 three different geographical origins (9 for each class) were analysed by HS-SPME/GC-MS  
317 technique under the optimized experimental conditions. A total of 65 VOCs have been  
318 identified belonging to a wide range of chemical classes including aldehydes (5), ketones (4),  
319 esters (16), acids (3), alcohols (8), terpenes (20), heterocyclic compounds (4), saturated,  
320 unsaturated and aromatic hydrocarbons (5) (Table 3). In particular, aldehydes as hexanal and  
321 terpenes as D-limonene and  $\beta$ -linalool are responsible for characteristic orange juice flavour.  
322 Conversely, fruity notes are mainly due to ethyl butanoate with minor contributions from ethyl  
323 2-methylpropanoate and ethyl 2-methylbutanoate (Perez-Cacho et al., 2008b). Moreover,  
324 high content of ester compounds have been shown to discriminate oranges obtained under  
325 conventional procedures from oranges cultivated under organic ones, characterized by high  
326 content of some terpenes and neryl acetate and geranyl acetate (Cuevas et al, 2017).

327 Accordingly, the composition (calculated as ratio of the analyte peak area relative to (E)-3-  
328 hexen-1-ol peak area) of fresh squeezed orange juice headspaces was investigated by one-  
329 way ANOVA analysis followed by a *post hoc* Tukey's HSD test in order to detect molecules  
330 discriminating samples in relation to their geographical origins. Among the identified  
331 compounds, the content of 28 molecules was significantly different among samples belonging  
332 to the three classes ( $p < 0.05$ ). As reported in Table 3, 13 analytes allowed to discriminate the  
333 geographical differences between Italian oranges and samples of the other two classes while  
334 only 2 volatile compounds discriminated South African oranges from the others. Moreover, 8  
335 and 3 molecules differentiated Italian samples from Spanish and South African samples,  
336 respectively. Moreover, Spanish oranges were distinguished from the samples of the other  
337 two classes by the high content of methyl butanoate and only from South African samples by  
338 the presence of 1-terpinen-4-ol.

339 Among the selected analytes, (E)-2-hexen-1-ol, (Z)-3-hexen-1-ol, (E)- $\beta$ -ionone and 6-methyl-  
340 5-hepten-2-one in the volatile fraction of Italian oranges showed the highest increase of  
341 contents from 7 to 19 times with respect to those measured for South African and Spanish  
342 oranges. (E)-2-hexen-1-ol and (Z)-3-hexen-1-ol have been already reported as volatile  
343 compounds of fresh prepared orange juice (Perez-Cacho et al., 2008b). Moreover, 4-methyl-  
344 heptane and 1-penten-3-one were not detected in the headspace of Italian samples while  
345 ethyl octanoate, 1-octen-3-ol and  $\beta$ -caryophyllene were absent in volatile fraction of Spanish  
346 oranges. Most molecules in the selected pattern have already been associated to fresh  
347 orange juice (Cuevas et al, 2017; Bai et al., 2014; Perez-Cacho et al., 2008b; Sádecká et al.,  
348 2014) with the exception of 2-methyl-1-pentene, 4-methyl-heptane, o-cymene and 4-methyl-2-  
349 heptanone, which were related for the first time to orange fresh squeezed juices.

350 However, although the HS-SPME/GC–MS analysis clearly highlighted difference in the VOCs  
351 profile of the oranges from the three geographical origins, a pattern of analytes able to  
352 discriminate simultaneously samples of the three different origins was not found.  
353 Consequently, the application of the multivariate analysis to the whole dataset obtained by  
354 MS-eNose analysis was confirmed to be the most appropriate approach to permit the rapid  
355 prediction of geographical origin of oranges.

356

#### 357 **4. Conclusion**

358 In this study, a rapid and inexpensive method based on MS-eNose analysis in combination  
359 with chemometrics was successfully used to classify orange samples of three different  
360 geographical origins, i.e. Italy, South Africa and Spain.

361 In particular, three multivariate statistical approaches, i.e. PCA-LDA, SELECT-LDA and PLS-  
362 DA, were tested. Although, all tested statistical models permitted acceptable recognition and  
363 prediction abilities, the SELECT/LDA model showed the highest percentages in terms of  
364 prediction ability in cross-validation and external validation, with average values of 97.8% and  
365 95.7%, respectively. The performances of the proposed method makes it suitable as powerful  
366 tool to assess the authenticity of oranges.

367 Although, HS-SPME/GC–MS analysis showed the absence of specific markers, differences in  
368 the pattern and content of VOCs of orange samples of the three different geographical origins  
369 were observed confirming the validity of the multivariate statistical approach used in this  
370 study.

371

372



373 **References**

- 374 Bai, J., Baldwin, E., Hearn, J., Driggers, R., & Stover, E. (2014). Volatile profile comparison of  
375 usda sweet orange-like hybrids versus 'Hamlin' and 'Ambersweet'. *Hortscience*, *49*, 1262–  
376 1267.
- 377 Berrueta, L. A., Alonso-Salces, R. M., & Héberger, K. (2007). Supervised pattern recognition  
378 in food analysis. *Journal of Chromatography A*, *1158*, 196–214.
- 379 Casale, M., Casolino, C., Oliveri, P., & Forina, M. (2010). The potential of coupling information  
380 using three analytical techniques for identifying the geographical origin of Liguria extra  
381 virgin olive oil. *Food Chemistry*, *118*(1), 163-170.
- 382 Cefola, M., Damascelli, A., Lippolis, V., Cervellieri, S., Linsalata, V., Logrieco, A., & Pace, B.  
383 (2018). Relationships among volatile metabolites, quality and sensory parameters of 'Italia'  
384 table grapes assessed during cold storage in low or high CO<sub>2</sub> modified atmospheres.  
385 *Postharvest Biology and Technology*, *142*, 124-134.
- 386 Cerdán-Calero M., María Sendra J., & Sentandreu E., (2012). Gas chromatography coupled  
387 to mass spectrometry analysis of volatiles, sugars, organic acids and aminoacids in  
388 Valencia Late orange juice and reliability of the Automated Mass Spectral Deconvolution  
389 and Identification System for their automatic identification and quantification. *Journal of*  
390 *Chromatography A*, *1241*, 84-95.
- 391 Cerdán-Calero M., Izquierdo L., & Sentandreu E., (2013). Valencia Late orange juice  
392 preserved by pulp reduction and high pressure homogenization: Sensory quality and gas  
393 chromatography–mass spectrometry analysis of volatiles. *LWT - Food Science and*  
394 *Technology*, *51*, 476–483.

395 Citrus Annual, (2017). USDA Foreign Agricultural Service – Global Agricultural Information  
396 Network. 31.12.2017. <https://www.fas.usda.gov/data/eu-28-citrus-annual-2>. Accessed  
397 17.07.2018.

398 Cubero-Leon, E., Peñalver, R., & Maquet, A., (2014). Review on metabolomics for food  
399 authentication. *Food Research International*, 60, 95–107.

400 Cuevas F. J., Moreno-Rojas J. M., & Ruiz-Moreno M. J., (2017). Assessing a traceability  
401 technique in fresh oranges (*Citrus sinensis*L. Osbeck) with an HS-SPME-GC-MS method.  
402 Towards a volatile characterisation of organic oranges. *Food Chemistry*, 221,1930–1938.

403 Danezis G. P., Tsagkaris A. S., Brusica V., & Georgiou C. A., (2016). Food authentication: state  
404 of the art and prospects. *Current Opinion in Food Science*, 10, 22–31.

405 Díaz R., Pozo O. J., Sancho J. V., & Hernández F., (2014). Metabolomic approaches for  
406 orange origin discrimination by ultra-high performance liquid chromatography coupled to  
407 quadrupole time-of-flight mass spectrometry. *Food Chemistry*, 157, 84–93.

408 DOOR (Database of Origin and Registration), (2018). The Database of Origin and  
409 Registration (DOOR). European Commission, Agricultural and Rural Development, online  
410 database. URL <http://ec.europa.eu/agriculture/quality/door/list.html>. Accessed 17.07.2018

411 Esslinger, S., Riedl J., & Fauhl-Hassek C., (2014). Potential and limitations of non-targeted  
412 fingerprinting for authentication of food in official control. *Food Research International*, 60,  
413 189–204.

414 FAOSTAT (Food and Agriculture Organization of the United Nations). Crops data, 2016. URL  
415 [1http://www.fao.org/faostat/en/#data/QC](http://www.fao.org/faostat/en/#data/QC). Accessed 17.07.2018

416 Harper, A. M., Duewer, D. L., Kowalski, B. R., & Fashing J. L. (1977). ARTHUR and  
417 experimental data analysis: The Heuristic use of a polyalgorithm. In: B. R. Kowalski (Ed.),

418 ACS Symposium series 52. Chemometrics: Theory and Application (pp. 14–52).  
419 Washington: American Chemical Society.

420 ISTAT (Istituto Nazionale di Statistica). Tavola C24 - Superficie (ettari) e produzione (quintali):  
421 arancio, mandarino, clementina, limone, 2016. URL  
422 [http://agri.istat.it/jsp/dawinci.jsp?q=plC240000010000012000&an=2016&ig=1&ct=506&id=](http://agri.istat.it/jsp/dawinci.jsp?q=plC240000010000012000&an=2016&ig=1&ct=506&id=15A|21A|31A)  
423 [15A|21A|31A](http://agri.istat.it/jsp/dawinci.jsp?q=plC240000010000012000&an=2016&ig=1&ct=506&id=15A|21A|31A). Accessed 17.07.2018

424 Jandric, Z., & Cannavan, A., (2017). An investigative study on differentiation of citrus fruit/fruit  
425 juices by UPLC-QToF MS and chemometrics. *Food Control*, 72, 173-180.

426 Jeleń, H. H.; Ziolkowska, A.; & Kaczmarek, A., (2010). Identification of the botanical origin of  
427 raw spirits produced from rye, potato, and corn based on volatile compounds analysis  
428 using a SPME-MS method. *Journal of Agricultural and Food Chemistry*, 58, 12585–12591.

429 Jolliffe, I. T., (2002). *Principal component analysis* (2nd ed.). New York: Springer.

430 Kim, H. K., & Verpoorte, R. (2010). Sample preparation for plant metabolomics.  
431 *Phytochemical Analysis*, 21, 4–13.

432 Liberto E., Ruosi M. R, Cordero C., Rubiolo P., Bicchi C. & Sgorbini B., (2013). Non-  
433 separative Headspace Solid Phase Microextraction–Mass Spectrometry Profile as a  
434 Marker To Monitor Coffee Roasting Degree. *Journal of Agricultural and Food Chemistry*,  
435 61, 1652–1660.

436 Massart, D. L., Vandeginste, B. G. M., Buydens, L. M. C., De Jong S., Lewi P. J. & Smeyers-  
437 Verbeke J., (1997). *Handbook of Chemometrics and Qualimetrics: Part A*. Amsterdam:  
438 Elsevier.

439 Oliveri, P. & Downey, G. (2012). Multivariate class modelling for the verification of food-  
440 authenticity claims. *Trends in Analytical Chemistry*, 35, 74–86.

- 441 Oliveri, P., (2017). Class-modelling in food analytical chemistry: Development, sampling,  
442 optimisation and validation issues – A tutorial. *Analytica Chimica Acta*, 982, 9-19.
- 443 Perez-Cacho P. R. & Rouseff R. (2008a). Processing and storage effects on orange juice  
444 aroma: A review. *Journal of Agricultural and Food Chemistry*, 56, 9785–9796.
- 445 Perez-Cacho, P. R. & Rouseff, R. L. (2008b). Fresh Squeezed Orange Juice Odor: A Review.  
446 *Critical Reviews in Food Science and Nutrition*, 48, 681–695.
- 447 Pérez Pavón, J. L. P., del Noyal Sanchez, M., Pinto, C. G., Laespada, E. F., Cordero, B. M.,  
448 & Peña, A. G., (2006). Strategies for qualitative and quantitative analyses with mass  
449 spectrometry-based electronic noses. *Trends in Analytical Chemistry*, 25, 257–266.
- 450 Reid W.J., (2003). Instrumental methods in detecting taints and off-flavours. In B. Baigrie  
451 (Eds.) *Taints and Off-Flavours in Foods, A volume in Woodhead Publishing Series in Food*  
452 *Science, Technology and Nutrition* (pp 31–63), Cambridge : Woodhead Publishing.
- 453 Reinhard H., Sager F., & Zoller O., (2008). Citrus juice classification by SPME-GC-MS and  
454 electronic nose measurements. *LWT - Food Science and Technology*, 41, 1906–1912.
- 455 Sádecká, J., Polovka, M., Kolek, E., Belajová, E., Tobolková, B., Daško, L. & Durec, J.  
456 (2014). Orange juice with pulp: impact of pasteurization and storage on flavour,  
457 polyphenols, ascorbic acid and antioxidant activity. *Journal of Food & Nutrition Research*,  
458 53, 371-388.
- 459 Smyth, H., & Cozzolino, D., (2013). Instrumental Methods (Spectroscopy, Electronic Nose,  
460 and Tongue) As Tools To Predict Taste and Aroma in Beverages: Advantages and  
461 Limitations. *Chemical Reviews*, 113, 1429–1440.
- 462 Vandeginste, B.G.M., Massart, D.L., Buydens, L.M.C., De Jong, S., Lewi, P.J., & Smeyers-  
463 Verbeke, J. (1998). Handbook of Chemometrics and Qualimetrics: Part. B. In B.G.M.

464 Vandeginste, & S.C. Rutan (Eds.), *Supervised Pattern Recognition* (pp. 207-241).

465 Amsterdam: Elsevier.

466 Zellner, B. D. A., Bicchi, C., Dugo, P., Rubiolo, P., Dugo, G., & Mondello, L. (2008). Linear

467 retention indices in gas chromatographic analysis: a review. *Flavour and fragrance journal*,

468 23, 297–314.

469 Zierler, B., Siegmund, B., Pfannhauser, W., (2004). Determination of off-flavour compounds in

470 apple juice caused by microorganisms using headspace solid phase microextraction–gas

471 chromatography–mass spectrometry. *Analytica Chimica Acta*, 520, 3–11.

472

473

474

475

476

477

478

479 **Figure captions**

480

481 **Figure 1.** *PC1 vs PC2 scatter plot for orange samples. Geographical origins: Italy (o), Africa*

482 *(□), Spain (+).*

483

**Table 1**

Orange samples for each geographic area

<b>Geographical origin</b>	<b>Italy</b>	<b>South Africa</b>	<b>Spain</b>
N. Samples	49	52	36
Cultivars	Washington Navel, Newhall Navel, Navel Foglia, Navelina Bionda IGP, Duretta IGP, Ovale Valencia, Lane Late	Navel, Navelina navel	Lane Late, Navel Pawel Navel

**Table 2**

Recognition, CV prediction abilities and external prediction for all models built classifying oranges samples according to their geographical origin.

Model performance (%)												
	Recognition ability (Modelling)				Prediction ability (CV <sup>d</sup> 10)				External Prediction			
	ITA <sup>a</sup>	S.A. <sup>b</sup>	SPA <sup>c</sup>	Mean	ITA	S.A.	SPA	Mean	ITA	S.A.	SPA	Mean
<b>PCA/LDA</b> (13 Principal Components)	76.7 (23/30)	80.0 (24/30)	90.0 (27/30)	82.2	70.0 (21/30)	73.3 (22/30)	93.3 (28/30)	78.9	68.4 (13/19)	86.4 (19/22)	100.0 (6/6)	80.9
<b>SELECT/LDA</b> (29 variables)	100.0 (30/30)	100.0 (30/30)	100.0 (30/30)	100.0	93.3 (28/30)	100.0 (30/30)	100.0 (30/30)	97.8	89.5 (17/19)	100.0 (22/22)	100.0 (6/6)	95.7
<b>PLS-DA</b> (12 Latent Variables)	100.0 (30/30)	93.3 (28/30)	100.0 (30/30)	97.8	83.3 (25/30)	76.7 (23/30)	96.7 (29/30)	85.6	84.2 (16/19)	81.8 (18/22)	100.0 (6/6)	85.1

<sup>a</sup>: Italy; <sup>b</sup>: South Africa; <sup>c</sup>: Spain; <sup>d</sup>: Cross Validation.

**Table 3**

Volatile compounds (n = 65) identified by HS-SPME/GC–MS analysis of Italian, South African and Spanish oranges.

<b>Volatile Compound</b>	<b>Code</b>	<b>LRI<sub>It</sub>/LRI<sub>Sp</sub><sup>d</sup></b>	<b>Volatile Compound</b>	<b>Code</b>	<b>LRI<sub>It</sub>/LRI<sub>Sp</sub></b>
<b>Hydrocarbons</b>			<b>Alcohols</b>		
2-methyl-1-pentene <sup>b</sup>	1	- <sup>e</sup> /644	2-methyl-2-propanol	4	900/904
4-methyl-heptane <sup>b</sup>	2	790/765	ethanol	5	937/937
undecane	15	1100/1100	(Z)-2-penten-1-ol	30	1329/1330
1,3-bis(1,1-dimethylethyl)-benzene	36	1423/1433	(Z)-3-hexen-1-ol <sup>a,b</sup>	34	1393/1394
4-acetyl-1-methylcyclohexene	43	1568/1566	(E)-2-hexen-1-ol <sup>a,b</sup>	35	1417/1416
<b>Terpenes</b>			1-octen-3-ol <sup>b</sup>	38	1459/1459
1R- $\alpha$ -pinene <sup>b</sup>	9	1022/1022	1-octanol	44	1575/1577
$\alpha$ -thujene	11	1030/1027	2,4-bis(1,1-dimethylethyl)-phenol	62	2321/2321
$\beta$ -phellandrene	17	1183/1118	<b>Acids</b>		
3-carene	19	1146/1146	acetic acid	39	1468/1467
D-limonene <sup>a,b</sup>	23	1200/1203	nonanoic acid <sup>b</sup>	61	2184/2184
$\gamma$ -terpinen	27	1249/1249	dodecanoic acid	63	2503/2495
o-cymene <sup>b</sup>	28	1276/1273	<b>Esters</b>		
$\alpha$ -terpinolene <sup>a,b</sup>	29	1285/1286	ethyl acetate	3	895/896
$\beta$ -linalool	42	1563/1562	ethyl propanoate	6	961/961
$\beta$ -caryophyllene <sup>a,b</sup>	46	1603/1602	ethyl 2-methylpropanoate	7	966/969
1-terpinen-4-ol <sup>c</sup>	47	1616/1611	methyl butanoate <sup>b,c</sup>	8	993/992
4,11-selinadiene <sup>a,b</sup>	52	1656/1703	ethyl butanoate	12	1041/1041
$\alpha$ -terpineol	53	1707/1708	ethyl 2-methylbutanoate	13	1057/1056
valencene <sup>a,b</sup>	54	1728/1736	diethyl carbonate	16	1102/1114
(S)-(+)-carvone <sup>a,b</sup>	55	1748/1752	ethyl (E)-2-butenoate <sup>a,b</sup>	20	1164/1168
nerol <sup>a,c</sup>	56	1817/1819	methyl hexanoate	22	1192/1194
(Z)-carveol <sup>b</sup>	57	1856/1857	ethyl hexanoate	26	1239/1239
(E)-carveol	58	1882/1888	2-hexenyl acetate	31	1323/1340
$\beta$ -ionone <sup>a,b</sup>	59	1967/1969	ethyl 2-hexenoate	33	1357/1351
nootkatone <sup>b</sup>	65	2563/2590	ethyl octanoate <sup>a,b</sup>	37	1442/1442
<b>Aldehydes</b>			methyl benzoate	48	1636/1636
hexanal <sup>a,b</sup>	14	1086/1086	butyrolactone	49	1643/1647
heptanal <sup>a,b</sup>	21	1191/1191	ethyl 3-hydroxyhexanoate	51	1673/1693
(E)-2-hexenal	25	1222/1222	<b>Heterocyclic compounds</b>		
benzaldehyde <sup>a,c</sup>	41	1535/1535	furfural	40	1475/1475
3-methyl-benzaldehyde	50	1624/1663	5-methyl-2-furfural	45	1587/1586



**Ketones**

			2,5-furandicarboxaldehyde	60	2022/2017
			5-(hydroxymethyl)-2-furfural	64	2530/2543
1-penten-3-one <sup>a</sup>	10	1025/1025			
4-heptanone	18	1162/1128			
4-methyl-2-heptanone <sup>a</sup>	24	1206/1211			
6-methyl-5-hepten-2-one <sup>a</sup>	32	1343/1343			

---

<sup>a</sup>: compound selected by *post hoc* Tukey's HSD test and discriminating between Italian and South African oranges.

<sup>b</sup>: compound selected by *post hoc* Tukey's HSD test and discriminating between Italian and Spanish oranges.

<sup>c</sup>: compound selected by *post hoc* Tukey's HSD test and discriminating between South African and Spanish oranges.

<sup>d</sup>: LRI<sub>lit</sub>: Linear Retention Index reported in literature by [www.pherobase.com](http://www.pherobase.com), [www.flavornet.org](http://www.flavornet.org), [www.chemspider.com](http://www.chemspider.com) and [www.nist.gov](http://www.nist.gov); LRI<sub>sp</sub>: Linear Retention Index calculated against n-alkanes (C<sub>5</sub>-C<sub>14</sub> and C<sub>8</sub>-C<sub>40</sub>) on VF-WAXms column .

<sup>e</sup> not available.

Figure 1

