

Document downloaded from:

<http://hdl.handle.net/10251/60344>

This paper must be cited as:

Vitale, R.; Bevilacqua, M.; Bucci, R.; Magrì, A.; Magri, A.; Marini, F. (2013). A rapid and non-invasive method for authenticating the origin of pistachio samples by NIR spectroscopy and chemometrics. *Chemometrics and Intelligent Laboratory Systems*. 121:90-99.
doi:10.1016/j.chemolab.2012.11.019.



The final publication is available at

<http://dx.doi.org/10.1016/j.chemolab.2012.11.019>

Copyright Elsevier

Additional Information

1 **A RAPID AND NON-INVASIVE METHOD FOR AUTHENTICATING THE ORIGIN OF**
2 **PISTACHIO SAMPLES BY NIR SPECTROSCOPY AND CHEMOMETRICS**

3 Raffaele Vitale, Marta Bevilacqua, Remo Bucci, Andrea D. Magri, Antonio L. Magri, Federico
4 Marini*

5 *Dept. Chemistry, University of Rome "La Sapienza", Rome, Italy.*

6

7

8

9

***Corresponding author:**

10 dr. Federico Marini

11 Dept. of Chemistry

12 University of Rome "La Sapienza"

13 P.le Aldo Moro 5

14 I-00185 Rome

15 Italy

16 Tel +39 06 4991 3680

17 Fax +39 06 4457 050

18 e-mail: fmmonet@hotmail.com

19

20 **A RAPID AND NON-INVASIVE METHOD FOR AUTHENTICATING THE ORIGIN OF**
21 **PISTACHIO SAMPLES BY NIR SPECTROSCOPY AND CHEMOMETRICS**

22 Raffaele Vitale, Marta Bevilacqua, Remo Bucci, Andrea D. Magri, Antonio L. Magri, Federico
23 Marini

24 *Dept. Chemistry, University of Rome "La Sapienza", Rome, Italy.*

25
26
27 **Abstract**

28 In this study, near-infrared spectroscopy coupled to chemometrics is used to build an analytical
29 protocol to authenticate the origin of pistachio nuts (*Pistacia vera* L.), a high value-added food
30 product.

31 In particular, 483 samples from six different origins (Sicily, India, Iran, Syria, Turkey and U.S.A.)
32 were analyzed by NIR spectroscopy. Spectra were recorded on half seeds cut longitudinally in
33 reflectance mode. Spectral data were then processed by chemometrics to build classification models
34 by SIMCA and PLS-DA. The discriminant approach resulted in classification accuracies higher
35 than 90% for most of the classes. On the other hand, SIMCA built class-models with high
36 sensitivity and specificities, the only exception being the two categories Turkey and Iran, whose
37 heterogeneity resulted in a poorer specificity (anyway higher than 80%). In particular, the results
38 obtained for the samples coming from Bronte (Sicily), the only PDO pistachio production in Europe
39 – 95.5% non error rate in PLS-DA, 90% sensitivity and 97% specificity in SIMCA, as evaluated on
40 the external test set – are very promising from the viewpoint of the authentication of this product.

41 In general, the results show that the coupling of NIR spectroscopy to chemometric classification
42 techniques can be a valuable tool for tracing the origin of pistachio nuts, providing a reliable
43 authentication in a rapid, relatively cheap and non invasive way.

44
45 **Keywords:** Pistachio (*Pistacia vera* L.) nuts, Near Infrared Spectroscopy (NIR), Classification,
46 Partial Least Squares-Discriminant Analysis (PLS-DA), Soft Independent Modeling of Class
47 Analogies (SIMCA).

48
49
50 **1. Introduction**

51 Pistachio (*Pistacia vera* L.) is a nut having peculiar organoleptic characteristics. It is widely
52 consumed as a raw or toasted snack or ingredient of many desserts, ice cream, cakes, pastry and for
53 the production of some sausages such as mortadella [1]. The genus *Pistacia* L. is a member of the

54 *Anacardiaceae* family and consists of at least 11 species. Among these species, *Pistacia vera* is the
55 only cultivated and economically important one [2]. It can grow in dry and hot areas and under
56 saline conditions [3]. Because of its marked resistance to extreme environmental (pedoclimatic and
57 hydrologic) conditions, it is cultivated in Europe and Asia on soils that are unsuitable for other fruit
58 crops [4].

59 The *Pistacia vera* tree is native to arid zones of central and west Asia [5]. Nowadays, only a few
60 major growing areas exist worldwide: the principal pistachio-producing countries of the world are,
61 in order, Iran, USA (California), Turkey, Syria, but to lesser extent, other countries, such as, Italy
62 and India [2,3,6], cultivate pistachios as well [7].

63 In Italy, only a single variety (Bianca) is grown [4,8] and its cultivation is concentrated mainly in
64 Bronte, an area around the Etna volcano, where the lava and climate allow the production of an
65 intensely green nut with a very aromatic taste that is highly prized on the international markets.
66 Italian production is very poor in comparison to Asian and American ones; however, it is
67 compensated by the very high quality of the final products [5,9]. Moreover, fee rates and national
68 laws on commodities of each producing country vary dramatically [7,10]. So, pistachios variation in
69 quality, food safety (*e.g.*, contamination by aflatoxins), import/export fees, legal implications, and
70 financial concerns makes determining the country of origin for pistachios important to protect the
71 consumers against potential fraud [7,10,11]. Moreover, given the fundamental economic
72 implications of any fraud, not only the consumers but also pistachio producers and traders are
73 moved to discover objective chemical techniques that can confirm food labels identifying
74 geographic indications [7,10]. As a consequence, there is the need to develop analytical procedures,
75 which can provide a reliable authentication of the geographical origin of this product.

76 In this framework, some works have already been published in the literature concerning the
77 possibility of differentiating the geographical origin of samples coming from various producing
78 countries, using different chemical indices and analytical techniques. For example, Dyszel & Pettit
79 used the triglycerol profile determined by HPLC and the areas of some DSC peaks to discriminate
80 nuts coming from California from the Iranian and Turkish ones [11]. Furthermore, Anderson &
81 Smith proposed the use of stable isotope analysis to distinguish pistachio samples from the three
82 main growing areas (USA, Asia and Mediterranean countries) [10,12]. Other researches from
83 different groups hypothesized that the variations in the fatty acid composition of pistachio nuts,
84 determined by various techniques (HPLC, GC, DSC, NMR), could be related to the different
85 geographical origin of the product [1,7,11]. Lastly, differences in the profiles of inorganic anions,
86 organic acids, and in color among pistachios of different origins and varieties have also been
87 reported in the literature [4,12].

88 On the other hand, there is a great number of researches where the problem of assessing the
89 authenticity of a wide range of other food commodities, and in particular the problem of tracing
90 their geographical origin, is tackled and solved by the use of near infrared spectroscopy (NIR)
91 coupled with the application of chemometric classification methods for data processing. The
92 possibility of using this spectroscopic technique to address problems connected to the authentication
93 of foodstuff has attracted extensive attention by scientists due to its being rapid, relatively cheap
94 and non-polluting, characteristics which perfectly fit the concept of “green analytical chemistry”
95 [13]. Moreover, in many cases, the use of NIR spectroscopy allows the operator to analyze samples
96 without the need to perform any previous chemical or physical treatment [14]. In the framework of
97 the authentication of food origin, these methods have already been successfully applied, for
98 example, to discriminate geographical origin of olive oils [15-19], meat [20], cheese [21], and
99 honey [22]. However, to our knowledge this approach has never been tried before to trace the origin
100 of pistachio nut samples.

101 Therefore, the aim of the present work is to investigate the possibility of using NIR spectroscopy
102 coupled to chemometric classification methods to build a rapid, relatively cheap and non-invasive
103 analytical procedure for the authentication of the geographical origin of pistachio samples and, in
104 particular, for the recognition of the PDO samples from Bronte (Italy). To this purpose, an
105 experimental setup allowing the spectroscopic determinations to be carried out directly on the nuts,
106 without any sample pretreatment steps was designed. Furthermore, from a data processing
107 standpoint, both a discriminant and class-modeling chemometric approaches (by means of the
108 algorithms PLS-DA and SIMCA, respectively) were used. Special care was also taken in the choice
109 of the suitable method of signal spectral pretreatment prior to the construction of the models.

110

111 **2. Materials and methods**

112 *2.1 Samples*

113 Pistachio nut samples from 6 different countries – the four main producers (Iran, USA, Turkey, and
114 Syria) and two of the smaller but still relevant ones (India and Italy) – were collected and analyzed.
115 In particular, the Italian samples were all coming from the Protected Designation of Origin (PDO)
116 “Pistacchio Verde di Bronte” (Bronte, Sicily). In all cases, samples were obtained from different
117 sources and suppliers, chosen to be as representative as possible of the different production areas.
118 Pistachio nuts were stored in a refrigerator at 4 °C and protected from light until the day prior to
119 analysis, to prevent any kind of surface modification and photodegradation of their molecular
120 constituents. In total, 483 pistachio samples - 41 from Bronte (Italy), 41 from India, 121 from Iran,
121 40 from Syria, 120 from Turkey, 120 from USA (California) - were analyzed.

122

123 *2.2 Acquisition of NIR spectra*

124 For the acquisition of spectra a Nicolet 6700 FT-NIR instrument (Thermo Scientific Inc., Madison,
125 WI), equipped with a tungsten-halogen source and an InGaAs detector, was used. The signals were
126 recorded between 10000 and 4000 cm^{-1} , collecting 82 scans at a nominal resolution of 4 cm^{-1} . All
127 the spectra were acquired at room temperature in the interval on individual pistachio nuts, without
128 any further sample treatment, in reflectance mode, through the use of an integrating sphere (Thermo
129 Scientific Inc., Madison, WI). Operationally, each nut was split in two by a longitudinal cut, so that
130 a flat surface was produced on both halves. Two different NIR spectra were then recorded on each
131 half nut, aligning the pistachio first parallel and then perpendicular to the axis of the optical slit of
132 the integrating sphere, and the four spectra corresponding to a single sample (two halves at two
133 orientations) were averaged prior to the successive elaboration. The data were then exported from
134 Omnicare Suite software (Thermo Fisher Scientific Inc., Waltham, MA) as ASCII files, which were
135 then imported into MATLAB (release R2011a, The MathWorks Inc., Natick, MA), for the
136 successive chemometric analysis. In the data analytical stage, 7 different signal pre-processing
137 techniques were evaluated and compared: MSC (Multiplicative Scatter Correction) [23,24],
138 detrending [25], first and second derivatives, computed according to the Savitzky-Golay method (15
139 points window and third-degree interpolating polynomial) [26], and the combinations of MSC with
140 each of the other three; the possibility of no pretreatment was taken into account, too.

141

142 *2.3 Statistical data analysis*

143 Since the aim of this study is to develop a method to predict the geographical origin of pistachio
144 nuts and, in particular, to build a traceability model for the PDO “Pistacchio verde di Bronte”, the
145 measured NIR data were processed by statistical pattern recognition techniques. In particular, two
146 different techniques were chosen, PLS-DA [27,28] and SIMCA [29,30], as examples of
147 discriminant and class-modeling approaches, respectively. Discriminant techniques focus on the
148 differences between samples coming from different classes and operate by dividing the hyperspace
149 of the variables in as many regions as the number of available categories, while class-modeling
150 techniques are rather focused on the similarities among samples of the same class than on the
151 differences among the classes and act by defining the category space of one class at a time.

152

153 *2.3.1 Partial least squares-discriminant analysis (PLS-DA) [27,28]*

154 Building a classification model can be viewed as finding the best relationship between a
155 multivariate independent matrix \mathbf{X} , whose i^{th} row contains the spectral fingerprint recorded on the

156 i^{th} sample, and a qualitative vector of responses. Accordingly, if a suitably designed dummy
 157 response matrix \mathbf{Y} is introduced, traditional regression methods can be used also to tackle with
 158 classification problems. In particular, when dealing with a classification problem involving m
 159 classes, each training sample is associated with a dummy binary-coded m -dimensional \mathbf{y} vector
 160 having all entries equal to zero except for the component corresponding to the class the sample
 161 belongs to, which is equal to 1. For instance, in a problem involving 6 classes, like the one
 162 considered in the present study, samples belonging to the first category will be described by the
 163 dependent vector [1 0 0 0 0], samples belonging to the second by the vector [0 1 0 0 0] and so
 164 on. Under these assumptions, it is possible to use traditional regression methods to operate
 165 classification, computing a calibration model relating the matrix of predictors and this dummy
 166 matrix of responses. As the name itself suggests, the core of the PLS-DA approach is the use of
 167 Partial Least Squares regression [31], which operates a bilinear decomposition of both the X- and
 168 Y-spaces, under the assumption that a relationship between the two internal spaces exists, to
 169 compute the model parameters. The result is a linear classifier that has proved to be statistically
 170 equivalent to Linear Discriminant Analysis (LDA, [32]), but that is also applicable to all the cases
 171 when LDA cannot be used (low number of samples with respect to variables and/or correlated
 172 indices) [27].

173 In order to interpret the results in terms of the most significant spectral regions, it is important to
 174 check which of the measured variables contribute the most to the definition of the model. In the
 175 case of PLS-based techniques, this kind of information can be summarized in an index called
 176 Variable Importance in Projection (VIP [33]), a value that expresses whether a predictor is
 177 significant in the definition of the F latent vectors model for the prediction of a particular response.
 178 Mathematically, it is defined according to the formula:

$$VIP_j = \sqrt{N_{vars} \frac{\sum_{k=1}^F (b_k^2 \mathbf{t}_k^T \mathbf{t}) (w_{jk} / \|\mathbf{w}_k\|)^2}{\sum_{k=1}^F (b_k^2 \mathbf{t}_k^T \mathbf{t})}} \quad (1)$$

182 where \mathbf{t}_k is the vector of sample scores along the k^{th} latent variable, \mathbf{b}_k is the coefficient of the k^{th}
 183 PLS inner relationship, N_{vars} is the number of experimental variables and w_{jk} and \mathbf{w}_k are the weight
 184 of the j^{th} variable for the k^{th} LV and the weight vector for the k^{th} LV, respectively. Since the average
 185 of squared VIP scores equals 1, ‘greater than one rule’ is generally used as a criterion to identify the
 186 most significant variables. Interpretation of the results can be further improved by inspection of the
 187 regression coefficients of the PLS model which, if opportunely examined, can indicate whether the

188 values of the different variables measured for samples coming from a specified category are higher
189 or lower than those recorded on samples from all the other classes.

190

191

192 2.3.2 Soft Independent Modeling of Class Analogies (SIMCA) [29,30]

193 As stated above, the core of the class modeling approach is that each category is modeled
194 independently on the others. In particular, SIMCA describes each class on the basis of a principal
195 component model of opportune dimensionality, according to the equation:

196

$$197 \mathbf{X}_i = \mathbf{T}_A \mathbf{P}_A^T + \mathbf{E} \quad (2)$$

198

199 where \mathbf{X}_i is the sub-matrix of the original data set obtained by selecting only the samples from the
200 i^{th} class, \mathbf{T}_A and \mathbf{P}_A are the matrices containing the first A scores and loading vectors, respectively,
201 and \mathbf{E} is the matrix of the residuals. Once the principal component model is computed, the class
202 space is defined according to some statistically defined criterion for outlier detection. In particular,
203 two statistical variables are used to express the degree of outlyingness of a sample with respect to
204 the computed principal component model: T^2 which accounts for the distance of the sample within
205 the model space and Q which represents its distance from the model space. The values of these two
206 statistics for the analyzed samples are estimated from the scores matrix \mathbf{T}_A and the residual matrix
207 \mathbf{E} , respectively. The distance between each sample and the model of a category is then computed as
208 “reduced distance” according to the equation:

$$209 d_i = \sqrt{(T_{red,i}^2)^2 + (Q_{red,i})^2} = \sqrt{\left(\frac{T_i^2}{T_{lim}^2}\right)^2 + \left(\frac{Q_i}{Q_{lim}}\right)^2} \quad (3)$$

210 where T_{lim}^2 and Q_{lim} are threshold values for the two statistics corresponding to a selected percentile
211 of the distributions, usually 95%, under the null hypothesis. Commonly, if the reduced distance of a
212 sample exceeds $\sqrt{2}$, the sample is considered as an outlier and rejected by the class model;
213 otherwise, if the distance is lower than this value, it is accepted and recognized as being part of that
214 class.

215

216 3. Results and discussion

217 As anticipated in the introductory section, aim of this study was to build and validate reliable
218 classification models for the traceability of pistachio nuts, coming from 6 different countries,

219 coupling near-infrared spectroscopy and chemometrics. To this purpose, a set of 483 pistachio
220 samples from the 6 investigated countries was collected and analyzed by NIR spectroscopy, as
221 described in section 2.2; the corresponding spectra (after averaging the 4 signals measured for each
222 sample) are reported in Figure 1.

223 To relate the spectral fingerprints to the origin of the samples, discriminant (PLS-DA) and modeling
224 (SIMCA) classification approaches were used and compared. However, since several instrumental
225 effects can hinder or worsen the performances of the classification models, different kinds of
226 spectral pretreatment were tested. As mentioned before, the chosen pretreatments were
227 Multiplicative Scatter Correction (MSC), detrending, first and second derivatives, and the
228 combination of MSC with any of the latter three (with MSC being applied first as recommended by
229 Rinnan *et al.* [34]). Models built after these preprocessings were also compared with the ones
230 calculated from raw spectra.

231 As, when dealing with supervised methods, validation of the models on an independent test set is of
232 paramount importance to unbiasedly assess their predictive ability and performances, the whole
233 data set made of 483 samples was then divided into training and test sets (the former to build the
234 models, the latter to validate them). In order to maintain the same diversity in both sets, a sample
235 splitting scheme based on the Duplex algorithm [35] was adopted. Indeed, the Duplex
236 algorithm starts selecting the two objects in the data matrix that are farthest away from each other
237 according to their Euclidean distance and putting them into the training set. Then, among the
238 remaining candidates, the two objects farthest from each other are put into the test subset. At the
239 next step, consecutive objects are selected and put alternatively in the training and test sets, the
240 object added being the one farthest away from all the objects of the data matrix already selected in
241 the considered set. To determine which object is the farthest one, a so-called maximin criterion,
242 which is the same as in the Kennard and Stone algorithm [36], is used: the Euclidean distance
243 between each candidate object and its closest neighbor already in the considered subset is computed
244 and the object for which this distance is maximal is added.

245 In order to ensure that each of the 6 classes was adequately represented, the selection was
246 performed separately for each category. Moreover, to account for the fact that different
247 pretreatments had to be tested and that as much as possible of the variation after scatter or baseline
248 removal was covered in the selection, at the same time having a unique sample splitting scheme to
249 be able to compare the outcomes after the different preprocessings, a procedure recently designed
250 by our group for another study [19] was adopted. In detail, Duplex algorithm was applied separately
251 on each of the 8 data matrices corresponding to the different pretreatments (7 preprocessings and
252 the raw spectra). The selection was performed class-wise using a 2:1 training:test splitting ratio and

253 working on the principal component representation of the data matrices (considering 20 PCs per
254 category). Accordingly, the frequency of selection of each sample as part of the test set was
255 computed so that, eventually, all the individuals selected more than 50% of the times (i.e. at least 5
256 times out of 8) were included in the final test set (made, in total, of 163 samples: 19 from Bronte, 18
257 from India, 38 from Iran, 15 from Syria, 34 from Turkey, 39 from USA). All the remaining 320
258 samples constituted the training set (22 from Bronte, 23 from India, 83 from Iran, 25 from Syria, 86
259 from Turkey, 81 from USA). The effectiveness of the splitting procedure in keeping a comparable
260 diversity among the two sets can be graphically evaluated in Figure 2, where the projection of the
261 training and the test samples onto the space spanned by the first two principal components is shown
262 for the different spectral preprocessing. It is evident that the proposed procedure allows to select test
263 samples spanning the same space of the training objects, irrespectively of the pretreatment
264 considered, so that the chosen splitting scheme can be used to properly compare the results obtained
265 with the different approaches.

266

267

268 *3.1 PLS-DA analysis*

269 In a first stage of our study, classification models were built according to a discriminant approach
270 using the PLS-DA algorithm. In particular, each of the 8 data sets corresponding to the different
271 spectral pretreatments was processed individually, after mean centering. In each case, the optimal
272 complexity of the classification models was chosen as the one which led to the minimum overall
273 classification error in cross-validation (10 cancellation groups). The results are summarized in Table
274 1, where the correct classification rates for each of the 6 classes and the overall one in calibration
275 and cross-validation are reported for the different spectral preprocessing. It can be observed from
276 the Table that MSC followed by detrending is the pretreatment leading to the best results in cross-
277 validation, the corresponding PLS-DA model (built using 18 LVs) resulting in a classification error
278 of 4.48% in calibration and 5.77% in cross-validation. Moreover, investigation on the classification
279 accuracy for the different classes shows that the model built on data pretreated by MSC and
280 detrending results in non error rates in cross-validation higher than 93% with the only exception of
281 Iran, for which it is anyway slightly less than 90%.

282 This optimal model was then validated on the external test set, and the results are reported in Figure
283 3, where the values of the 6 components of the predicted y vector are reported for each sample.
284 Assigination of an unknown sample to one of the 6 investigated categories is made based on the
285 values of this predicted response vector: the sample is predicted to belong to the category
286 corresponding to the highest value of the component. For the sake of easier visualization, horizontal

287 lines were added to the Figure to indicate the y values above which a sample was assigned to the
288 particular category. It can be seen that the optimal PLS-DA model built on the training set was able
289 to correctly predict the country of origin of most of the validation samples, thus confirming the
290 effectiveness of the approach. The validation results are also reported in Table 2 in terms of non
291 error rate in prediction both on individual categories and on the whole test set. These results
292 indicate that the class belonging of unknown samples can be predicted very accurately (prediction
293 error is lower than 10% for all categories with the only exception of Iran, analogously to what
294 already observed in cross-validation).

295 The goodness of what resulted to be the best PLS-DA classification model, i.e. the one built on
296 spectral data pretreated with MSC and detrending, can also be graphically visualized in Figure 4,
297 where the projection of the training (cross-validated scores) and test samples onto the space spanned
298 by the first three PLS-DA latent variables is displayed. It is apparent from the Figure that the test
299 samples lie well into the space spanned by the training ones. Moreover, it is also possible to observe
300 relatively well the grouping of samples from the different categories even if, due to the complexity
301 of the model (18 LVs), it is difficult to appreciate the separation between the classes on this three
302 dimensional representation.

303 In order to identify the spectral frequencies which contribute the most to the discriminant model,
304 VIP scores [33] were computed and examined: as already described in section 2.3.1, VIP is an
305 index accounting for the contribution of individual experimental variables to the bilinear model and
306 it is scaled in such a way that indices having VIP larger than 1 are considered to be significant.
307 Information on the VIP score was integrated with the examination of the regression coefficient for
308 the interpretation of the model. Indeed, even if the presence of nonorthogonal contributions to the
309 signal can perturb the shape of the regression coefficient vector so that it no longer looks like the
310 pure spectrum [37], inspection of its values, although not straightforward, still can provide useful
311 information. Accordingly, for the sake of interpretation variable significance estimated by VIP
312 scores and the values of the regression coefficients were graphically represented in Figure 5,
313 superimposed to the average spectral profile recorded on the samples after the optimal
314 preprocessing. It can be seen from the Figure that the features identified as relevant by the model
315 (which are also summarized in Table 3) correspond to spectrally meaningful frequencies, and that
316 while most of the spectral intervals are common to all the investigated categories, there are also
317 some significant differences. In particular, the spectral regions that appear to be relevant for all the
318 categories involve the peaks at around $4500\text{--}5000\text{ cm}^{-1}$, which may be attributed to combination
319 bands of C-C and C-H stretching vibration, the signals between 5650 and 6000 cm^{-1} , due to the
320 combination bands and first overtone of C-H bonds, and those between 7074 and 7180 cm^{-1} , which

321 can be ascribed to C-H bonds combination band. On the other hand, portions of the band between
322 8000 and 9000 cm^{-1} (second overtone of methylenic stretching vibrations) are significant only in the
323 definition of the categories Turkey and USA.

324

325 3.2 SIMCA analysis

326 In a second stage, class-modeling approach was also used to process the same data set. Indeed, even
327 if very satisfactory results were obtained using PLS-DA, as described in the previous sub-section,
328 the asymmetry in the number of available samples for each category together with the difference in
329 heterogeneity of the classes due, for instance, to the uneven geographical distribution of the
330 productive areas in the investigated countries (production in Bronte being limited to a very narrow
331 area compared to e.g. Iran, where numerous cultivation sites scattered over the country exist), could
332 be the cause of the non perfect classification rates obtained. Moreover, in problems like those
333 concerning the traceability of foodstuff, where the question to be answered is “Is the product
334 coming from country X as declared?”, one is more interested in assessing whether the investigated
335 sample is compatible with the model of a specific category, which is exactly what class-modeling
336 does. In this investigation, class-modeling on the data set containing the spectral fingerprints of
337 pistachio samples was performed by means of the SIMCA algorithm. Accordingly, independent
338 class models were built for each of the 6 investigated categories, whose optimal complexity was
339 chosen as the number of principal components corresponding to the highest geometrical average of
340 sensitivity and specificity in 6-fold row-wise cross-validation. As in the case of PLS-DA, the effect
341 of the spectral pretreatments on the classification ability was evaluated, by comparing the cross-
342 validated results of SIMCA modeling on the 8 matrices corresponding to the different
343 preprocessings considered. Consistently with what already observed for the discriminant approach,
344 also with SIMCA the best spectral preprocessing resulted to be the coupling of MSC with
345 detrending, which allowed to achieve the highest values of sensitivity and specificity for all the
346 categories. Indeed, most of the class models result in a perfect sensitivity and a rather good
347 specificity in calibration while, in the cross-validation phase, sensitivity decreases significantly for
348 many of the considered pretreatments. Besides, irrespectively of the pretreatment, the models for
349 the categories Iran and Turkey have a significantly lower specificity than the others. This could be
350 explained by the fact that the samples coming from these two regions are produced in various areas
351 scattered in the country, far away from one another and characterized by different climates,
352 environmental conditions and latitude values: to take into account this heterogeneity, keeping
353 reasonable values of sensitivity, the class models have to be wide at the expense of specificity. This
354 hypothesis was also supported by the observation that, for each of the two countries (in particular,

355 for Iran), there is a marked difference among the spectral fingerprints of the producing regions
356 represented in the data set. When the models built after the best spectral pretreatment
357 (MSC+detrending) were applied to the external test set, very good results were obtained (Table 4).
358 Indeed, sensitivity and specificity values for the validation samples were over 80% for most of the
359 categories (the only exceptions being the sensitivities of Syria and Turkey models), being in many
360 cases even higher than 90%. These results can be graphically evaluated in Figure 6, where the
361 projections of the training (cross-validated predictions) and test samples onto the model space of
362 each single category are reported. This representation allows to easily visualize which samples are
363 accepted or rejected by the different class models: the dotted lines in the graphs correspond to the
364 threshold values of reduced distance, below which the samples are accepted by the model of the
365 considered category, as described in Section 2.3.2. The Figure shows clearly that all the models are
366 very sensitive both in cross-validation and on the external test set but that some of them,
367 particularly Turkey and Iran, have lower specificity.

368 When more than one category is modeled, it is possible to check whether the samples are accepted
369 by one, more than one or none of them. This could be useful in order to turn SIMCA into a
370 discriminant classifier by assigning each sample to the category it is closer to. This information can
371 be easily visualized building a so-called Coomans plot [38], a graph where the two axes represent
372 the distance of the samples to each of the two class models under study. As an example, **Error!**
373 **Reference source not found.** shows the Coomans plot for the Bronte and USA class models built
374 on the spectral data pretreated by MSC + detrending algorithms. The dotted black lines correspond
375 to the threshold distances values (in our case $\sqrt{2}$) and cut the plot in four different regions: the
376 uppermost left and the lowermost right will correspond to unambiguous acceptance by a single
377 category model (respectively Bronte and USA), the lowermost left to acceptance by both classes
378 while the uppermost right to rejection by both category models. Most of the samples coming from
379 the two different geographical areas lie inside the space of the corresponding class model. Only few
380 samples coming from the other four zones were accepted by the two different class models, while
381 the remaining samples are found in the uppermost right region of the plot. Moreover, none of the
382 samples are accepted by both models. The diagonal line bisecting the plot represents the
383 discriminant classification boundary so that all the samples lying above it are classified as from
384 Bronte, while all the samples lying below it are predicted as from USA. Based on these
385 considerations, discriminant classification based on SIMCA models would result in 100% (Bronte)
386 and 100% (USA) both in cross-validation and on the external test set.

387

388 *3.3 A closer look on Bronte*

389 As reported in the Introduction, one of the 6 classes investigated in this study, “Pistacchio verde di
390 Bronte”, is the only pistachio product with a Protected Designation of Origin. Therefore, from the
391 standpoint of the analytical control of frauds on certified foodstuff it would be of outmost
392 importance to have an accurate method for checking the authenticity of this product. In this respect,
393 the results of this study appear quite promising, both if the discriminant or if the class-modeling
394 approach are concerned. Indeed, by looking at the results of PLS-DA reported in Section 3.1 one
395 could see that for the class Bronte a non-error rate in prediction higher than 95% was obtained.
396 Additionally, by inspecting the plots in Figure 3 it can be observed that only a very low number of
397 samples from other categories are erroneously predicted as belonging to Bronte.
398 As far as SIMCA is concerned, the results in Table 3 indicate that the model of the category Bronte
399 has very high sensitivity and specificity in prediction, as requested for a reliable traceability model.

400

401 **4. Conclusions**

402 In this study, the potential of NIR spectroscopy coupled to chemometric discriminant and class-
403 modeling pattern recognition techniques for the traceability of pistachio nuts samples was
404 demonstrated. Classification models with high accuracy were built to recognize the geographical
405 origin of the samples, as evaluated on an external test set. The outcomes for both the different
406 classification approaches are very satisfying: the origin of over 95% of validation samples was
407 correctly predicted using PLS-DA and these results were confirmed by SIMCA modeling of the
408 same data, which allowed to build very sensitive and highly specific models for authenticating the
409 provenance of pistachio nuts. In particular, the results obtained for the category Bronte, which is the
410 only type of pistachio having a Protected Denomination of Origin, appear really promising in the
411 light of the possibility of building a traceability model for this product.

412 Furthermore, a first attempt of interpreting the observed differences in terms of significant spectral
413 bands by means of the inspection of VIP scores was made: discriminant information was found to
414 be associated to meaningful signals corresponding to methylene overtones and C-H and C=C
415 combination bands.

416 Finally, it can be concluded that NIR spectroscopy coupled to chemometric classification
417 techniques is a powerful tool to trace pistachio nuts samples, allowing a fast, cheap and non-
418 invasive/non-destructive analysis.

419

420 **References**

- 421 [1] E. Arena, S. Campisi, B. Fallico, E. Maccarone, Distribution of fatty acids and phytosterols as a
422 criterion to discriminate geographic origin of pistachio seeds, *Food Chem.* 104 (2007) 403-408.
- 423 [2] S. Kafkas, H. Ozkan, B.E. Ak, I. Akar, H.S. Atli, S. Koyuncu, Detecting DNA polymorphism
424 and genetic diversity in a wide pistachio germplasm: comparison of AFLP, ISSR and RAPD
425 markers, *J. Am. Soc. Hortic. Sci.* 131 (2006) 522-529.
- 426 [3] E. Tsantili, C. Takidell, M. Christopoulos, E. Lambrinea, D. Rouskas, P. Roussos, Physical,
427 compositional and sensory differences in nuts among pistachio (*Pistachia vera* L.) varieties, *Sci.*
428 *Hortic.* 125 (2010) 562–568.
- 429 [4] M. Bellomi, B. Fallico, Anthocyanins, chlorophylls and xanthophylls in pistachio nuts (*Pistacia*
430 *vera*) of different geographic origin, *J. Food Compos. Anal.* 20 (2007) 352-359.
- 431 [5] G. Ballistreri, E. Arena, B. Fallico, Characterization of triacylglycerols in *Pistacia Vera* L. oils
432 from different geographic origins, *Ital. J. Food Sci.* 22 (2010) 69-75.
- 433 [6] Food and Agriculture Organization. Food and Agricultural commodities production, 2011.
434 <http://faostat.fao.org/site/339/default.aspx> Accessed 22/04/12
- 435 [7] K. Zur, A. Heier, K.W. Blaas, C. Fauhl-Hassek, Authenticity control of pistachios based on ¹H-
436 and ¹³C-NMR spectroscopy and multivariate statistics, *Eur. Food Res. Technol.* 227 (2008) 969–
437 977.
- 438 [8] A. Fabbri, C. Valenti, The sicilian pistachio industry: an overview, *Acta Hort.* 470 (1998) 43-49.
- 439 [9] L. Di Marco, Il pistacchio in Sicilia: situazione e prospettive, *Agricoltura Ricerca* 79 (1987) 9.
- 440 [10] K.A. Anderson, B.W. Smith, Effect of Season and variety on the differentiation of geographic
441 growing origin of pistachios by stable isotope profiling, *J. Agric. Food Chem.* 54 (2006) 1747-1752.
- 442 [11] S.M. Dyszel, B.C. Pettit, Determination of the country of origin of pistachio nuts by DSC and
443 HPLC, *J. Am. Oil Chem. Soc.* 67 (1990) 947–951.
- 444 [12] K.A. Anderson, B.W. Smith, Use of chemical profiling to differentiate geographic growing
445 origin of raw pistachios, *J. Agric. Food Chem.* 53 (2005) 410-418.
- 446 [13] S. Armenta, S. Garrigues, M. De La Guardia, Green Analytical Chemistry, *Trends Anal. Chem.*
447 27 (2008) 497-511.
- 448 [14] D.H. Ma, X.C. Wang, L.P. Liu, Y. Liu, Current progress in food geographical origin
449 traceability by near infrared spectroscopy technology, *Guangpuxue Yu Guangpu Fenxi* 31 (2011)
450 877-881.
- 451 [15] M. Casale, C. Casolino, P. Oliveri, M. Forina, The potential of coupling information using
452 three analytical techniques for identifying the geographical origin of Liguria extra virgin olive oil,
453 *Food Chem.* 118 (2009) 163-170.

- 454 [16] T. Woodcock, G. Downey, C.P. O'Donnell, Confirmation of declared provenance of European
455 extra virgin olive oil samples by NIR spectroscopy, *J. Agric. Food Chem.* 56 (2008) 11520-11525.
- 456 [17] N. Sinelli, E. Casiraghi, D. Tura, G. Downey, Characterization and classification of Italian
457 virgin olive oils by near- and mid-infrared spectroscopy, *J. Near Infrared Spectrosc.* 16 (2008) 335-
458 342.
- 459 [18] O. Galtier, N. Dupuy, Y. Le Dreau, D. Ollivier, C. Pinatel, J. Kister, J. Artaud, Geographic
460 origins and compositions of virgin olive oils determined by chemometric analysis of NIR spectra,
461 *Anal. Chim. Acta* 595 (2007) 136-144.
- 462 [19] M. Bevilacqua, R. Bucci, A.D. Magri, A.L. Magri, F. Marini, Tracing the origin of extra virgin
463 olive oils by infrared spectroscopy and chemometrics: A case study, *Anal. Chim. Acta* 717 (2012)
464 39-51.
- 465 [20] S.M. Sun, B.L. Guo, Y.M. Wei, M.T. Fan, Application of near infrared spectral fingerprint
466 technique in lamb meat origin traceability, *Guangpuxue Yu Guangpu Fenxi* 31 (2011) 937-941.
- 467 [21] R. Karoui, A.M. Mouazen, E. Dufour, L. Pillonel, E. Schaller, J. De Baerdemaeker, J.O. Bosset,
468 Chemical characterization of European Emmental cheeses by near infrared spectroscopy using
469 chemometric tools, *Int. Dairy J.* 16 (2006) 1211-1217.
- 470 [22] T. Woodcock, G. Downey, J.D. Kelly, C. O'Donnell, Geographical classification of honey
471 samples by near-infrared spectroscopy: a feasibility study, *J. Agric. Food Chem.* 55 (2007) 9128-
472 9134.
- 473 [23] H. Martens, S.A. Jensen, P. Geladi, Multivariate linearity transformations for near infrared
474 reflectance spectroscopy. In: *Proc. Nordic Symposium on Applied Statistics*, Stokkland Forlag,
475 Stavanger, Norway, 1983, 205–234.
- 476 [24] P. Geladi, D. MacDougall, H. Martens, Linearization and scatter-correction for near-infrared
477 reflectance spectra of meat, *Appl. Spectrosc.* 39 (1985) 491-500.
- 478 [25] R.J. Barnes, M.S. Dhanoa, S. Lister, Standard normal variate transformation and de-trending of
479 near-infrared diffuse reflectance spectra, *Appl. Spectrosc.* 43 (1989) 772-777.
- 480 [26] A. Savitzky, M. Golay, Smoothing and differentiation of data by simplified least squares
481 procedures, *Anal. Chem.* 36 (1964) 1627-1639.
- 482 [27] S. Wold, C. Albano, W.J. Dunn III, K. Esbensen, S. Hellberg, E. Johansson, M. Sjöström,
483 Pattern recognition: finding and using regularities in multivariate data, in: H. Martens, H.
484 Russwurm Jr. (Eds.), *Food research and data analysis*, Elsevier Applied Science, Barking, UK, 2009,
485 pp. 147-188.
- 486 [28] M. Barker, W. Rayens, Partial least squares for discrimination, *J. Chemom.* 17 (2003) 166-173.

- 487 [29] S. Wold, Pattern recognition by means of disjoint principal components models, *Pattern*
488 *Recogn.* 8 (1976) 127-139.
- 489 [30] S. Wold, M. Sjöström, SIMCA: a method for analysing chemical data in terms of similarity
490 and analogy, in: B.R. Kowalski (Ed.), *Chemometrics, theory and application*, 52nd ed., American
491 Chemical Society, Washington, DC, 1977, pp. 243-282.
- 492 [31] S. Wold, H. Martens, H. Wold, *Matrix Pencils: Proceedings of a Conference Held at Pite*
493 *Havsbad*. Ruhe, A., Kagstrom, B., Eds.; Springer-Verlag: Heidelberg, Germany, 1983; pp. 286-293.
- 494 [32] R.A. Fisher, The use of multiple measurements in taxonomic problems, *Ann. Hum. Genet.* 7
495 (1936) 179–188.
- 496 [33] S. Wold, E. Johansson, M. Cocchi, 3D QSAR, in: H. Kubinyi (Ed.), *Drug Design: Theory,*
497 *Methods and Applications*, Escom Science Publishers, Leiden, The Netherlands, 1993, pp. 523–550.
- 498 [34] Å. Rinnan, F. van den Berg, S. Engelsen, Review of the most common pre-processing
499 techniques for near-infrared spectra, *Trends Anal. Chem.* 28 (2009) 1201-1222.
- 500 [35] R. Snee, Validation of regression models: methods and examples, *Technometrics* 19 (1977)
501 415-428.
- 502 [36] R.W. Kennard, L.A. Stone, Computer aided design of experiments, *Technometrics* 11 (1969)
503 137–148.
- 504 [37] M.B. Seasholtz, B. Kowalski, Qualitative information from multivariate calibration models.
505 *Appl. Spectrosc.* 44 (1990) 1337-1348.
- 506 [38] D. Coomans, I. Broeckaert, M.P. Derde, A. Tassin, D.L. Massart, S. Wold, Use of a
507 microcomputer for the definition of multivariate confidence regions in medical diagnosis based on
508 clinical laboratory profiles, *Comput. Biomed. Res.* 17 (1984) 1-14.
- 509
- 510

511 **Figure captions:**

512 Figure 1 – Raw spectra of the 483 pistachio samples analyzed in this study (after averaging the 4
513 signals recorded on each nut).

514 Figure 2 - Representation of the data splitting between the training and test set as a function of the
515 different spectral pretreatments. Data are projected on the space spanned by the first two principal
516 components. (●) Bronte; (■) India; (◆) Iran; (▼) Syria; (▲) Turkey; (★) USA; empty symbols
517 correspond to training samples and filled symbols correspond to the test samples.

518 Figure 3 – PLS-DA analysis after MSC+detrending: predicted values of the dummy vector
519 components corresponding to the different categories for test set samples. Horizontal lines indicate
520 the threshold above which a sample is assigned to that particular class. (●) Bronte; (■) India; (◆)
521 Iran; (▼) Syria; (▲) Turkey; (★) USA.

522 Figure 4 – PLS-DA analysis after MSC+detrending: projection of the training (cross-validated
523 scores) and test samples onto the space spanned by the first three latent variables. (●) Bronte; (■)
524 India; (◆) Iran; (▼) Syria; (▲) Turkey; (★) USA; empty symbols correspond to training samples
525 and filled symbols correspond to the test samples.

526 Figure 5 – PLS-DA analysis after MSC+detrending: graphical representation of the regression
527 vectors for the 6 categories superimposed to the average pretreated spectrum. The regression vector
528 components are colored according to their VIP score: significant variables ($VIP > 1$) are colored in
529 red, while those estimated as not relevant in green).

530 Figure 6: SIMCA on near-infrared data after MSC+detrending: projection of the training (cross-
531 validated estimates) and test samples onto the model spaces of the 6 investigated categories. (●)
532 class Bronte; (■) class India; (◆) class Iran; (▼) class Syria; (▲) class Turkey; (★) class USA.
533 Empty symbols correspond to training samples and filled symbols correspond to the test samples.

534 Figure 7: SIMCA on near-infrared data after MSC+detrending: Coomans plot comparing the
535 models of classes Bronte and USA. (●) class Bronte; (■) class India; (◆) class Iran; (▼) class Syria;
536 (▲) class Turkey; (★) class USA. Empty symbols correspond to training samples and filled
537 symbols correspond to the test samples.

538

Table 1- PLS-DA analysis: comparison of the correct classification rates in calibration and cross-validation obtained after the different spectral pretreatments.

| LV | Bronte | | India | | Iran | | Syria | | Turkey | | USA | | Overall | |
|----|--------|--------|--------|--------|--------|--------|--------|--------|--------|--------|--------|--------|---------|--------|
| | Cal. | CV | Cal. | CV | Cal. | CV | Cal. | CV | Cal. | CV | Cal. | CV | Cal. | CV |
| 9 | 95.04% | 95.04% | 94.61% | 89.59% | 76.38% | 76.59% | 93.63% | 93.46% | 86.09% | 85.3% | 86.66% | 86.25% | 85.53% | 84.90% |
| 7 | 94.37% | 94.37% | 74.12% | 69.43% | 82.95% | 80.45% | 91.63% | 91.29% | 90.8% | 90.8% | 90.22% | 89.80% | 87.73% | 86.61% |
| 7 | 96.14% | 95.97% | 86.73% | 83.71% | 80.84% | 78.76% | 90.78% | 90.61% | 87.74% | 86.36% | 87.72% | 88.13% | 86.69% | 85.64% |
| 15 | 98.49% | 97.82% | 98.82% | 94.31% | 92.98% | 89.91% | 93.97% | 91.46% | 88.96% | 84.99% | 97.71% | 96.67% | 93.97% | 91.28% |
| 14 | 97.32% | 94.87% | 91.78% | 90.94% | 85.39% | 83.97% | 89.12% | 86.95% | 87.47% | 85.45% | 91.28% | 88.58% | 89.01% | 87.02% |
| 18 | 97.48% | 97.32% | 96.82% | 95.3% | 90.54% | 89.48% | 96.47% | 93.97% | 95.17% | 93.15% | 99.79% | 99.17% | 95.52% | 94.23% |
| 15 | 97.82% | 97.15% | 94.14% | 93.13% | 94.46% | 91.81% | 95.63% | 89.29% | 86.95% | 84.71% | 97.93% | 95.43% | 93.62% | 91.08% |
| 16 | 98.66% | 93.61% | 99.66% | 96.82% | 96.93% | 93.64% | 95.80% | 91.63% | 89.18% | 87.22% | 97.92% | 96.05% | 95.32% | 92.59% |

Table 2: PLS-DA analysis on NIR data after MSC + detrending pretreatment: modeling and validation results

| Classes | Non error rate in calibration | Non error rate in cross-validation | Non error rate in prediction |
|----------------|--------------------------------------|---|-------------------------------------|
| Bronte | 97.48% | 97.32% | 95.14% |
| India | 96.82% | 95.30% | 90.29% |
| Iran | 90.54% | 89.48% | 83.59% |
| Syria | 96.47% | 93.97% | 93.63% |
| Turkey | 95.17% | 93.15% | 91.71% |
| USA | 99.79% | 99.17% | 99.19% |

Table 3 – PLS-DA model: features identified as relevant for all the classes

| Spectral Region | Vibrational modes |
|------------------------------|---|
| 4500–5000 cm ⁻¹ | combination bands of C-C and C-H stretching vibration |
| 5650–6000 cm ⁻¹ | combination bands and first overtone of C-H bonds |
| 7074–7180 cm ⁻¹ | C-H bonds combination band |
| 8000–9000 cm ⁻¹ * | second overtone of methylenic stretching vibrations |

*Meaningful only for the categories Turkey and USA

Table 4: Results of SIMCA analysis on NIR data after MSC + detrending pretreatment for both training and test sets

| Classes | LV | Calibration | | Cross-validation | | Prediction | |
|---------|----|-------------|-------------|------------------|-------------|-------------|-------------|
| | | Sensitivity | Specificity | Sensitivity | Specificity | Sensitivity | Specificity |
| Bronte | 4 | 95.45% | 95.64% | 72.73% | 95.30% | 89.47% | 96.53% |
| India | 7 | 100.00% | 90.57% | 65.22% | 93.60% | 83.33% | 98.62% |
| Iran | 5 | 98.80% | 68.35% | 93.98% | 67.93% | 92.11% | 76.80% |
| Syria | 1 | 88.00% | 88.81% | 88.00% | 87.46% | 73.33% | 82.43% |
| Turkey | 13 | 95.35% | 79.06% | 84.88% | 80.77% | 73.53% | 80.62% |
| USA | 10 | 93.83% | 100.00% | 85.19% | 100.00% | 87.18% | 99.19% |

Figure 1

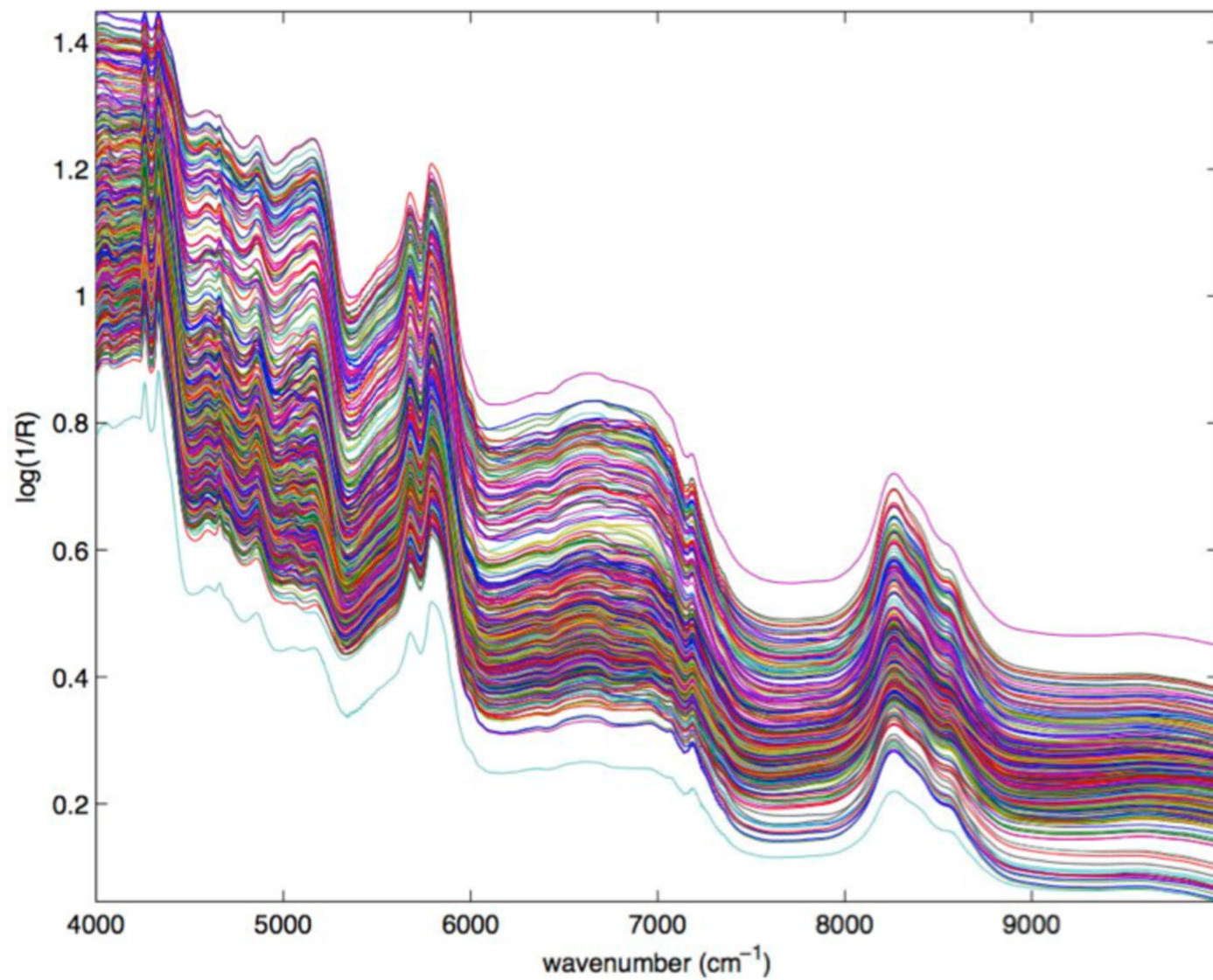


Figure 2a

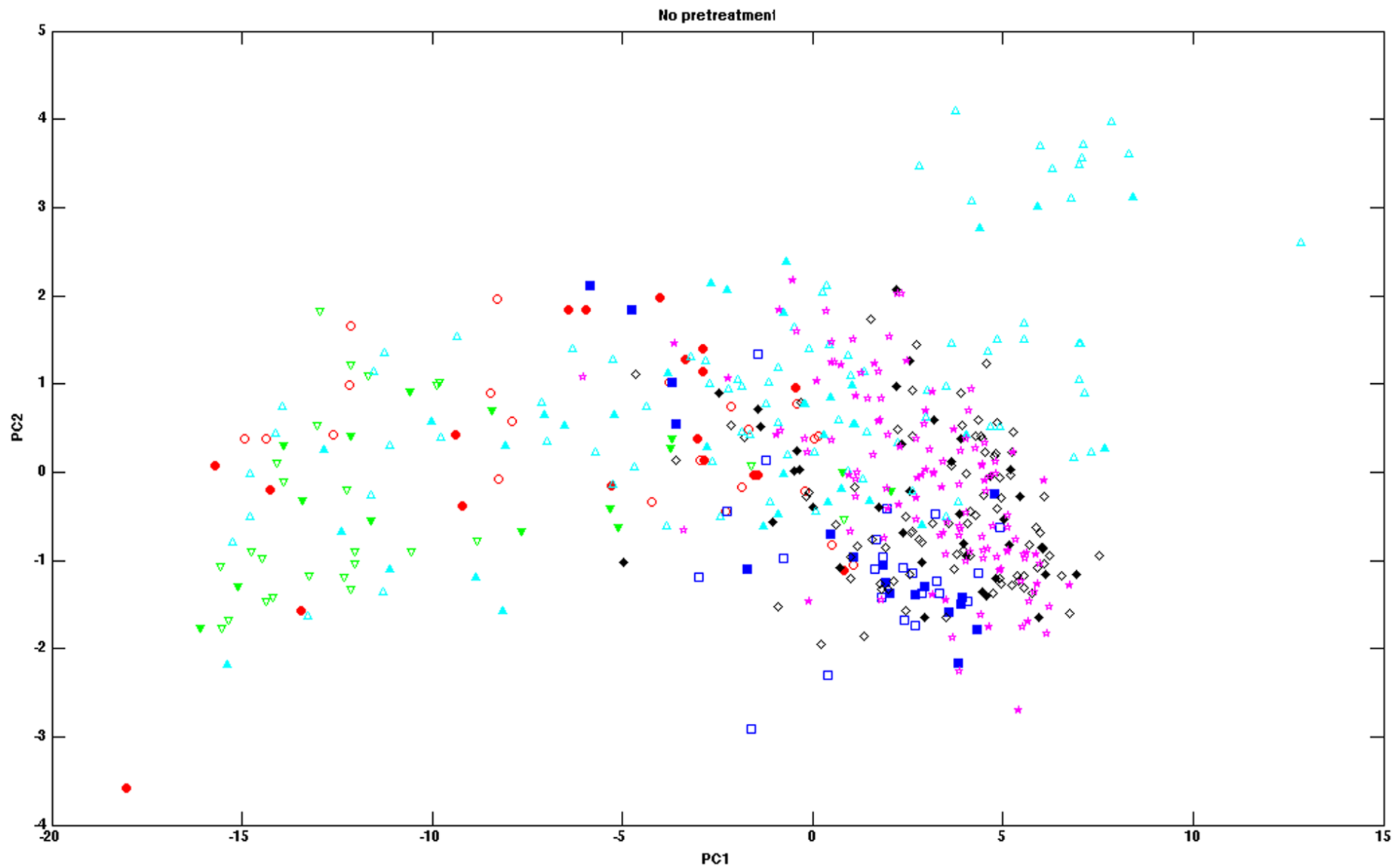


Figure 2b

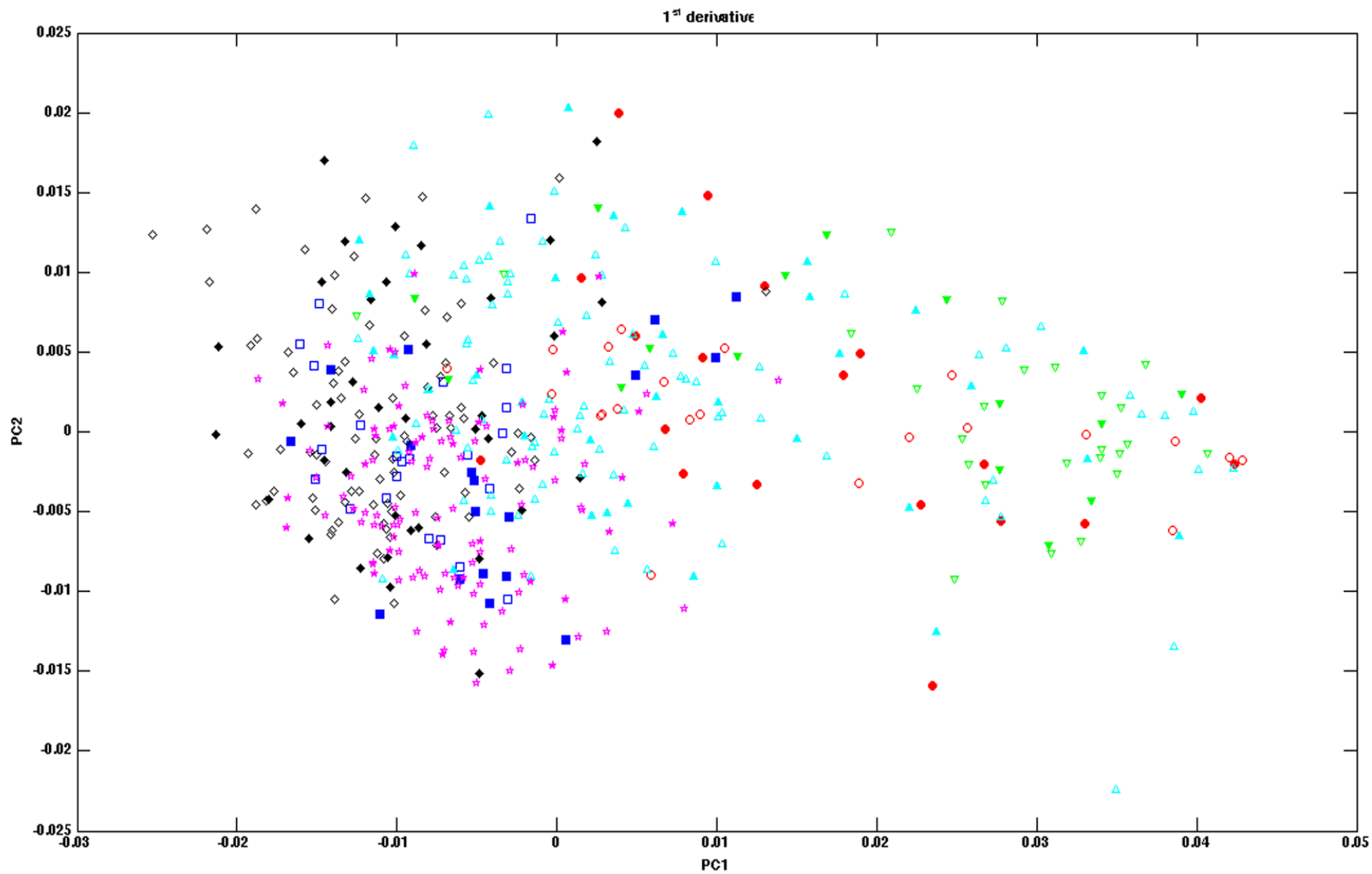


Figure 2c

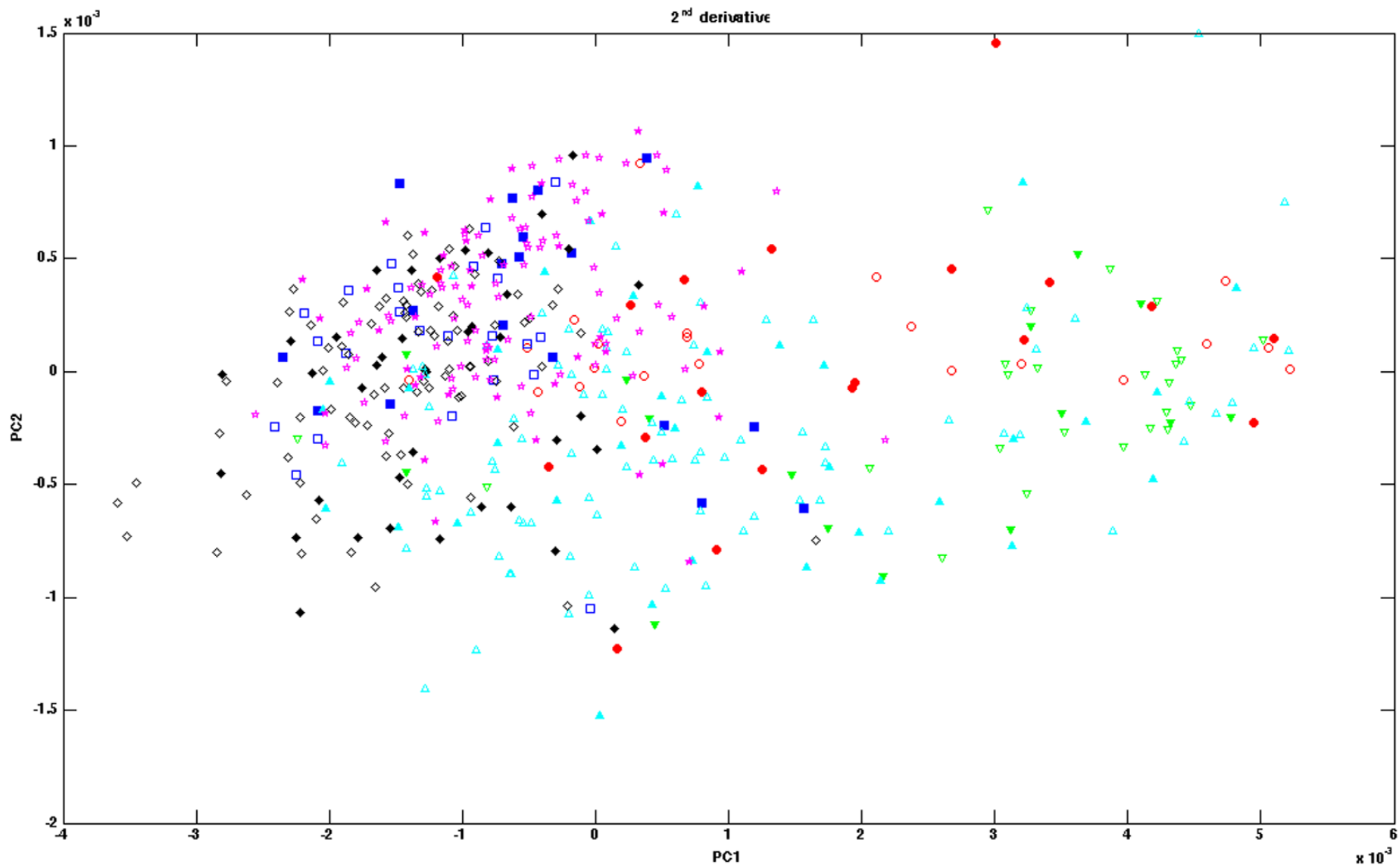


Figure 2d

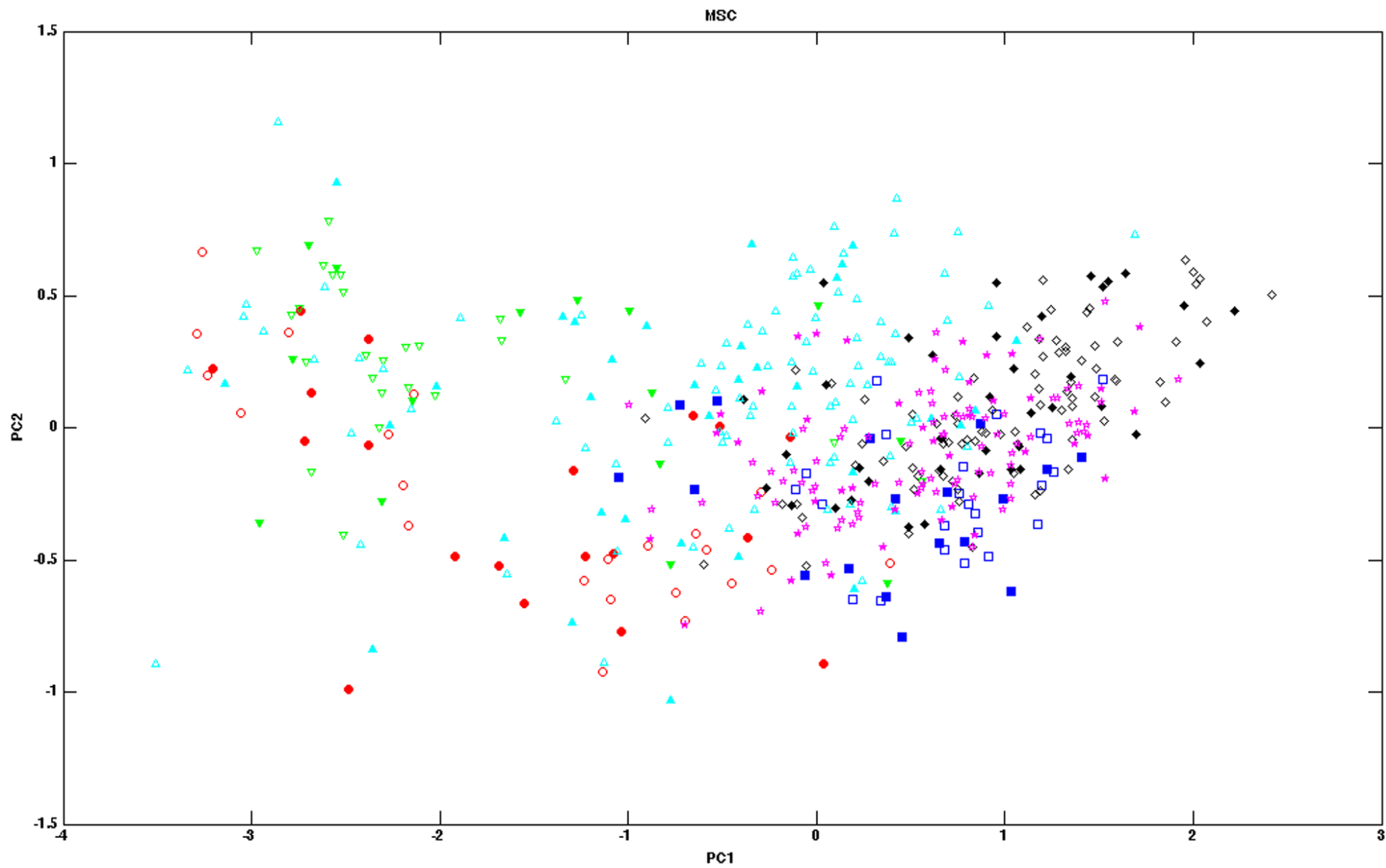


Figure 2e

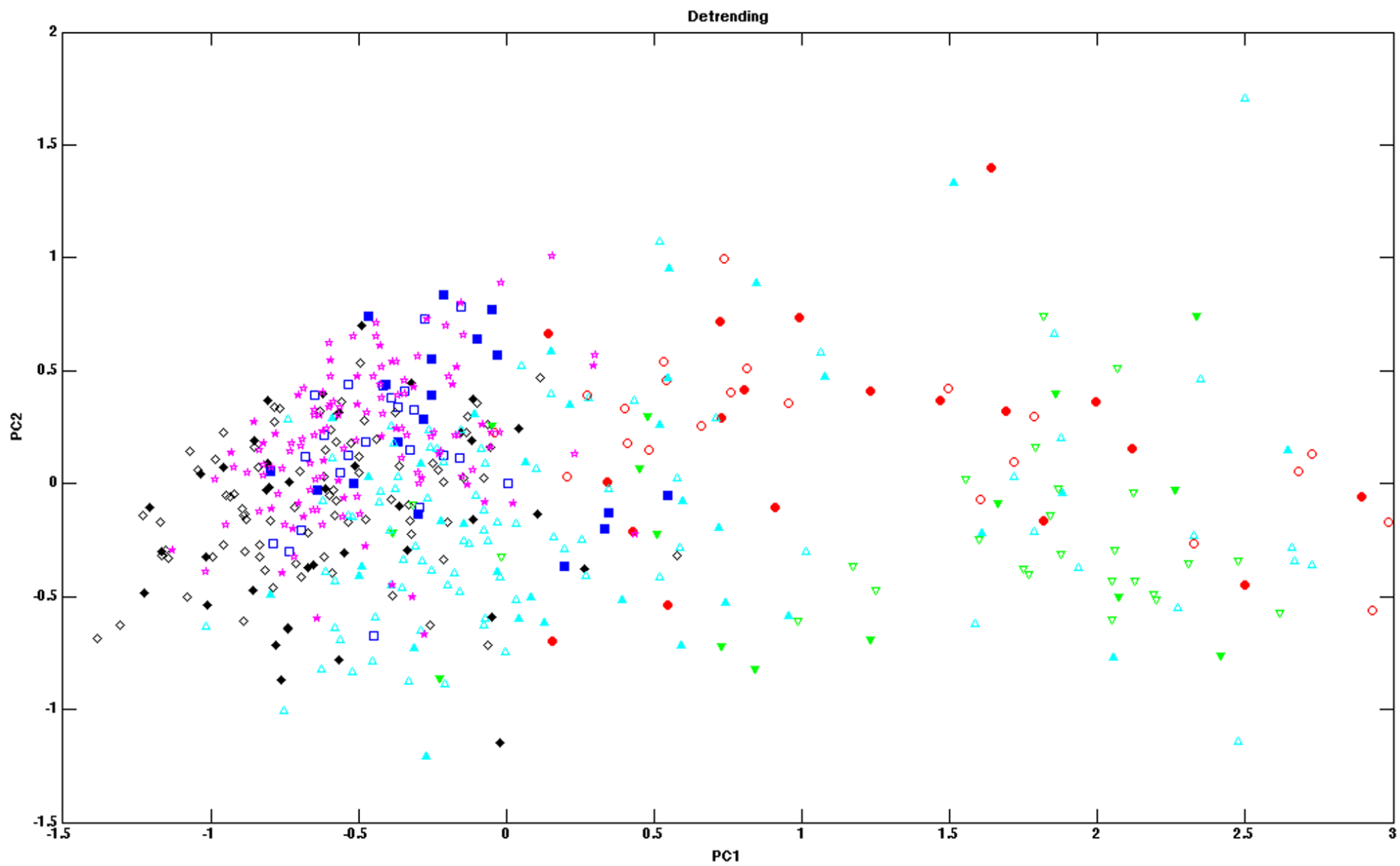


Figure 2f

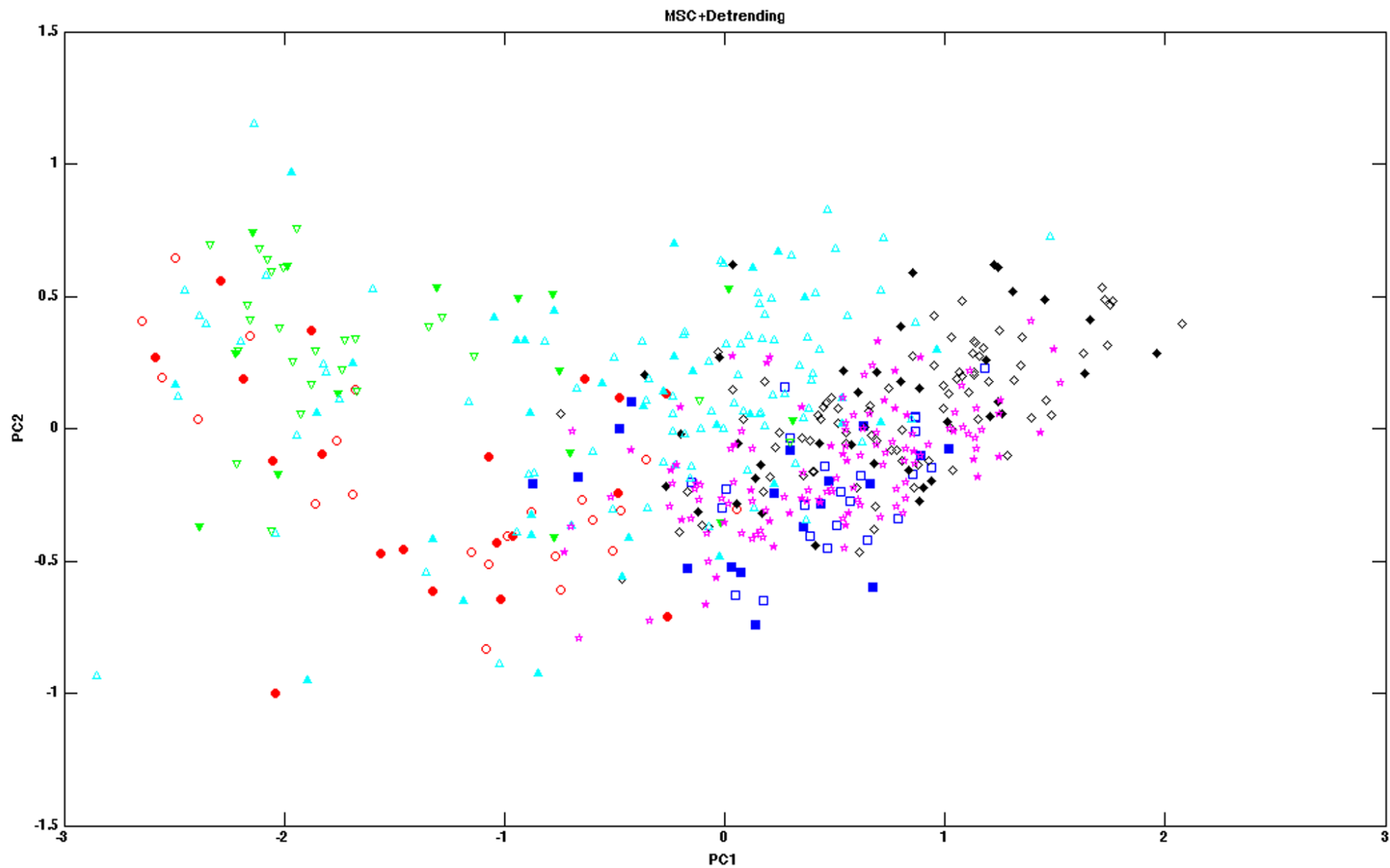


Figure 2g

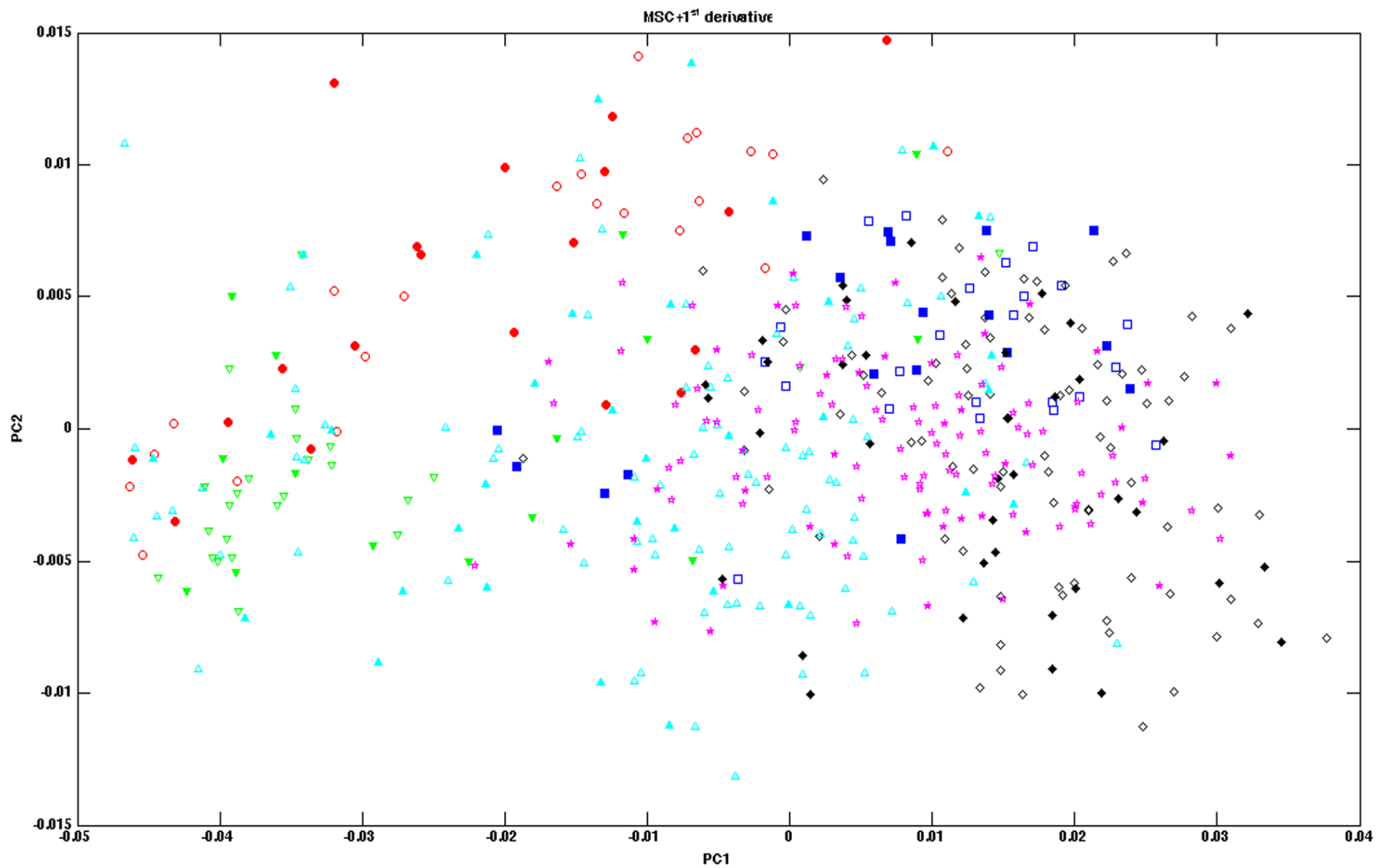


Figure 2h

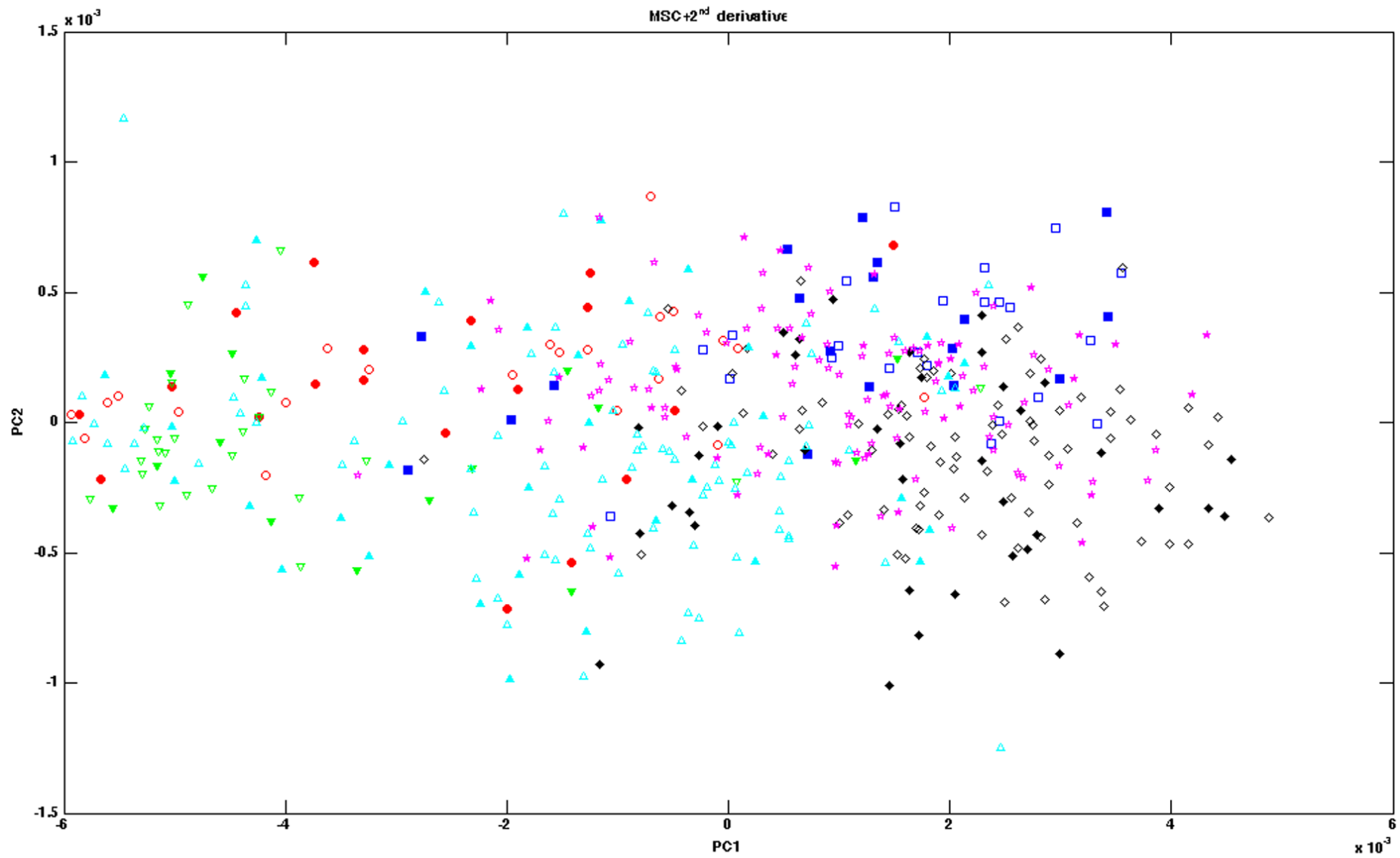


Figure 3

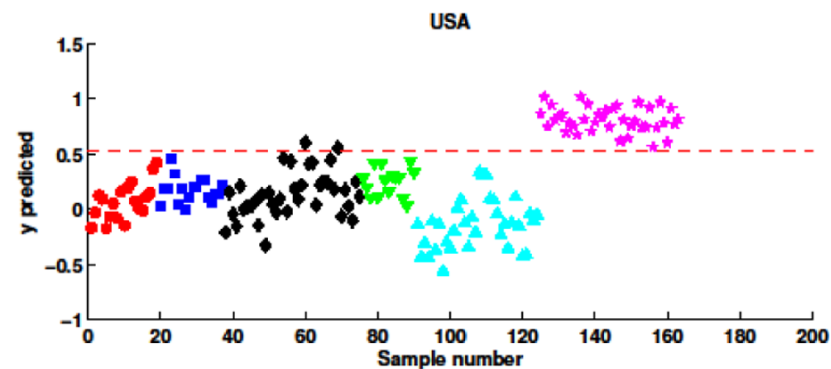
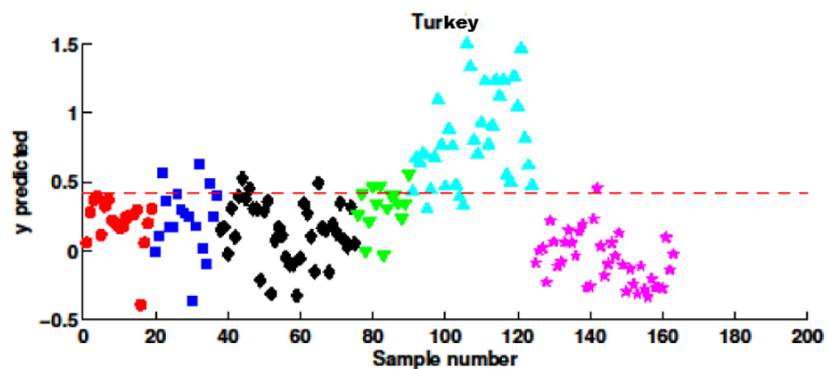
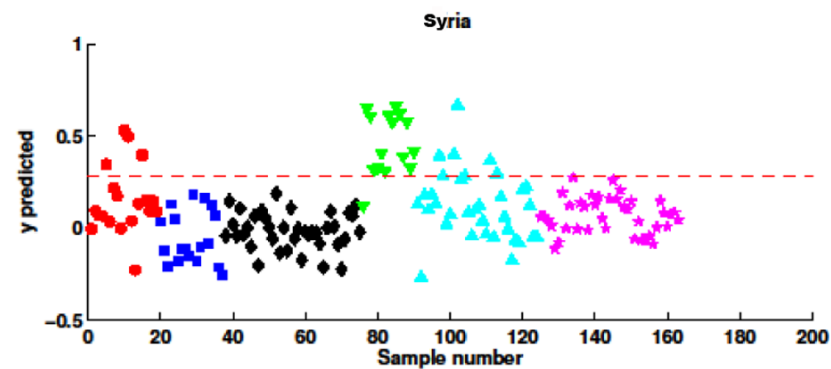
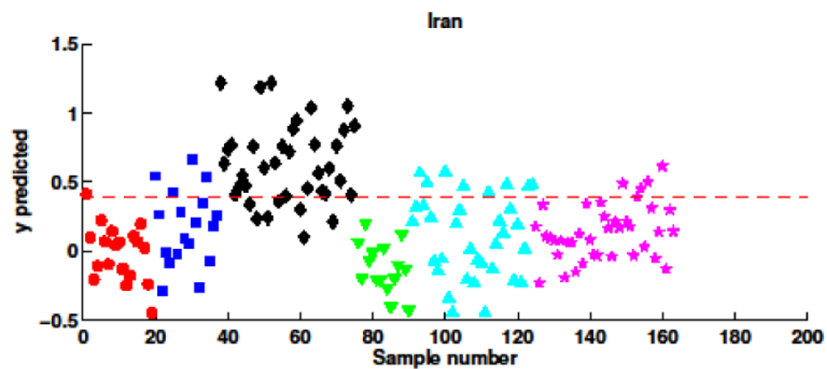
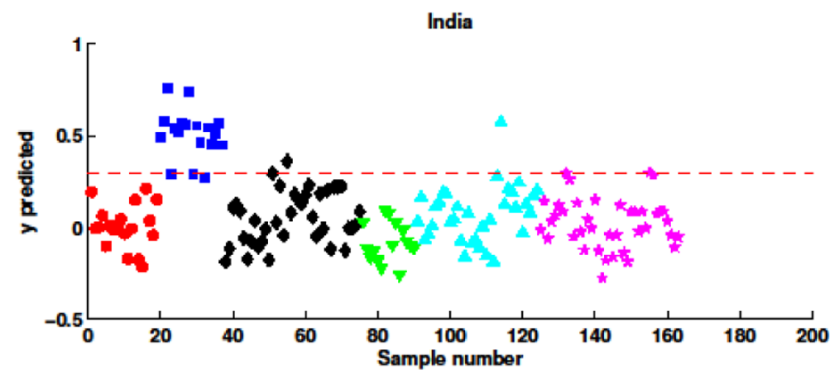
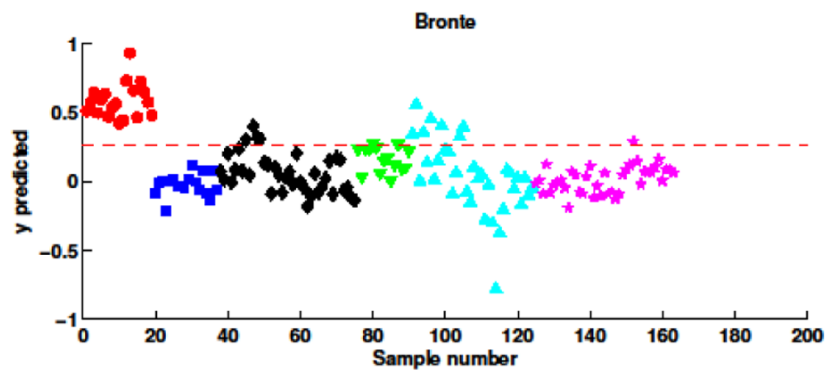


Figure 4

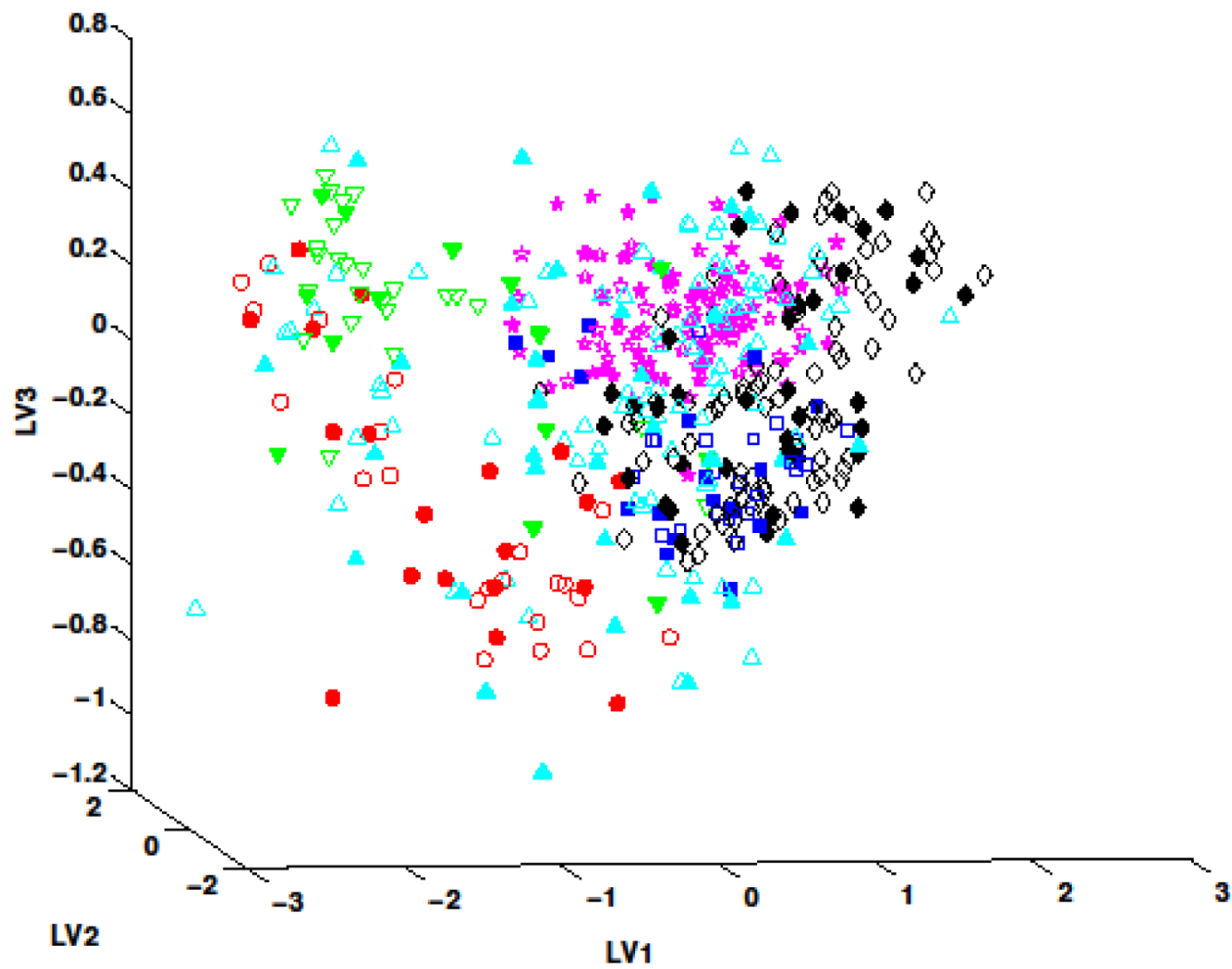


Figure 5

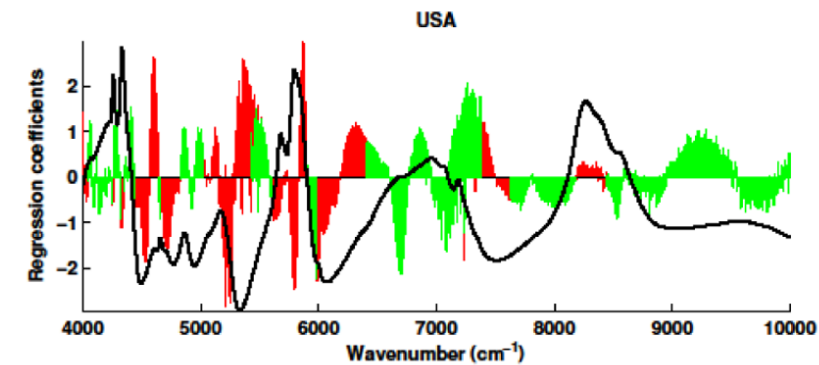
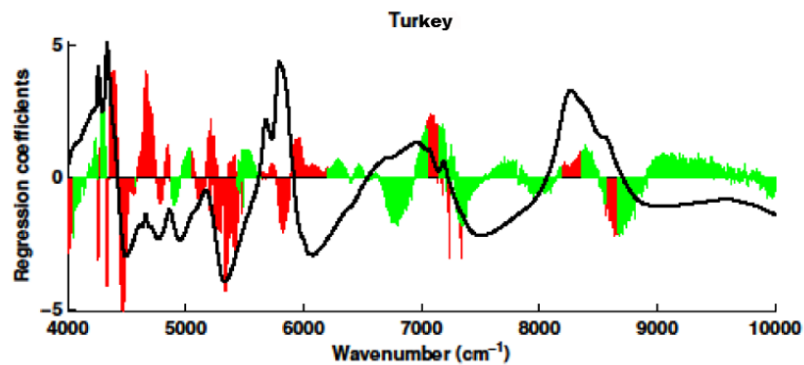
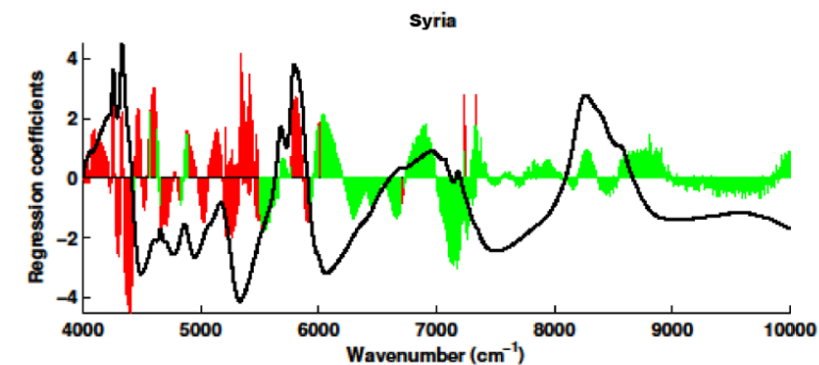
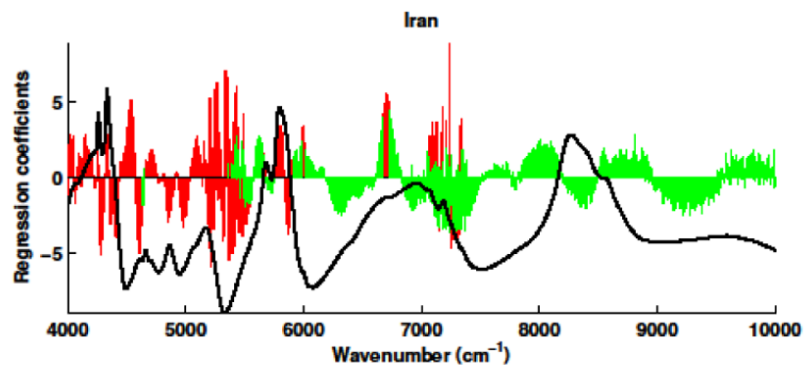
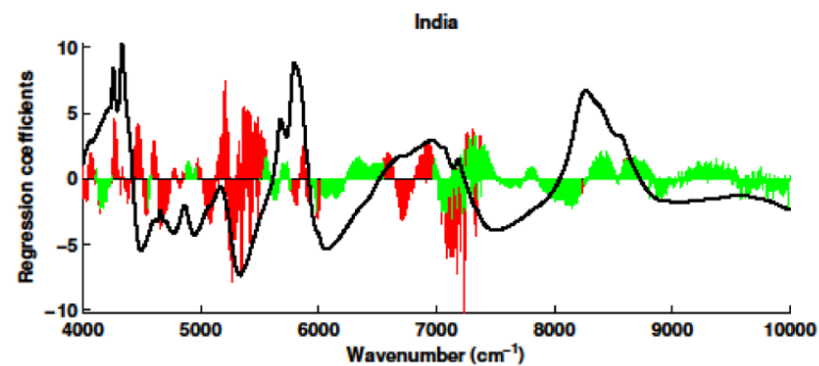
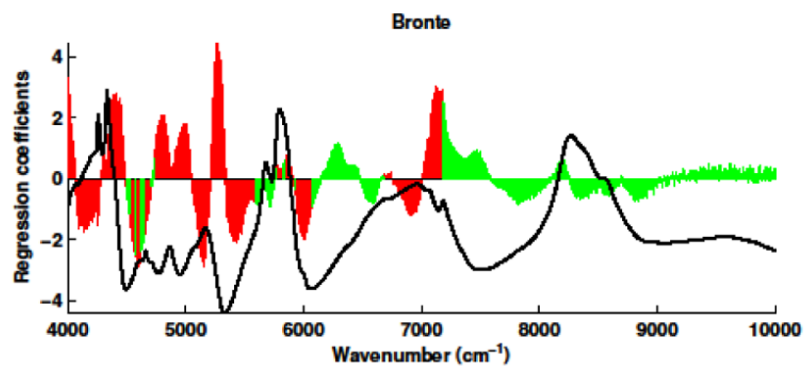


Figure 6

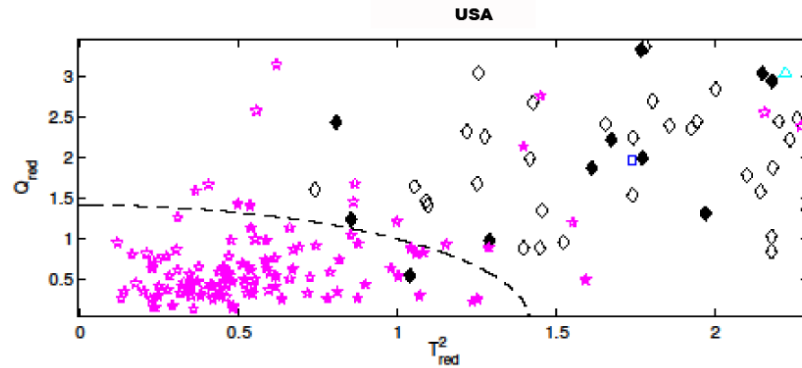
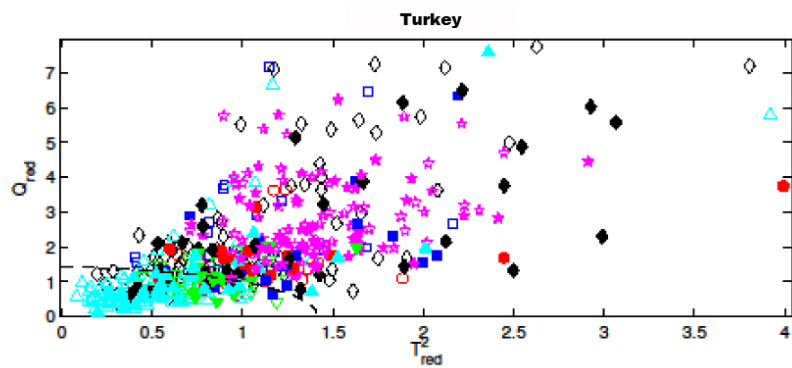
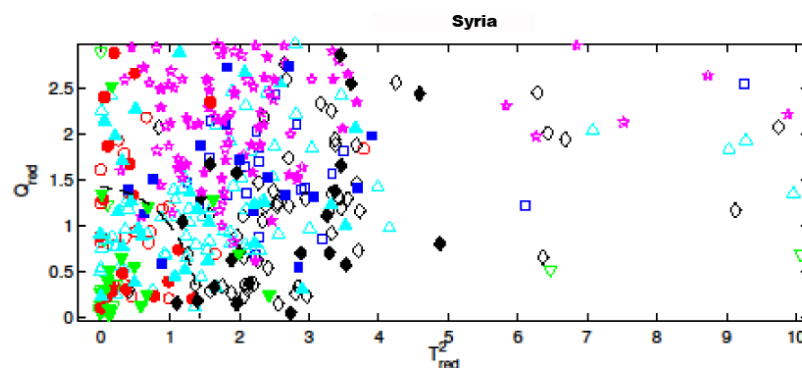
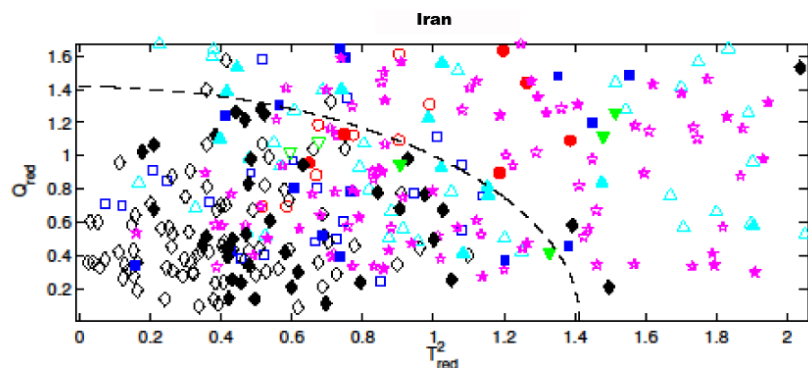
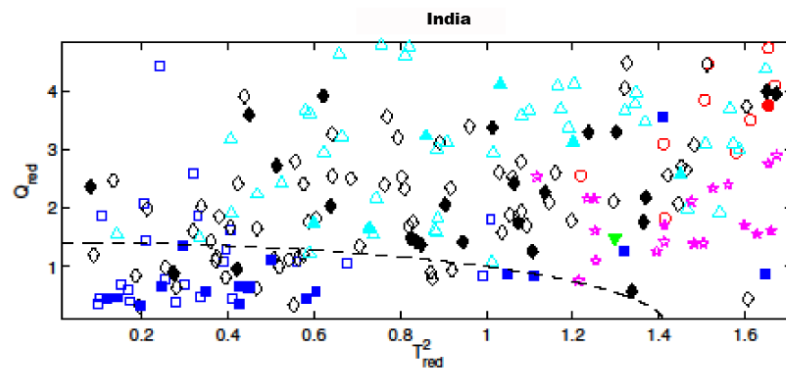
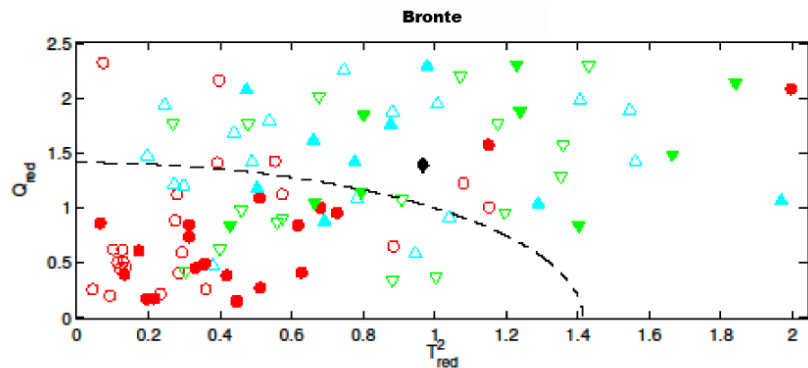


Figure 7

