

# Analytical Methods

Accepted Manuscript



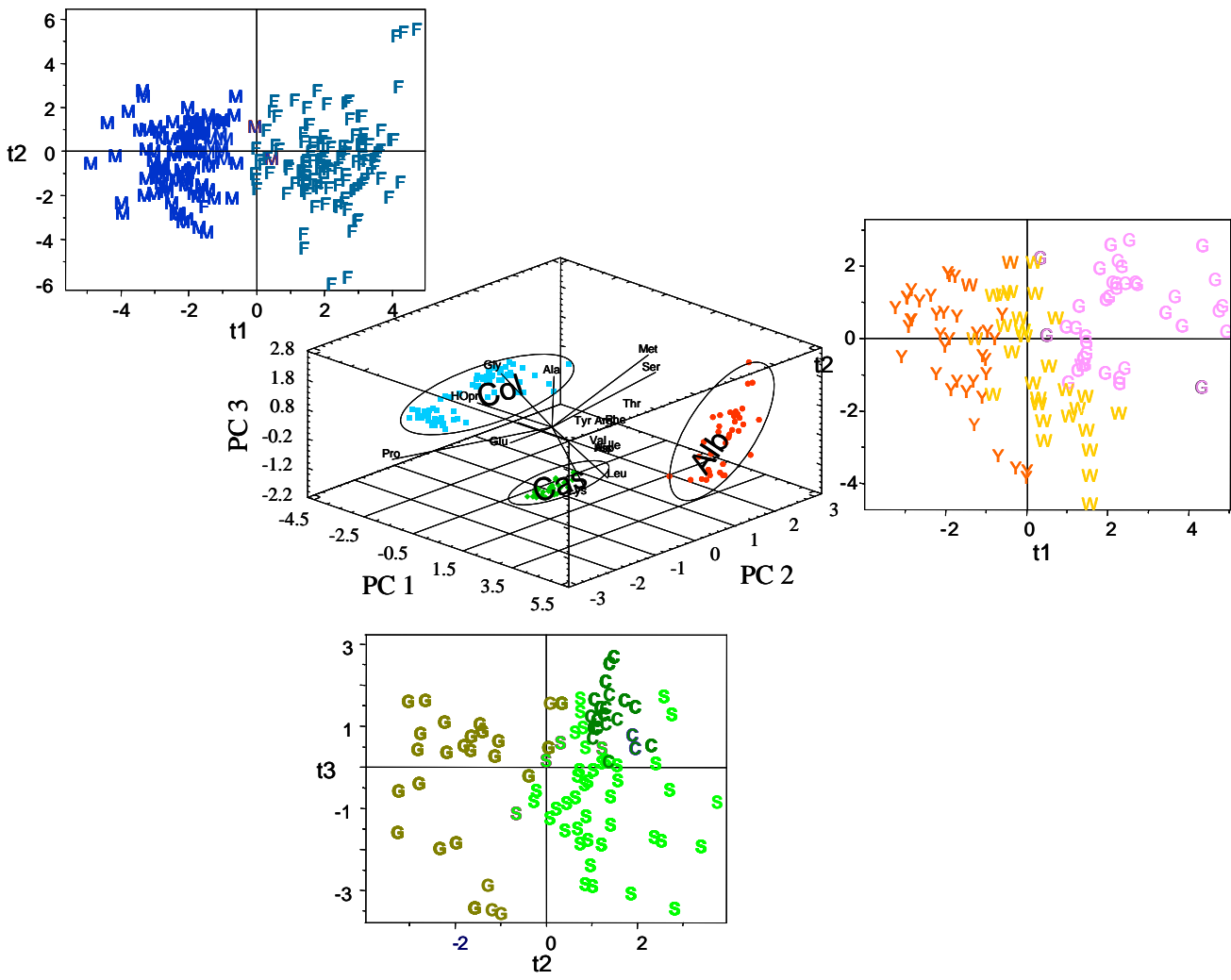
This is an *Accepted Manuscript*, which has been through the Royal Society of Chemistry peer review process and has been accepted for publication.

*Accepted Manuscripts* are published online shortly after acceptance, before technical editing, formatting and proof reading. Using this free service, authors can make their results available to the community, in citable form, before we publish the edited article. We will replace this *Accepted Manuscript* with the edited and formatted *Advance Article* as soon as it is available.

You can find more information about *Accepted Manuscripts* in the [Information for Authors](#).

Please note that technical editing may introduce minor changes to the text and/or graphics, which may alter content. The journal's standard [Terms & Conditions](#) and the [Ethical guidelines](#) still apply. In no event shall the Royal Society of Chemistry be held responsible for any errors or omissions in this *Accepted Manuscript* or any consequences arising from the use of any information it contains.

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43



# CHARACTERISATION AND CLASSIFICATION OF BINDERS USED IN ART MATERIALS AT CLASS AND SUBCLASS LEVEL

R. Checa-Moreno<sup>1</sup>, E. Manzano<sup>2</sup>, L.F. Capitán-Vallvey<sup>2\*</sup>

<sup>1</sup> *Laboratorio Central de Sanidad Animal, Ministerio de Medio Ambiente y Medio Rural y Marino, Camino del Jau s/n, E-18320 Santa Fe, Granada, Spain.*

<sup>2</sup> *ECsens. Department of Analytical Chemistry, Campus Fuentenueva, Faculty of Sciences, University of Granada, E-18071 Granada, Spain.*

## Abstract

SIMCA pattern recognition is used with amino acid chromatographic profiles in a large homemade collection of natural protein binders obtained following old recipes traditionally used by painters and considered here as the standard of classification. An initial cluster analysis of the full data set made it possible to distinguish three main classes of protein binders: albumin, casein and collagen-like substances. An additional iterative study of each class revealed a new subclass, i.e., glair, yolk and whole egg for the albumin class; goat, sheep and cow for the casein class; and mammals and fish for the collagen class. Optimized SIMCA models for each class and subclass were obtained with good results in terms of sensitivity (90-100 %), specificity (73-100 %) and interclass distance (>1.4), providing identification of the protein binder present in a set of samples of different origins such as natural products, commercial binders and works of art considered cultural heritage.

**Keywords:** Amino acids, Liquid Chromatography, Soft Independent Modelling of Class Analogy Pattern Recognition, Protein binder, Two-level classification.

## 1. Introduction

1  
2  
3  
4  
5  
6 27 The organic binders used by artists in the preparation of a painting determine the artist's  
7  
8 28 technique, and differentiate painting styles <sup>1</sup>. Moreover, the knowledge of the kind of  
9  
10 29 binder present can help specialists to authenticate or refute questionable works of art.  
11  
12 30 Artists over time have used a wide variety of procedures preserved in recipes to improve  
13  
14 31 and/or modify the painting properties of materials. The origin of the binding media present  
15  
16 32 in the pictorial layer of artworks is a question in the analysis of cultural heritage materials  
17  
18 33 that has not been resolved. This information is necessary to establish the historical  
19  
20 34 provenance of materials from among schools of art and even to authenticate or refute  
21  
22 35 questionable works of art. The substances used include drying oils, resins as components of  
23  
24 36 varnishes, sugars, proteinaceous materials and waxes, among many others, and also  
25  
26 37 complex types of mixtures of them. Since ancient times, the proteinaceous materials used  
27  
28 38 as binders in the colour layers of old paintings have been found in nature and include:  
29  
30 39 animal glues prepared from animal skin or bones containing several types of collagen, egg  
31  
32 40 white, egg yolk and casein <sup>2</sup>.

33  
34  
35  
36  
37 41 The identification of both the chemistry and origin of proteinaceous binders is not an  
38  
39 42 easy task for several reasons as reported in the literature: a) they are natural products and  
40  
41 43 artists obtain them using old recipes usually without any prior purification steps; b) the  
42  
43 44 proteinaceous materials found in paintings are used either alone, in combination with oils  
44  
45 45 or with other organic materials such as impurities resulting from their preparation; c) the  
46  
47 46 organic materials tend to suffer degradation, chemical transformations and oxidation  
48  
49 47 processes with the environment, pigments and others substances that can change their  
50  
51 48 initial chemical composition by aging and degradation processes; <sup>3-5</sup> d) the small amount of  
52  
53 49 sample available and occasionally the small percentage of binder. In addition, the difficulty  
54  
55 50 in identifying them is exacerbated by the fact that the artists might have used mixtures of  
56  
57  
58  
59  
60

1  
2  
3  
4  
5  
6 51 several types of organic materials and sometimes undocumented formulations in their  
7  
8 52 search for artistic effects and mechanical behaviours which they use to give shape to their  
9  
10 53 work <sup>2</sup>.

11  
12  
13 54 The great variety of analytical methods proposed in the literature, sample treatment  
14  
15 55 procedures, strategies and mathematical tools for data treatment have made it possible to  
16  
17 56 discriminate among oils, proteins and other classes of binders <sup>6</sup> and, although with more  
18  
19 57 difficulty, between the three types of proteinaceous materials used as paint media, i.e. egg,  
20  
21 58 casein and collagen <sup>1,7,8</sup>. The earliest works that identified protein binders were based on  
22  
23 59 the use of observational methods with stratigraphic cuts of pictorial samples based on  
24  
25 60 coloured or fluorescent reactions <sup>9</sup>, solubility tests <sup>10</sup>, immunological techniques <sup>11,12</sup>, and  
26  
27 61 more recently immunodetection-based methods <sup>13</sup>, although these have not yet been  
28  
29 62 adapted to routine analysis in conservation laboratories. The classic analytical methods  
30  
31 63 make it possible to discriminate between the general categories of binding media (oil, gum,  
32  
33 64 protein, wax and terpenic resin) by qualitative means. Different optical instrumental  
34  
35 65 techniques such as FT-IR <sup>4,14</sup>, diffuse reflection infrared spectroscopy <sup>4</sup>, Raman and micro-  
36  
37 66 Raman spectroscopy <sup>4,5,15-18</sup>, and NMR<sup>19</sup> have proven useful in the study of artworks  
38  
39 67 because of their versatility in obtaining analytical information from both inorganic and  
40  
41 68 organic materials and also performing ageing studies. Nevertheless, so far the  
42  
43 69 characterization of organic binders, in particular proteinaceous materials, has been  
44  
45 70 essentially performed using chromatographic techniques <sup>20-22</sup>. The first chromatographic  
46  
47 71 techniques, both paper (PC) and thin layer (TLC), have been progressively replaced by  
48  
49 72 high performance liquid chromatography with fluorescence or UV-Vis detection<sup>20</sup>, gas  
50  
51 73 chromatography (most commonly used with mass spectrometry detection (GC-MS),<sup>1,23,24</sup>  
52  
53 74 coupling analytical pyrolysis (Py-GC-MS) or a wet-chemical treatment of the samples prior  
54  
55  
56  
57  
58  
59  
60

1  
2  
3  
4  
5  
6 75 to CG-MS analysis (chemolysis and derivatization reactions) and capillary  
7  
8 76 electrophoresis<sup>25</sup>. Recently, proteomic techniques mainly used for biological sample  
9  
10 77 analysis have been introduced for painting analysis<sup>7,26-28</sup> once they were adapted to handle  
11  
12 78 the requirements presented by the specific situation. Proteomic approaches based on mass  
13  
14 79 spectrometry applied in conservation science have promising results for identification of  
15  
16 80 the binder protein in mixtures mainly at a group level, i.e. with egg, animal glues and milk  
17  
18 81 products. Only limited results in conservation science have recently been published: the  
19  
20 82 distinction between egg yolk and egg glair temperas<sup>29</sup>, different milk species<sup>30</sup>, and animal  
21  
22 83 glues<sup>31</sup> have been studied to some extent. This method also solves the outstanding problem  
23  
24 84 of the identification of the mixtures of proteinaceous binders, which is typical for the other  
25  
26 85 commonly used analytical methods but not that of identifying/discriminating the source of  
27  
28 86 proteinaceous binders.

29  
30  
31  
32  
33 87 In the field of cultural heritage, the identification of the categories of proteinaceous  
34  
35 88 materials through their amino acid composition is based on the evaluation of the some  
36  
37 89 chromatographic amino acid profiles or the presence of specific markers, making it possible  
38  
39 90 to differentiate between eggs, casein and collagen used as paint media. Over time, the  
40  
41 91 strategies have increased in the number of amino acids used to make the identification and  
42  
43 92 consequently the complexity of data treatment<sup>6</sup> and has become more robust. Several  
44  
45 93 strategies have been developed: (a) amino acid ratio flow charts<sup>21</sup>; (b) bidimensional plots  
46  
47 94 of amino acid ratios<sup>32</sup>; (c) joint amino acid profiles of the sample using a correlation index  
48  
49 95 estimated with amino acid profiles of samples and standard databases<sup>33</sup>, (d) multivariate  
50  
51 96 statistical analysis such as principal components analysis (PCA)<sup>34,35</sup>, factor analysis (FA)<sup>1</sup>  
52  
53 97 and neural networks<sup>36</sup>, (e) use of multivariate approaches based on the SIMCA technique<sup>8</sup>.

1  
2  
3  
4  
5  
6 98 All the strategies considered have in common the fact that they use reference  
7  
8 99 proteinaceous standards. This is another important aspect be taken into account. The  
9  
10 100 correct identification of protein binders by comparison with reference proteinaceous  
11  
12 101 standards will depend on selecting the standard used well. Many researchers have used  
13  
14 102 chemical standards of purified proteins to perform identification but, as mentioned above,  
15  
16 103 the protein binder material present in an artwork sample is an entirely natural product and  
17  
18 104 consequently a very complex substance, so it is important to use standards of natural  
19  
20 105 products similar to those used by artists in the past. Additionally, the intra-specie variability  
21  
22 106 must be taken into account by considering the protein binder standards from different  
23  
24 107 individuals belonging to the same or different species.  
25  
26  
27

28 108 This paper presents some significant results obtained from the use of the soft independent  
29  
30 109 modelling of the class analogy classification technique (SIMCA) <sup>37</sup> on the profile of amino  
31  
32 110 acids collected by HPLC-DAD analysis. Both reference materials and samples from works  
33  
34 111 of art have been analyzed using phenylisothiocyanate (PITC) as the derivatization reagent  
35  
36 112 <sup>20</sup>. Amino acid profiles were obtained from a collection of reference proteinaceous binders  
37  
38 113 prepared by us and a test set from paintings, manuscripts and sculptures from the 15–18<sup>th</sup>  
39  
40 114 centuries. With SIMCA, more than with traditional strategies, it is possible to use software  
41  
42 115 to know the confidence level for each classification made. This is performed by an  
43  
44 116 appropriate statistical F-test. The strategy is important to differentiate the painting  
45  
46 117 technique adopted by different artists and is useful for classification purposes and  
47  
48 118 provenance studies.  
49  
50  
51  
52

53 119

## 54 120 **2. Materials and Methods**

### 55 121 **2.1. Reagents and solutions.**

56  
57  
58  
59  
60

1  
2  
3  
4  
5  
6 122 All chemicals were of analytical grade. Individual standard amino acids analyzed were  
7  
8 123 purchased from Sigma (Deisenhofen, Germany), phenylisothiocyanate (PITC) and  
9  
10 124 triethylamine (TEA), hydrochloric acid, acetonitrile (HPLC quality) and acetic acid were  
11  
12 125 obtained from Panreac (Montcada i Reixac, Barcelona, Spain) and absolute ethanol from  
13  
14 126 Merck (Darmstadt, Germany). Standard stock solutions of each amino acid were prepared  
15  
16 127 by adequate weighing and dissolution in 0.1 M hydrochloride acid (HCl). Reverse osmosis  
17  
18 128 quality water was produced by a Milli-RO and Milli-Q 185 Plus purification system  
19  
20 129 (Millipore Co., Bedford, MA, USA).  
21  
22  
23  
24  
25

130

## 131 **2.2. Standards of natural protein binders.**

132 A collection of 143 natural proteinaceous binders traditionally employed by past artists was  
133 prepared (Table 1) and used as classification standards<sup>38</sup>. Egg protein standards were  
134 prepared from whole eggs or by physically separating the glair and yolk. Standards of  
135 casein were prepared from previously skimmed milks by centrifugation at 30000 R.F.C.  
136 and subsequently acid precipitation to pH 4 with hydrochloric acid, at room temperature.  
137 Collagen standards were obtained of fish skins, backbones and air bladders and mammal  
138 skins, bones and cartilage from different species by lixiviation in boiling water.  
139 Approximately 2 mL of each protein binder natural standard were individually aliquoted in  
140 5 mL vials, dry-frozen and conserved by freezing at -20°C for the correct long  
141 conservation.

142

Table 1

## 143 **2.3. Apparatus and software.**

144 A Pico Tag workstation from Waters (Milford, MA, USA) for protein hydrolysis and amino  
145 acid derivatization provided with an oven (100-150° C) was used. A Hewlett-Packard HP



1  
2  
3  
4  
5  
6 146 1090 liquid chromatograph (Palo Alto, CA, USA) provided with a Diode Array Detector  
7  
8 147 (DAD) and an Aminoquant ODS column (5  $\mu\text{m}$  200 x 2.1 mm i. d.) was used. The PITC  
9  
10 148 derivatives were identified by their retention time at 254 nm. The chromatographic  
11  
12 149 conditions for amino acid determination as PITC-derivatives were those previously  
13  
14 150 optimized by us <sup>20</sup>; column heater: 40°C; flow-rate 0.5 ml/min; buffer A: 0.28 M sodium  
15  
16 151 acetate, 0.075 % (v/v), TEA, and 6% acetonitrile (pH 6.38); buffer B: 60% acetonitrile.  
17  
18 152 Mobile phase gradient was: 0% B at 0 min, first linear gradient 2% B at 2 min, second  
19  
20 153 linear gradient 43% B at 9 min; 50 % B at 13 min.

21  
22  
23  
24 154 The Mettler AE 160 analytical balance used (Mettler-Toledo AG, Greifensee,  
25  
26 155 Switzerland) was regularly checked with certified type E2 weights (5 mg, 100 mg and 100  
27  
28 156 g). The fixed volume micropipettes (Biohit, Helsinki, Finland) were periodically controlled  
29  
30 157 through gravimetry to ensure the traceability of the results.

31  
32  
33 158 For treatment and later data analysis, the software packages Statgraphics Plus for  
34  
35 159 Windows by Statistical Graphics Corp. and SIMCA-S for Windows ver. 5.1 (1994) by  
36  
37 160 Umetri AB (Umea, Sweden) were used in a Pentium 300MHz personal computer. The  
38  
39 161 SIMCA-S software package included modules to define a data file; to scale, weigh and  
40  
41 162 transform data; to edit and list the files; to input the data, define classes and perform  
42  
43 163 principal component analysis for classes; to test the fit of data to defined classes; to  
44  
45 164 perform several plots as PC-scores, loadings, etc.

46  
47  
48  
49 165

#### 50 51 166 **2.4. Analytical procedure.**

52  
53 167 A small amount of standard protein binder or test sample (1-10 mg) was dissolved in 0.05  
54  
55 168 M pH 12.3 phosphate buffer solution and 25  $\mu\text{l}$  of this solution subjected to hydrolysis and  
56  
57 169 PITC derivatization according to the Waters Picotag© method. Before sealing the samples  
58  
59  
60

1  
2  
3  
4  
5  
6 170 in a vacuum for hydrolysis at 110°C for 16 h, the dry samples in small tubes (6 x 50 mm)  
7  
8 171 were placed in the reaction vial with 200 µl of 6 M HCl. The hydrolyzed samples were  
9  
10 172 dried and redried by adding 20 µl of ethanolic solution (ethanol-water-TEA) to ensure that  
11  
12 173 a trace amount of ammonia was left. For derivatization, the samples were coupled with 20  
13  
14 174 µl of PITC solution (ethanol-water-TEA-PITC, 7:1:2:1) for 10 min, dried again in the  
15  
16 175 workstation, and reconstituted for analysis in sample diluent (0.5 M sodium phosphate  
17  
18 176 buffer, pH 7.4, and 5% acetonitrile). The total amount of each amino acid for each standard  
19  
20 177 or test sample was determined (in picomoles) by a weighted calibration based on the peak  
21  
22 178 area to internal standard ratio. Each protein binder standard was analyzed by three-five  
23  
24 179 replicates in conditions of reproducibility in order to consider the variability of the sample  
25  
26 180 treatment method.  
27  
28  
29  
30  
31  
32

181

### 182 3. Rationale

33  
34  
35 183 SIMCA (Soft Independent Modelling of Class Analogy) <sup>39</sup> is a supervised classification  
36  
37 184 technique that builds a distinct confidence region around each class. A principal component  
38  
39 185 analysis (PCA) is performed on each separated class in the data set, and a sufficient number  
40  
41 186 of principal components are retained to account for most of the variation within each class.  
42  
43 187 New objects are considered to belong to the class if their Euclidean distance towards the  
44  
45 188 constructed PC space is not significantly larger than the Euclidean distance of the class  
46  
47 189 objects towards their PC space. The variance that is explained by the class model is called  
48  
49 190 the modelling variance, which describes the signal, whereas the noise in the data is  
50  
51 191 described by the residual variance or the variance not accounted for by the model. By  
52  
53 192 comparing the residual variance of an unknown  $S_x^2(q)$  to the average residual variance of  
54  
55  
56  
57  
58  
59  
60

1  
2  
3  
4  
5  
6 193 those samples that make up the class  $S_o^2(q)$  by an F-test it is possible to obtain a direct  
7  
8 194 measure of the similarity of the unknown sample to the class.  
9

10  
11  
12 195 
$$F = \frac{S_x^2(q)}{S_o^2(q)} \quad \text{Eq. 1}$$
  
13  
14

15 196 An advantage of SIMCA is that an unknown is only assigned to the class for which it has  
16  
17 197 a high probability. If the residual variance of a sample exceeds the upper limit for every  
18  
19 198 modelled class in the data set, the sample would not be assigned to any of the classes  
20  
21 199 because it is either an outlier or comes from a class that is not represented in the data set.  
22  
23 200 There are diagnostics to assess the quality of the data, such as the modelling power (MP)  
24  
25 201 and the discriminatory power (DP). The modelling power describes how well a variable  
26  
27 202 helps the principal components to model variation, and discriminatory power describes how  
28  
29 203 well the variable helps the principal components to classify the samples in the data set.  
30  
31 204 Variables with low modelling and discriminatory power are usually deleted from the data  
32  
33 205 because they only contribute noise to the principal component models and new models with  
34  
35 206 lower variables are developed again.  
36  
37  
38  
39

40 207 When several classes are present, it is of interest to have a measure of the distance  
41  
42 208 between each pair of classes. This can be calculated as, for instance, the pooled variance of  
43  
44 209 the residuals obtained when objects of class “one” are fitted to the class model “two”,  
45  
46 210 divided by vice versa the pooled residual variance obtained when the objects are fitted to  
47  
48 211 their “own” class model. Suppose two class  $q$  and  $r$ , are to be studied. Two different models  
49  
50 212 will be constructed for each class. The distance between two classes,  $d_{(r-q)}$  is calculated by  
51  
52 213 eq. 2, where  $S_r^2(q)$  is the residual variance of class  $r$  fitted to class  $q$  and  $S_o^2(q)$  is the  
53  
54 214 variance within class  $q$ . When the distance between two classes is close to zero, the classes  
55  
56 215 are very similar; values near to 1 indicate poor separation and values larger than 2 good  
57  
58  
59  
60

1  
2  
3  
4  
5  
6 216 resolution. This distance can be compared with F-statistics to judge the significance of the  
7  
8 217 class separation<sup>37</sup>.

10  
11  
12  
13 218 
$$d_{rq}^2 = \frac{\sum_{j=1}^p [s_r^2(q) + s_q^2(r)]}{\sum_{j=1}^p [s_o^2(q) + s_o^2(r)]}$$
 Eq. 2  
14  
15  
16  
17

18 219 In this work, the class models have been developed with a higher number of objects and  
19  
20 220 using the interclass separation as criterium of optimization, and consequently with a  
21  
22 221 different combination of variables/objects to those employed in our previously published  
23  
24 222 paper<sup>38</sup>. Data analysis is performed in a few steps: a) preliminary univariate data analysis to  
25  
26 223 detect possible outliers, information about the relevance of variables, etc.; b) cluster or  
27  
28 224 principal component analysis of the complete data set to establish classes, groups, clusters,  
29  
30 225 etc.; c) SIMCA model development of the emerging groups; c) Optimization of SIMCA  
31  
32 226 models by deleting outlier objects and noise variables. This can be achieved by choosing  
33  
34 227 variables which contain the largest amount of modelling or discriminant information for the  
35  
36 228 classification. After deleting irrelevant variables or outliers, the new PC models are refitted.  
37  
38  
39  
40 229

## 42 230 **4. Results and discussion**

### 43 231 **4.1. Homemade protein binder collection.**

44  
45  
46 232 The starting condition to build a model to classify protein binders by origin is an  
47  
48 233 arrangement of a set of samples with enough specimens. To cover a wide variety of  
49  
50 234 traditionally used protein binders, several albumin, casein and collagen-like species were  
51  
52 235 considered to build a collection of reference substances. At least two specimens belonging  
53  
54 236 to the same species was obtained whenever possible in order to consider the intra-specie  
55  
56 237 variability. The standard preparation for the proteins was done using old recipes which  
57  
58  
59  
60

1  
2  
3  
4  
5  
6 238 produced the standard substances similar to those used by past artists. Current knowledge  
7  
8 239 about the techniques used over the centuries in the creation of artworks comes mainly from  
9  
10 240 historical treatises that provide an overall view of the techniques used in different places  
11  
12 241 and ages. The book *Il Libro dell'Arte* by Cennino Cennini <sup>40</sup>, written at the beginning of the  
13  
14 242 15<sup>th</sup> century and considered a practical handbook describing common techniques from the  
15  
16 243 late 13<sup>th</sup> and mid-14<sup>th</sup> centuries, was used as reference for preparing the samples.  
17  
18  
19 244 Eggs and natural milks were obtained from local farms while collagen-like substances were  
20  
21 245 collected from different parts of fish and mammals that had been previously purchased in  
22  
23 246 different slaughterhouses and supermarkets in Granada (Spain). At least two eggs  
24  
25 247 belonging to each of the species considered were collected and used as egg protein  
26  
27 248 standards, one egg to prepare the whole egg standards and the other to obtain the glair and  
28  
29 249 yolk standards. Skimmed milk samples from different species and origins were acidified to  
30  
31 250 precipitate the casein fraction. This more artificial way was preferred to the classic way of  
32  
33 251 naturally skimming milk by letting it settle followed by lactic-induced precipitation of the  
34  
35 252 casein fraction because it is avoided some microbiological contamination and later  
36  
37 253 degradation. Human and donkey milk samples were collected from lactating mothers to  
38  
39 254 increase the variability of casein group. Human samples were provided by fully lactating,  
40  
41 255 healthy mothers during the first stage (1-5 days) of lactation. Finally, eighty-one collagen-  
42  
43 256 like standard samples were obtained of skins, bones and cartilages from mammals and  
44  
45 257 skins, backbones and air bladders from several species of fish, by a lixiviation process in  
46  
47 258 water.  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

## 260 **4.2. Data analysis.**

1  
2  
3  
4  
5  
6 261 Figure 1 shows the chromatogram types of the albumin, casein and collagen samples. There  
7  
8 262 are several clear differences among the three kinds of substances. HOpr is an amino acid  
9  
10 263 only present in collagen-like substances and is therefore useful when differentiating  
11  
12 264 between collagen and albumin or casein substances. The contents of Asp, Glu, Ser, Phe,  
13  
14 265 etc. are also interesting in terms of discriminating albumin from casein substances. The  
15  
16 266 problem is more complex when distinguishing between several substances containing the  
17  
18 267 same principal protein, for example glair, yolk or whole egg substances, since all of them  
19  
20 268 belong to the albumin-like complex where the amino acid profiles are very similar.  
21  
22 269 Therefore, using a few chromatographic peaks for differentiation may not provide enough  
23  
24 270 confidence due to the similarity of these proteins in structure and properties. The  
25  
26 271 application of multivariate statistical methods is thus helpful since it works with the overall  
27  
28 272 amino acids (peaks) and their rates, establishing the differences.  
29  
30  
31  
32

33 Figure 1  
34

35 274 All standards/specimens from the natural collection of proteinaceous binders were  
36  
37 275 analyzed using 3-5 replicates. In this way, data obtained contained the variability in each  
38  
39 276 natural species resulting from its genotype differences and also to observe any error in the  
40  
41 277 analytical method. The amino acid composition of samples has been discussed in the vast  
42  
43 278 literature in several forms including column mass injected, molar and mass percentage, etc.  
44  
45 279 Here the raw data obtained were described in a pMol-injected on column basis but the data  
46  
47 280 generated were subject to a process of internal normalization consisting of the expression  
48  
49 281 of the contents of each individual sample as a percentage of the sum of its amino acids.  
50  
51 282 (The full data for all the samples are available as Electronic Supplementary Information,  
52  
53 283 ESI Table S1). This process is appropriate for many characterization problems in which the  
54  
55 284 shape of the profile signal, and not the intensity, contains the relevant information.  
56  
57  
58  
59  
60

1  
2  
3  
4  
5  
6 285 However, the quantitative information is lost. We selected the molar percentage, although  
7  
8 286 the absolute amino acid content in several replicates samples may vary depending on the  
9  
10 287 error in weighting the sample and solution volumes used in the sample treatment.  
11

12 288 Univariate analysis. Univariate analysis was performed on the full raw data set (455 rows  
13  
14 289 x 18 columns). The Box-and-Whiskers plot analysis highlighted one outlier in the collagen-  
15  
16 290 like class (No. 112, ESI Table S1). Since none of the samples showed outliers for more  
17  
18 291 than three-four variables of the seventeen used, neither sample was rejected *a priori*. To  
19  
20 292 establish the discriminant capacity of each amino acid, a one-way ANOVA using Fisher's  
21  
22 293 least significant difference criterion (LSD) at 95% confidence level was performed using  
23  
24 294 the species as the criterion to compare the mean values (ESI Table S2). It concluded that  
25  
26 295 there were many amino acids, making it possible to completely differentiate the three main  
27  
28 296 classes: albumin, casein and collagen. Only HOpr, His, Leu and Lys showed no statistical  
29  
30 297 difference in distinguishing between albumin and casein classes and, analogously Arg and  
31  
32 298 Pro in distinguishing between albumin-casein and casein-collagen, respectively. The case  
33  
34 299 of HOpr for differentiation between albumin and casein classes is obvious because this  
35  
36 300 amino acid is not present in these kinds of proteins, but the HOpr composition in these  
37  
38 301 protein standards was written as 0.2 pMol, i.e., the detection limit. There were several  
39  
40 302 amino acids that completely differentiated the three kind of proteins considered.  
41  
42 303 Additionally, the ANOVA analysis was performed to check for the possibility of  
43  
44 304 differentiation between subclasses according to their origin but good results were not  
45  
46 305 obtained and consequently we resort to pattern recognition techniques.  
47  
48  
49  
50  
51  
52  
53  
54

55 307 **4.3. Principal component analysis.**  
56  
57  
58  
59  
60

1  
2  
3  
4  
5  
6 308 PCA makes it possible to visualize data set information in a few principal components,  
7  
8 309 retaining the maximum possible variability. Scores for each sample and loadings are  
9  
10 310 represented in the three first principal components in Figure 2. In the first principal  
11  
12 311 component, the collagen samples to the left of the graph that have a negative score are  
13  
14 312 completely separated from the remaining samples. To the right of the graph, albumins and  
15  
16 313 caseins can be separated along the second component. The casein samples have higher  
17  
18 314 scores than the albumins. The albumin class shows higher score dispersion than the  
19  
20 315 remaining two classes. Similar results were obtained from the clustering analysis.  
21  
22  
23

24 316 Figure 2

25  
26 317 Other important information obtained from the principal component analysis is the  
27  
28 318 loading plot. The variables responsible for the separation of two classes can be directly  
29  
30 319 identified. An examination of the variable loadings from the principal component analysis  
31  
32 320 showed that the contents of Gly, HOpr, Glu, Ala and Pro were the most responsible for the  
33  
34 321 formation of the collagen class, whereas the greater contents of Ser, Thr and Met were for  
35  
36 322 albumin class. Finally, the casein class can be differentiated from the albumin by the  
37  
38 323 content of Lys.  
39  
40  
41

42 324

#### 43 325 **4.4. SIMCA class modelling.**

44  
45  
46 326 SIMCA is a modelling technique that builds a model for each category or class. The centre  
47  
48 327 of the model is the mean value of the objects and the space orientation is defined by the  
49  
50 328 principal components. A range for each component is built on the basis of the score  
51  
52 329 distribution. A scale effect in raw data can be avoided by scaling the variables. The most  
53  
54 330 common way of doing this is using the z-transform, also called autoscaling. This refers to  
55  
56 331 mean-centring followed by dividing by the standard deviation for each sample. This  
57  
58  
59  
60



1  
2  
3  
4  
5  
6 332 produces a feature with zero mean and a unit variance. Multivariate analysis was performed  
7  
8 333 on the autoscaled data. The good separation of the three classes observed in the PCA plot  
9  
10 334 made it possible to construct the SIMCA models. The first objective was to find PC models  
11  
12 335 that would separate the three kinds of protein substances. The SIMCA analysis with all the  
13  
14 336 variables showed that the classes (albumin, casein and collagen) can be well described by  
15  
16 337 PC models with two, three and three components, respectively. The explained variance for  
17  
18 338 each model is 56, 67 and 64 % and the sensitivities are 90, 85 and 89 %, respectively with  
19  
20 339 an excellent specificity of 100 %. On the basis of low modelling power (MP) and low  
21  
22 340 discriminatory power (DP), several variables and objects showing a high leverage effect  
23  
24 341 were deleted from each class. The new models obtained on the basis of the remaining  
25  
26 342 variables and objects are described in Table 2.  
27  
28  
29  
30

31 343 Table 2

32  
33 344 Nine, thirteen and eleven objects were considered as outliers of the albumin, casein and  
34  
35 345 collagen classes, respectively. HOpr, an amino acid not found in the casein protein, was  
36  
37 346 used due to the z-score scale transformation of the data employed. The three class models  
38  
39 347 showed very good sensitivities and full specificity. SIMCA also provides differentiated  
40  
41 348 information about the variables through the modelling and discriminant power. The  
42  
43 349 modelling power is the contribution of each variable to the model and the discriminant  
44  
45 350 power is the capacity to differentiate among classes. All amino acids have a similar  
46  
47 351 modelling power around 0.5 in the albumin class, as do Glu, Pro, Val and Leu in the case of  
48  
49 352 casein, and the most hydrophilic as acids, Asp, Glu, Ser and Gly, are the highest modelling  
50  
51 353 variables in the collagen class. Regarding discriminant power, Glu and Gly are the most  
52  
53 354 discriminant amino acids between the casein and collagen classes; HOpr, Gly and Leu  
54  
55  
56  
57  
58  
59  
60

1  
2  
3  
4  
5  
6 355 between the albumin and collagen, and Glu and Arg between the albumin and casein. These  
7  
8 356 three class models were perfectly separated.

9  
10 357 4.4.1. Subclass analysis.

11  
12 358 Figure 2 shows that the three categories of binders are well separated in the principal  
13  
14 359 component space, but it is not clear if any distinction can be made according their origin in  
15  
16 360 each category. To investigate this, and taking in account the great number of objects  
17  
18 361 available for each category, a new separate one-to-one PC analysis was performed for each  
19  
20 362 class. Figure 3 presents the results obtained. Three new subclasses can be distinguished in  
21  
22 363 the albumin class according to the egg fraction prepared: glair, yolk and whole egg; for  
23  
24 364 casein according to the taxonomical family Bovidae: bovinæ, caprinae and genus ovis  
25  
26 365 (goat, cow and sheep); and for collagen two new subclasses related to the class of  
27  
28 366 Subphylum Vertebrata (mammals or fish). PCA can find new subgroups when a high  
29  
30 367 number of objects are available. Obviously, this ability is not due to the data treatment  
31  
32 368 systems, i.e, this is not something inherent in PC analysis, but is due to the proper nature of  
33  
34 369 the problem. New classes, namely subclasses, can appear because the objects belonging to  
35  
36 370 the subclasses have singular properties. These new properties make it possible to  
37  
38 371 differentiate among subclasses. The real virtue of the methodology (PCA, SIMCA, etc.) is  
39  
40 372 finding the new subclasses on the basis of the data set available.

41  
42  
43  
44  
45  
46  
47 373 Figure 3

48  
49 374 4.4.2. Subclassification of the albumin class.

50  
51 375 Figure 3a shows that new subclasses are well differentiated in first principal component.  
52  
53 376 Glair objects have positive scores whereas the yolk subclass is negative; obviously the  
54  
55 377 whole egg, as a mixture of glair and yolk, lies between them. New PC models for these  
56  
57 378 subclasses were built by deleting noisy aminoacids and following an iterative optimization  
58  
59  
60

1  
2  
3  
4  
5  
6 379 process based on specificity, sensitivity and especially inter-class distance criteria in order  
7  
8 380 to guarantee a good separation. The submodels are summarized in Table 2. The amino  
9  
10 381 acids used to model the yolk subclass were the most hydrophobic which fits with the fact  
11  
12 382 that the lipoproteins such as phosvitin, livetin, lipovitellin, etc. present in yolk egg are  
13  
14 383 made up of lipophilic amino acids. The amino acids HOpr, His and Met were not used to  
15  
16 384 generate either of the models of the subclass because their low modelling and  
17  
18 385 discriminating power. This was in agreement with the fact that the cluster analysis  
19  
20 386 previously performed on the variables applying Ward's clustering method, with the  
21  
22 387 Euclidean distance as the similarity measure, presented two groups (Figure 4): the first  
23  
24 388 brings together the principal amino acids used in the modelling and the second contains the  
25  
26 389 amino acids HOpr, His, Met and Tyr with no participation in any of these submodels. Note  
27  
28 390 that HOpr is not present in these proteins and His, Met and Tyr are the most irreproducible  
29  
30 391 amino acids in the hydrolysis step in the used Pico-Tag method.  
31  
32  
33  
34

35 392 Figure 4

36  
37 393 A measure of the distance between two classes  $r$  and  $q$  is calculated from a) the total  
38  
39 394 residuals obtained when all objects in class  $r$  are fitted to class model  $q$  and vice versa all  
40  
41 395 objects in class  $q$  are fitted to class model  $r$  in comparison with b) the residuals when all  
42  
43 396 objects in class  $q$  and  $r$  are fitted to their "own" class models. Table 3 gives the class  
44  
45 397 distances for (i) the initial models with all the variables and (ii) the optimized models with  
46  
47 398 the retained variables. In both cases the subclasses are fairly well separated ( $d > 1$ ) and the  
48  
49 399 separation increases when the irrelevant variables are deleted. The whole egg subclass is  
50  
51 400 also visibly closer to glair than yolk. Acceptable distances were obtained between the  
52  
53 401 whole egg and both glair and yolk (2.5 and 3.4, respectively). The distance between the  
54  
55 402 glair and yolk was the highest (6.4) as well.  
56  
57  
58  
59  
60

1  
2  
3  
4  
5  
6 403 Table 3  
7

8 404 In order to find another criterion of classification inside albumin class news PC analysis  
9  
10 405 were performed. The PC projection of the egg glair samples codified according their  
11  
12 406 phylogenetic origin made it possible to distinguish them in the space of the first principal  
13  
14 407 component. Glair samples are phylogenetically separated at the level of order or family.  
15  
16 408 All birds belong to the Animalia Kingdom, Phylum of Chordata, and Class Aves (birds). At  
17  
18 409 the order level, the birds begin to diverge: Anseriformes (ducks, geese, screamers, swans,  
19  
20 410 and waterfowl), Coliiformes (mouse birds and colies), Columbiformes (pigeons and doves),  
21  
22 411 Galliformes (chickens, fowl), Pisciformes (woodpeckers) and so on up to at least twenty-  
23  
24 412 three orders. Glair sample projection shows the two well-defined groups: samples belong to  
25  
26 413 the orders Columbiformes and Galliformes (G&G) and, on the other hand Anseriformes  
27  
28 414 samples (A). Table 4 shows the features of the new SIMCA models developed. Threonine,  
29  
30 415 aspartic acid, serine and glutamic acid are the amino acids with the greatest discriminant  
31  
32 416 power between the A and G&C classes. In other words, the amino acids with high polarity  
33  
34 417 are responsible for distinguishing between the two classes considered here. The statistical  
35  
36 418 interclass distance was 4.2 (>1), therefore showing a good separation between the two  
37  
38 419 classes. Figure 5 shows the Coomans plot of the two SIMCA models. None of the models  
39  
40 420 built admitted samples from the other class and the specificity was 100 %.

41  
42  
43  
44  
45  
46  
47 421 Table 4  
48

49 422 Figure 5  
50

51 423 The same approach with egg glair was carried out with the yolk and whole egg samples.  
52  
53 424 It was not possible to find a similar behaviour as with the egg glair. A good separation for  
54  
55 425 yolk and whole egg according to the phylogenetic origin was not found, perhaps because  
56  
57 426 the amino acid composition of these egg fractions is influenced by the great number of  
58  
59  
60

1  
2  
3  
4  
5  
6 427 protein substances present in egg yolk. These protein substances may introduce a hidden  
7  
8 428 effect in the amino acid profile of the yolk and whole egg samples.  
9

10 429 4.4.3. Subclassification of the casein class.

11  
12 430 As with the albumin class, the PC analysis was performed with the objects from the casein  
13  
14 431 class, revealing the appearance of three new subgroups or subclasses: sheep, goats and  
15  
16 432 cows. Figure 3b shows that the new groups are well differentiated. The sheep objects have  
17  
18 433 positive scores in the second PC whereas the goat objects are negative; cows with positive  
19  
20 434 scores are separated from the sheep along the third PC. The new submodels optimized are  
21  
22 435 shown in Table 2. It can see that amino acids with high polarity such as Asp and Glu were  
23  
24 436 not used to model the cow subclass. The goat and cow models showed a full specificity  
25  
26 437 whereas the sheep model reported 81%; this was because six and two objects that belonged  
27  
28 438 to the goat and cow classes respectively were incorrectly assigned to the sheep class. The  
29  
30 439 best separation between the sheep and cow classes was obtained when the most polar  
31  
32 440 amino acid (Asp to Thr) was present only in one of them. For that reason, the polar amino  
33  
34 441 acids Asp and Glu were used in the sheep class but not in the cow class. In the same way,  
35  
36 442 HOpr, Ser, Gly and His were only used to model the cow subclass but not the sheep  
37  
38 443 subclass. The modelling power of the variables retained was very similar. None of the  
39  
40 444 amino acids proved to be especially significant in modelling these subclasses. With respect  
41  
42 445 to the discriminant power, the most significant amino acids were HOpr, Ser, His, and Tyr.  
43  
44 446 Separation between the goat-sheep subclasses is due to Tyr, Gly and Met. Ser, His and Tyr  
45  
46 447 were the most important amino acids in the goat-cow differentiation. Finally HOpr, Tyr and  
47  
48 448 His had a role in distinguishing between sheep and cow. Tyr is a very important amino acid  
49  
50 449 in the separation of the three casein subclasses.  
51  
52  
53  
54  
55  
56

57 450 4.4.4. Subclassification of the collagen class.

1  
2  
3  
4  
5  
6 451 As with the albumin and casein classes, PC analysis was performed on all the objects  
7  
8 452 belonging to the collagen class in order to find new subclasses. Even though the collagen  
9  
10 453 standards were obtained from several animal parts, such as skins, backbones and air  
11  
12 454 bladders for fish and skins, bones and cartilages for mammals, the new subgroups appeared  
13  
14 455 when the samples were codified according to their phylum membership: mammalian and  
15  
16 456 fish. No separation was observed when either fish or mammal samples were projected  
17  
18 457 individually on PC plots (in other words it was not possible to distinguish between fish and  
19  
20 458 mammal samples using the animal part as criteria). Figure 3 shows that the new subclasses  
21  
22 459 are well separated. Mammal objects have negative scores on the first PC whereas fish  
23  
24 460 objects are positive. The new PC models for these subclasses were optimized as shown in  
25  
26 461 Table 2. The non-polished models ( $A=3$ ) with all the variables and available objects for  
27  
28 462 mammalian and fish subclasses reported: 60 and 57 % of variance explained, 86 and 88 %  
29  
30 463 as sensitivity and 100 and 90 % specificity, respectively; with a good interclass distance  
31  
32 464 (2.1). The optimization performed for the sub-models produced tabulated results. His, Pro,  
33  
34 465 Val and the lowest polar amino acids Leu, Phe and Lys did not participate in the modelling.  
35  
36 466 The fish model was built using the most polar amino acids (Asp-Gly), Ala and Met. On the  
37  
38 467 other hand, Ser, Met and Thr were the more important amino acids and made it possible to  
39  
40 468 discriminate between mammal and fish subclasses. This fits the data published by Eastoe.<sup>41</sup>  
41  
42 469 Not enough information is obtained from the amino acid profile to differentiate the samples  
43  
44 470 inside the mammal subgroup according their species (bovine, porcine, ovine)<sup>42</sup> or in the  
45  
46 471 fish subgroup according species or animal part (skin, bone, air bladder).  
47  
48  
49 472 This methodological strategy supported by the development of robust class models  
50  
51 473 generated using natural standards similar to those used by the old masters presented here  
52  
53  
54  
55  
56  
57  
58  
59  
60

1  
2  
3  
4  
5  
6 474 makes it possible to identify the hierarchical nature of protein binders at different levels  
7  
8 475 (class and subclass) in a non-subjective way.  
9

10 476

## 11 12 477 **5. Applications**

13  
14  
15 478 A set of test samples from different origins including natural, commercial and restoration  
16  
17 479 field samples was selected to test the feasibility of the SIMCA models and sub-models. The  
18  
19 480 origin and kind of protein present in several of them were previously known, which made it  
20  
21 481 possible to validate them. All the samples were treated as in the *Analytical Procedures*  
22  
23 482 section; compositional amino acid profiles were obtained and their distances to the  
24  
25 483 models/sub-models constructed subsequently calculated. Table 5 shows the SIMCA  
26  
27 484 distances of each test sample to the models and sub-models established. Their variance  
28  
29 485 values ( $s_i^2$ ) were statistically compared to each model ( $s_o^2$ ) by means of the F-test.  
30  
31

32  
33 486 Four casein-like test samples were analyzed. Donkey and human colostrums milk-  
34  
35 487 samples (1-3) initially collected as casein samples were outliers in their class and  
36  
37 488 subsequently considered test samples in order to discover their nature. They were not  
38  
39 489 classified as belonging to any class considered here, but it is remarkable that these samples  
40  
41 490 showed a lower distance to the albumin class. This may mean that the protein present in  
42  
43 491 these samples is like albumin. This fits with the fact that alpha-lactoalbumin is the principal  
44  
45 492 protein in milk colostrums, even higher than casein whose content in colostrums is very  
46  
47 493 low<sup>43</sup> (virtually nil). Only milk sample No. 4 is correctly identified as casein and barely as  
48  
49 494 a cow casein (P=1%).  
50  
51

52  
53 495

Table 5

54  
55 496 Test samples 5-8 were called “only glue” because the only information available about  
56  
57 497 them was that they were glues. Test samples 6 and 8 were correctly classified as collagens  
58  
59  
60

1  
2  
3  
4  
5  
6 498 (77 and 98 %, respectively) and subsequently as belonging to the mammal class with a high  
7  
8 499 probability success (68 and 75 %, respectively); whereas 5 and 7 were not well classified as  
9  
10 500 collagens (2 and 1 % probability success, respectively), perhaps because these test samples  
11  
12 501 show a high content of Leu, Ser and Thr. Test samples 9-13 were all correctly assigned to  
13  
14 502 the collagen class ( $P > 70\%$ ) and subsequently identified as mammal samples. No  
15  
16 503 identification of these samples according their animal part (bone, skin) was possible  
17  
18 504 because the respective models could not be established. This is currently under study. Any  
19  
20 505 protein identification strategy is subjected to a good selection of the appropriate standard  
21  
22 506 reference set used. The advantage of the SIMCA method is that the identification of  
23  
24 507 samples is performed only on the classes considered; making is possible to classify a  
25  
26 508 sample as unknown or not belonging to any class.  
27  
28  
29  
30

31 509 Three artwork test samples from cultural heritage restoration works were considered (14-  
32  
33 510 16). They were classified as fish collagen with a high confidence level.  
34

35 511

36 512

## 37 513 **6. Conclusions**

38  
39  
40  
41 514 We have elaborated a set of standard protein samples used as a reference to identify  
42  
43 515 protein samples through their amino acid profile using the SIMCA pattern recognition  
44  
45 516 technique. Thirteen SIMCA models at both class and subclass levels were developed and  
46  
47 517 then optimized following variance and interclass distance criteria. We have improved the  
48  
49 518 performance of SIMCA models respect to our previous approaches because of the use of  
50  
51 519 interclass distance as optimization criteria and the increase in the number of available  
52  
53 520 objects. These models were used to identify the binders in a set of test samples of different  
54  
55 521 origins, showing the validity of models built. Successful identification was made possible  
56  
57  
58  
59  
60



1  
2  
3  
4  
5  
6 522 by the availability of various reference standards. The advantage of SIMCA is that the  
7  
8 523 identification of the binder present in the samples is done only with classes that have been  
9  
10 524 considered previously, making it possible to classify a sample as belonging to one of them  
11  
12 525 or as unknown, i.e., not belonging to any class studied. Additionally, with the SIMCA  
13  
14 526 method, it is possible to know what the proteinaceous binders are, not only at a class level  
15  
16 527 but at a subclass level as well. This methodology can be applied to identifying the origin of  
17  
18 528 protein binders in artworks in the field of conservation and restoration, may provide  
19  
20 529 information about the historical provenance of materials in schools of art and might help to  
21  
22 530 authenticate or refute questionable works of art. However, at this time it has one particular  
23  
24 531 handicap: identifying mixtures of binders is difficult. New research is currently underway  
25  
26 532 in this respect.  
27  
28  
29  
30  
31  
32

533

### 534 **7. Acknowledgements**

535 This study was supported by Research Group FQM118. The authors acknowledge financial  
536 support from the *Ministerio de Ciencia e Innovación, Dirección General de Investigación y*  
537 *Gestión del Plan Nacional de I+D+i* (Spain) (Project CTQ2009-14428-C02-01) and Junta  
538 de Andalucía (*Proyecto de Excelencia P10-FQM-5974*). These projects were partially  
539 supported by European Regional Development Funds (ERDF).  
540

541

### 542 **8. References**

543

544 1 R. Aruga, P. Mirti, A. Casoli and G. Palla, *Fresenius' J. Anal. Chem.*, 1999, **365** (6),  
545 559-566.

546 2 Pacheco, F. In *Arte De La Pintura*; Ed. Cátedra, Colección Arte: Madrid, Spain,  
1990.

- 1  
2  
3  
4  
5  
6 547 3 A. Karpowicz, *Stud. Conserv.*, 1981, **26** , 153-160.  
7  
8 548 4 E. Manzano, N. Navas, R. Checa-Moreno, L. Rodriguez-Simon and L. F. Capitan-  
9 549 Vallvey, *Talanta*, 2009, **77** (5), 1724-1731.  
10  
11 550 5 A. Nevin, I. Osticioli, D. Anglos, A. Burnstock, S. Cather and E. Castellucci, *Anal.*  
12 551 *Chem.*, 2007, **79** (16), 6143-6151.  
13  
14 552 6 M. P. Colombini, A. Andreotti, I. Bonaduce, F. Modugno and E. Ribechini, *Acc.*  
15 553 *Chem. Res.*, 2010, **43** (6), 715-727.  
16  
17 554 7 W. Fremout, J. Sanyova, S. Saverwyns, P. Vandenabeele and L. Moens, *Anal.*  
18 555 *Bioanal. Chem.*, 2009, **393** (8), 1991-1999.  
19  
20 556 8 G. Musumarra and M. Fichera, *Chemom. Intell. Lab. Syst.*, 1998, **44** (1,2), 363-372.  
21  
22 557 9 M. C. Gay, *Ann. Lab. Rech. Mus. France*, 1970, **1** , 8-24.  
23  
24 558 10 L. Masschelein-Kleiner, *PACT (Rixensart, Belg. )*, 1986, **13** , 185-207.  
25  
26 559 11 P. L. Jones, *Stud. Conserv.*, 1962, **7** , 10-16.  
27  
28 560 12 M. Johnson and E. Packard, *Stud. Conserv.*, 1971, **16** , 145-164.  
29  
30 561 13 L. S. Dolci, G. Sciutto, M. Guardigli, M. Rizzoli, S. Prati, R. Mazzeo and A. Roda,  
31 562 *Anal. Bioanal. Chem.*, 2008, **392** (1-2), 29-35.  
32  
33 563 14 A. Domenech-Carbo, F. B. Reig, J. V. G. Adelantado and V. P. Martinez, *Anal. Chim.*  
34 564 *Acta*, 1996, **330** (2-3), 207-215.  
35  
36 565 15 N. Navas, J. Romero-Pastor, E. Manzano and C. Cardell, *J. Raman Spectrosc.*, 2010,  
37 566 **41** (11), 1196-1203.  
38  
39 567 16 A. Nevin, I. Osticioli, D. Anglos, A. Burnstock, S. Cather and E. Castellucci, *J.*  
40 568 *Raman Spectrosc.*, 2008, **39** (8), 993-1000.  
41  
42 569 17 P. Vandenabeele, B. Wehling, L. Moens, H. Edwards, M. De Reu and G. Van  
43 570 Hooydonk, *Anal. Chim. Acta*, 2000, **407** (1-2), 261-274.  
44  
45 571 18 J. Romero-Pastor, C. Cardell, E. Manzano, A. Yebra-Rodriguez and N. Navas, *J.*  
46 572 *Raman Spectrosc.*, 2011, **42** (12), 2137-2142.  
47  
48 573 19 E. L. Ghisalberti and I. M. Godfrey, *Stud. Conserv.*, 1998, **43** (4), 215-230.  
49  
50 574 20 R. Checa-Moreno, E. Manzano, G. Miron and L. F. Capitan-Vallvey, *J. Sep. Sci.*,  
51 575 2008, **31** (22), 3817-3828.  
52  
53 576 21 M. P. Colombini and F. Modugno, *J. Sep. Sci.*, 2004, **27** (3), 147-160.  
54  
55 577 22 M. T. Domenech-Carbo, *Anal. Chim. Acta*, 2008, **621** (2), 109-139.  
56  
57  
58  
59  
60

- 1  
2  
3  
4  
5  
6 578 23 A. Casoli, P. C. Musini and G. Palla, *J. Chromatogr. , A*, 1996, **731** (1 + 2), 237-246.  
7  
8 579 24 E. Kenndler, K. Schmidt-Beiwil, F. Mairinger and M. Pöhm, *Fresenius' J. Anal.*  
9 580 *Chem.*, 1992, **342** (1), 135-141.  
10  
11 581 25 I. Kaml, K. Vcelakova and E. Kenndler, *J. Sep. Sci.*, 2004, **27** (3), 161-166.  
12  
13 582 26 W. Fremout, M. Dhaenens, S. Saverwyns, J. Sanyova, P. Vandenabeele, D. Deforce  
14 583 and L. Moens, *Anal. Chim. Acta*, 2010, **658** (2), 156-162.  
15  
16 584 27 S. Kuckova, R. Hynek and M. Kodicek, *Anal. Bioanal. Chem.*, 2007, **388** (1), 201-  
17 585 206.  
18  
19 586 28 G. Leo, L. Cartechini, P. Pucci, A. Sgamellotti, G. Marino and L. Birolo, *Anal.*  
20 587 *Bioanal. Chem.*, 2009, **395** (7), 2269-2280.  
21  
22 588 29 C. Tokarski, E. Martin, C. Rolando and C. Cren-Olive, *Anal. Chem.*, 2006, **78** (5),  
23 589 1494-1502.  
24  
25 590 30 A. Chambery, M. Di, C. Sanges, V. Severino, M. Tarantino, A. Lamberti, A. Parente  
26 591 and P. Arcari, *Anal. Bioanal. Chem.*, 2009, **395** (7), 2281-2291.  
27  
28 592 31 W. Fremout, S. Kuckova, M. Crhova, J. Sanyova, S. Saverwyns, R. Hynek, M.  
29 593 Kodicek, P. Vandenabeele and L. Moens, *Rapid Commun. Mass Spectrom.*, 2011, **25**  
30 594 (11), 1631-1640.  
31  
32 595 32 R. Pancella and R. Bart, *Z. Kunsttechnol. Konservierung*, 1989, **3** , 101-111.  
33  
34 596 33 Schilling, M. R.; Khanjian, H. P. James & James (Science Publishers eds.), London,  
35 597 2011; pp 211-219.  
36  
37 598 34 M. P. Colombini, F. Modugno, M. Giacomelli and S. Francesconi, *J. Chromatogr. ,*  
38 599 *A*, 1999, **846** (1 + 2), 113-124.  
39  
40 600 35 G. Gautier and M. P. Colombini, *Talanta*, 2007, **73** (1), 95-102.  
41  
42 601 36 R. Lleti, L. A. Sarabia, M. C. Ortiz, R. Todeschini and M. P. Colombini, *Analyst*,  
43 602 2003, **128** (3), 281-286.  
44  
45 603 37 P. J. Gemperline, L. D. Webber and F. O. Cox, *Anal. Chem.*, 1989, **61** (2), 138-144.  
46  
47 604 38 R. Checa-Moreno, E. Manzano, G. Miron and L. F. Capitan-Vallvey, *Talanta*, 2008,  
48 605 **75** (3), 697-704.  
49  
50 606 39 S. Wold, *Pattern Recognition*, 1976, **8** (3), 127-139.  
51  
52 607 40 Cennini, C. In *The Craftsman's Handbook. Il Libro Dell'Arte*; Dover Publications,  
53 608 Inc.: New York, 1960.  
54  
55  
56  
57  
58  
59  
60

- 1  
2  
3  
4  
5  
6 609 41 J. E. Eastoe, *Biochem. J.*, 1957, **65**, 363-368.  
7  
8 610 42 M. Nemati, M. R. Oveisi, H. Abdollahi and O. Sabzevari, *J. Pharm. Biomed. Anal.*,  
9 611 2004, **34** (3), 485-492.  
10  
11 612 43 L. Hambraeus, B. Lonnerdal, E. Forsum and M. Gebre-Medhin, *Acta Paediatr.*  
12 613 *Scand.*, 1978, **67** (5), 561-565.  
13  
14 614  
15 615  
16  
17 616  
18  
19 617  
20  
21 618  
22  
23  
24 619  
25  
26 620  
27  
28 621  
29  
30  
31 622  
32

**Table 1. Collection of natural standard protein**

| ALBUMINS-Like<br>poultry <sup>(a)</sup>   | CASEINS-Like<br>Milk samples of Spanish breeds  | COLLAGENS-Like<br><sup>(b)</sup>  |
|---|---|---|
| Chicken <i>Gallus gallus</i><br>Dwarf chicken <i>Gallus gallus</i> *<br>Pheasant <i>Phasianus colchicus</i> (2)<br>Goose <i>Anser anser</i><br>Turkey <i>Meleagris gallopavo</i> (2)<br>Peacock <i>Pavo cristatus</i> (2)<br>Pigeon <i>Columba livia domestica</i><br>Duck <i>Cairina moschata</i><br>Mallard <i>Anas platyrhynchos</i> | Caprine <i>Capra aegagrus hircus</i><br>Granadina (3)<br>Granadino-Murciana<br>Capra pyrenaica hispanica (2)<br>Malagueña (2)<br>Manchega (2)<br>Bovine <i>Bos Taurus</i><br>Friesian (2)<br>Holstein (2)<br>Brown Swiss<br>Jersey<br>Ovine <i>Ovis aries</i><br>Segureña (2)<br>Merino (3)<br>Red Majorcan<br>Manchega<br>Lacha (2)<br>Churro (2)<br>Castellana (2)<br>Awassi (2)<br>Black-eye | Fish<br>Flounder <i>Platichthys flesus</i> (3)<br>Cod <i>Gadus morruha</i> (5)<br>Sturgeon <i>Acipenser sturio</i> (8)<br>Sole <i>Solea solea</i> (8)<br>Hake <i>Merluccius merluccius</i> (6)<br>Blue whiting <i>Micromesistius poutassou</i> (10)<br>Turbot <i>Psetta maxima</i> (3)<br>Mammalian:<br>Rabbit <i>Oryctolagus cuniculus</i> (5)<br>Pigs <i>Sus</i> (5)<br>Bovine <i>Bos primigenius</i> (6)<br>Ovine <i>Ovis aries</i> (8)<br>Caprine <i>Capra hircus</i> (8) |

<sup>(a)</sup> whites, yolks and whole eggs, <sup>(b)</sup> skins, backbones and air bladders of fish and skins, bones and cartilages of mammals

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

\* *holland bredd*

The number of replicated samples coming from different origins likes farms, supermarkets, etc. appears between brackets.

623

624

625

626

627

**Table 2. Models of class and subclass**

| <b>Albumin ( N = 101, P = 9 A=2 )</b>     |                                  |      |      |       |      |      |      |      |      |      |      |      |      |      |      |      |      |       |
|---|----------------------------------|------|------|-------|------|------|------|------|------|------|------|------|------|------|------|------|------|-------|
| V: 78 %                                   | Variables:                       | Asp  | Glu  | HOpr  | Ser  | Gly  | His  | Arg  | Thr  | Ala  | Pro  | Tyr  | Val  | Met  | Ile  | Leu  | Phe  | Lys   |
| Res. SD: 0.77                             | Loading[1]:                      |      | 0.40 |       | 0.36 | 0.19 |      | 0.40 | 0.37 |      | 0.24 |      |      | 0.37 |      | 0.28 |      | -0.33 |
| S: 96 %                                   | Loading[2]:                      |      | 0.11 |       | 0.25 | 0.52 |      | 0.12 | 0.08 |      | 0.49 |      |      | 0.21 |      | 0.48 |      | -0.35 |
| Sp: 100 %                                 | MP:                              |      | 0.54 |       | 0.52 | 0.49 |      | 0.58 | 0.45 |      | 0.53 |      |      | 0.50 |      | 0.62 |      | 0.54  |
| <b>Casein ( N = 79, P = 10, A=2 )</b>     |                                  |      |      |       |      |      |      |      |      |      |      |      |      |      |      |      |      |       |
| V: 73 %                                   | Variables:                       | Asp  | Glu  | HOpr  | Ser  | Gly  | His  | Arg  | Thr  | Ala  | Pro  | Tyr  | Val  | Met  | Ile  | Leu  | Phe  | Lys   |
| Res. SD: 0.80                             | Loading[1]:                      | 0.34 | 0.42 | -0.25 |      | 0.07 | 0.27 | 0.35 | 0.14 |      | 0.37 |      | 0.41 |      |      | 0.34 |      |       |
| S: 94 %                                   | Loading[2]:                      | 0.03 | 0.06 | -0.44 |      | 0.53 | 0.27 | 0.08 | 0.51 |      | 0.28 |      | 0.05 |      |      | 0.33 |      |       |
| Sp: 100 %                                 | MP:                              | 0.36 | 0.67 | 0.50  |      | 0.42 | 0.31 | 0.39 | 0.44 |      | 0.63 |      | 0.60 |      |      | 0.57 |      |       |
| <b>Collagen ( N = 204, P = 9, A=2 )</b>   |                                  |      |      |       |      |      |      |      |      |      |      |      |      |      |      |      |      |       |
| V: 75 %                                   | Variables:                       | Asp  | Glu  | HOpr  | Ser  | Gly  | His  | Arg  | Thr  | Ala  | Pro  | Tyr  | Val  | Met  | Ile  | Leu  | Phe  | Lys   |
| Res. SD: 0.82                             | Loading[1]:                      | 0.17 | 0.18 | 0.38  | 0.40 | 0.42 |      |      | 0.38 |      | 0.36 |      |      | 0.37 |      | 0.20 |      |       |
| S: 92 %                                   | Loading[2]:                      | 0.62 | 0.59 | -0.13 | 0.17 | 0.01 |      |      | 0.01 |      | 0.28 |      |      | 0.22 |      | 0.30 |      |       |
| Sp: 100 %                                 | MP:                              | 0.63 | 0.57 | 0.49  | 0.60 | 0.64 |      |      | 0.44 |      | 0.54 |      |      | 0.49 |      | 0.20 |      |       |
|   | DP <sub>albumin-casein</sub> :   | 11   | 21   | 7     |      | 5    | 11   | 19   | 14   |      | 11   |      | 6    |      |      | 7    |      |       |
|   | DP <sub>albumin-collagen</sub> : | 4    | 13   | 37    | 23   | 65   |      | 1    | 21   |      | 15   |      |      | 10   |      | 35   |      | 6     |
|   | DP <sub>casein-collagen</sub> :  | 16   | 69   | 41    | 34   | 94   | 8    | 22   | 23   |      | 11   |      | 7    | 25   |      | 26   |      |       |
| <b>Glair ( N = 33, P = 8, A = 2 )</b>     |                                  |      |      |       |      |      |      |      |      |      |      |      |      |      |      |      |      |       |
| V: 80 %                                   | Variables:                       | Asp  | Glu  | HOpr  | Ser  | Gly  | His  | Arg  | Thr  | Ala  | Pro  | Tyr  | Val  | Met  | Ile  | Leu  | Phe  | Lys   |
| Res. SD: 0.81                             | Loading[1]:                      | 0.45 |      |       |      |      |      | 0.44 | 0.41 | 0.47 | 0.29 |      |      |      |      | 0.10 | 0.34 | -0.05 |
| S: 97 %                                   | Loading[2]:                      | 0.01 |      |       |      |      |      | 0.25 | 0.23 | 0.10 | 0.36 |      |      |      |      | 0.56 | 0.37 | 0.54  |
| Sp: 73 %                                  | MP:                              | 0.50 |      |       |      |      |      | 0.64 | 0.52 | 0.59 | 0.39 |      |      |      |      | 0.62 | 0.53 | 0.52  |
| <b>Whole egg ( N = 32, P = 7, A = 2 )</b> |                                  |      |      |       |      |      |      |      |      |      |      |      |      |      |      |      |      |       |
| V: 78 %                                   | Variables:                       | Asp  | Glu  | HOpr  | Ser  | Gly  | His  | Arg  | Thr  | Ala  | Pro  | Tyr  | Val  | Met  | Ile  | Leu  | Phe  | Lys   |
| Res. SD: 0.88                             | Loading[1]:                      | 0.15 | 0.27 |       | 0.32 |      |      |      |      |      | 0.52 | 0.46 |      |      |      | 0.44 | 0.35 |       |
| S: 97 %                                   | Loading[2]:                      | 0.65 | 0.50 |       | 0.47 |      |      |      |      |      | 0.09 | 0.05 |      |      |      | 0.30 | 0.07 |       |
| Sp: 84 %                                  | MP:                              | 0.71 | 0.5  |       | 0.53 |      |      |      |      |      | 0.79 | 0.48 |      |      |      | 0.62 | 0.23 |       |
| <b>Yolk ( N = 30, P = 8, A = 2 )</b>      |                                  |      |      |       |      |      |      |      |      |      |      |      |      |      |      |      |      |       |
| V: 83 %                                   | Variables:                       | Asp  | Glu  | HOpr  | Ser  | Gly  | His  | Arg  | Thr  | Ala  | Pro  | Tyr  | Val  | Met  | Ile  | Leu  | Phe  | Lys   |
| Res. SD: 0.74                             | Loading[1]:                      |      |      |       |      | 0.40 |      | 0.37 | 0.36 |      | 0.43 |      | 0.12 |      | 0.33 | 0.41 | 0.30 |       |
| S: 100 %                                  | Loading[2]:                      |      |      |       |      | 0.20 |      | 0.14 | 0.37 |      | 0.21 |      | 0.62 |      | 0.42 | 0.22 | 0.39 |       |
| Sp: 100 %                                 | MP:                              |      |      |       |      | 0.55 |      | 0.40 | 0.62 |      | 0.73 |      | 0.67 |      | 0.66 | 0.64 | 0.46 |       |
|   | DP <sub>glair-whole egg</sub> :  | 3.0  | 1.5  |       | 1.9  |      |      | 3.3  | 1.4  | 1.1  | 3.6  | 2.3  |      |      |      | 4.7  | 2.7  | 1.5   |

|  |                                |            |            |             |            |            |            |            |            |            |            |            |            |            |            |            |            |             |
|--|--------------------------------|------------|------------|-------------|------------|------------|------------|------------|------------|------------|------------|------------|------------|------------|------------|------------|------------|-------------|
|  | DP <sub>glair-yolk</sub> :     | 3.7        |            |             |            | 7.4        |            | 7.9        | 6.3        | 1.7        | 5.6        |            | 4.2        |            | 11.0       | 9.0        | 7.5        | 4.6         |
|  | DP <sub>whole egg-yolk</sub> : | 1.1        | 1.4        |             | 4.4        | 4.4        |            | 4.2        | 5.1        |            | 2.2        | 2.9        | 1.4        |            | 5.5        | 4.8        | 3.4        |             |
| <b>Goat ( N = 20, P =10, A =3)</b>     |                                |            |            |             |            |            |            |            |            |            |            |            |            |            |            |            |            |             |
| V: 95 %                                | Variables:                     | <b>Asp</b> | <b>Glu</b> | <b>HOpr</b> | <b>Ser</b> | <b>Gly</b> | <b>His</b> | <b>Arg</b> | <b>Thr</b> | <b>Ala</b> | <b>Pro</b> | <b>Tyr</b> | <b>Val</b> | <b>Met</b> | <b>Ile</b> | <b>Leu</b> | <b>Phe</b> | <b>L ys</b> |
| Res. SD: 0.58                          | Loading[1]:                    | -0.41      | 0.27       | 0.40        |            | -0.29      |            | 0.31       |            |            | 0.08       | 0.34       | 0.24       | 0.31       |            |            |            | 0.37        |
| S: 100 %                               | Loading[2]:                    | -0.13      | -0.43      | -0.15       |            | -0.11      |            | -0.38      |            |            | 0.49       | -0.27      | 0.47       | 0.22       |            |            |            | 0.20        |
| Sp: 100 %                              | Loading[3]:                    | -0.16      | 0.18       | -0.08       |            | 0.57       |            | -0.18      |            |            | -0.40      | -0.28      | -0.03      | 0.48       |            |            |            | 0.33        |
|  | MP:                            | 0.64       | 0.69       | 0.54        |            | 0.75       |            | 0.69       |            |            | 0.75       | 0.62       | 0.67       | 0.79       |            |            |            | 0.73        |
| <b>Sheep ( N =43, P =8, A =2)</b>      |                                |            |            |             |            |            |            |            |            |            |            |            |            |            |            |            |            |             |
| V: 86 %                                | Variables:                     | <b>Asp</b> | <b>Glu</b> | <b>HOpr</b> | <b>Ser</b> | <b>Gly</b> | <b>His</b> | <b>Arg</b> | <b>Thr</b> | <b>Ala</b> | <b>Pro</b> | <b>Tyr</b> | <b>Val</b> | <b>Met</b> | <b>Ile</b> | <b>Leu</b> | <b>Phe</b> | <b>L ys</b> |
| Res. SD: 0.66                          | Loading[1]:                    | 0.40       | 0.41       |             |            |            |            | 0.31       |            | -0.29      |            |            | -0.42      |            |            | -0.37      | 0.24       | 0.36        |
| S: 98 %                                | Loading[2]:                    | -0.11      | -0.11      |             |            |            |            | -0.47      |            | -0.45      |            |            | 0.06       |            |            | 0.32       | 0.59       | 0.32        |
| Sp: 81 %                               | MP:                            | 0.58       | 0.63       |             |            |            |            | 0.65       |            | 0.53       |            |            | 0.67       |            |            | 0.67       | 0.67       | 0.58        |
| <b>Cow ( N =15, P =11, A =4)</b>       |                                |            |            |             |            |            |            |            |            |            |            |            |            |            |            |            |            |             |
| V: 98 %                                | Loading[1]:                    |            |            | 0.29        | -0.20      | -0.03      | 0.28       |            |            | -0.34      |            | 0.26       | -0.39      |            | -0.36      | -0.32      | 0.32       | 0.36        |
| Res. SD: 0.39                          | Loading[2]:                    |            |            | -0.31       | 0.45       | 0.50       | 0.35       |            |            | -0.27      |            | 0.33       | 0.07       |            | -0.13      | 0.24       | -0.21      | 0.14        |
| S: 93 %                                | Loading[3]:                    |            |            | -0.39       | -0.02      | -0.38      | -0.17      |            |            | -0.14      |            | 0.24       | 0.06       |            | 0.24       | 0.42       | 0.49       | 0.34        |
| Sp: 100 %                              | Loading[4]:                    |            |            | 0.01        | 0.17       | -0.29      | -0.30      |            |            | 0.18       |            | 0.77       | 0.01       |            | -0.09      | -0.27      | -0.24      | -0.19       |
|  | MP:                            |            |            | 0.88        | 0.80       | 0.85       | 0.79       |            |            | 0.78       |            | 0.87       | 0.81       |            | 0.67       | 0.85       | 0.89       | 0.76        |
|  | DP <sub>goat-sheep</sub> :     | 2.7        | 2.3        | 5.2         |            | 6.0        |            | 3.9        |            | 1.9        | 3.1        | 6.9        | 3.7        | 5.6        |            | 3.5        | 4.7        | 1.7         |
|  | DP <sub>goat-cow</sub> :       | 4.7        | 1.4        | 6.2         | 9.8        | 6.7        | 10         | 4.9        |            | 4.8        | 3.4        | 8.4        | 5.8        | 2.5        | 7.9        | 4.8        | 5.4        | 4.2         |
|  | DP <sub>sheep-cow</sub> :      | 2.9        | 1.5        | 9.5         | 7.3        | 6.0        | 8          | 2.2        |            | 3.5        |            | 7.7        | 1.6        |            | 6.7        | 1.4        | 2.3        | 3.6         |
| <b>Mammalian ( N = 94, P =9, A =2)</b> |                                |            |            |             |            |            |            |            |            |            |            |            |            |            |            |            |            |             |
| V: 64 %                                | Variables:                     | <b>Asp</b> | <b>Glu</b> | <b>HOpr</b> | <b>Ser</b> | <b>Gly</b> | <b>His</b> | <b>Arg</b> | <b>Thr</b> | <b>Ala</b> | <b>Pro</b> | <b>Tyr</b> | <b>Val</b> | <b>Met</b> | <b>Ile</b> | <b>Leu</b> | <b>Phe</b> | <b>L ys</b> |
| Res. SD: 0.99                          | Loading[1]:                    | -0.44      |            | 0.43        | -0.35      |            |            | -0.04      | 0.14       | 0.23       |            | 0.29       |            | 0.41       | 0.41       |            |            |             |
| S: 92 %                                | Loading[2]:                    | -0.28      |            | 0.00        | -0.09      |            |            | 0.53       | 0.46       | 0.45       |            | -0.28      |            | -0.30      | -0.24      |            |            |             |
| Sp: 100 %                              | MP:                            | 0.58       |            | 0.36        | 0.22       |            |            | 0.46       | 0.36       | 0.42       |            | 0.26       |            | 0.51       | 0.44       |            |            |             |
| <b>Fish ( N =121, P =7, A =2)</b>      |                                |            |            |             |            |            |            |            |            |            |            |            |            |            |            |            |            |             |
| V: 67 %                                | Variables:                     | <b>Asp</b> | <b>Glu</b> | <b>HOpr</b> | <b>Ser</b> | <b>Gly</b> | <b>His</b> | <b>Arg</b> | <b>Thr</b> | <b>Ala</b> | <b>Pro</b> | <b>Tyr</b> | <b>Val</b> | <b>Met</b> | <b>Ile</b> | <b>Leu</b> | <b>Phe</b> | <b>L ys</b> |
| Res. SD: 1.04                          | Loading[1]:                    | -0.31      | -0.29      | 0.11        | 0.23       | 0.49       |            |            |            | -0.51      |            |            |            | 0.50       |            |            |            |             |
| S: 90 %                                | Loading[2]:                    | -0.48      | -0.51      | 0.41        | -0.48      | -0.27      |            |            |            | 0.19       |            |            |            | 0.02       |            |            |            |             |
| Sp: 100 %                              | MP:                            | 0.48       | 0.51       | 0.24        | 0.40       | 0.50       |            |            |            | 0.49       |            |            |            | 0.35       |            |            |            |             |
|  | DP <sub>mam-fish</sub> :       | 2.2        | 2          | 5           | 5.3        | 2.1        |            | 2.3        | 4.3        | 4.1        |            | 3.7        |            | 5.3        | 2.3        |            |            |             |

V: Variance explained by the model, Res. SD: Residual standard deviation of the model, S: Sensitivity, Sp: Specificity, MP: modelling power, DP: Discriminant power

629

630

631

632

**Table 3. Distance between subclasses**

| Models:       | albumin subclass |     |     | casein subclass |     |     | collagen subclass |
|---------------|------------------|-----|-----|-----------------|-----|-----|-------------------|
|               | G-W              | G-Y | W-Y | C-O             | C-V | O-V | M-F               |
| All variables | 1.5              | 2.9 | 1.8 | 1.7             | 2.3 | 1.4 | 1.4               |
| Optimized     | 2.5              | 6.2 | 3.4 | 3.0             | 5.2 | 2.7 | 3.4               |

G: glair, W: whole egg, Y: yolk; C: caprine, O: ovine, V: vac

633



634

635

636

**Table 4. Major class parameters for phylogenetic submodels.**

| <u>Model</u>   | <u>Variance explained</u> |      | <u>Objects</u> |      | <u>Variables</u> |      |      | <u>PCs</u> |      |      | <u>Res. SD</u> |      | <u>Sensibility</u> |      | <u>Specificity</u> |  |
|--|---------------------------|------|----------------|------|------------------|------|------|------------|------|------|----------------|------|--------------------|------|--------------------|--|
| Anseriformes   | 89 %                      |      | 10             |      | 14               |      |      | 4          |      |      | 0.85           |      | 100 %              |      | 100 %              |  |
| Galliformes & Columbiformes  | 73 %                      |      | 27             |      | 13               |      |      | 3          |      |      | 0.82           |      | 89 %               |      | 85 %               |  |
|  | Asp                       | Glu  | Ser            | Gly  | His              | Arg  | Thr  | Ala        | Tyr  | Val  | Met            | Ile  | Leu                | Phe  | Lys                |  |
| <b>Modelling power:</b>  |                           |      |                |      |                  |      |      |            |      |      |                |      |                    |      |                    |  |
| Anseriformes:  | 0,62                      | 0,52 | 0,61           | 0,63 | 0,79             | -    | 0,26 | 0,40       | 0,35 | 0,61 | 0,71           | 0,62 | 0,65               | 0,72 | 0,64               |  |
| Galliformes & Columbiformes:   | -                         | 0,54 | 0,32           | 0,58 | -                | 0,60 | 0,58 | 0,62       | 0,60 | 0,49 | 0,25           | 0,38 | 0,49               | 0,55 | 0,58               |  |
| Discrim. P. A/G&C:   | 7,9                       | 6,4  | 6,6            | 3,2  | 3,6              | 0,5  | 8,9  | 5,5        | 2,7  | 3,1  | 3,5            | 4,2  | 3,1                | 4,3  | 2,0                |  |
| Res. SD: Residual standard deviation of the class, Discrim. P.: Discriminant power |                           |      |                |      |                  |      |      |            |      |      |                |      |                    |      |                    |  |

637

638

639

**Table 5.** Results of the SIMCA method. Distance ( $s_i^2$ ) for the unknown samples to every class and subclass modelled. Critical distance to each model in brackets

| Sample                    | SIMCA distance to each class/subclass model |               |               |               |               |               |               |               |               |               |               |               |               |
|---------------------------|---|---------------|---------------|---------------|---------------|---------------|---------------|---------------|---------------|---------------|---------------|---------------|---------------|
|                           | whole                                       |               |               |               |               |               |               |               |               |               |               |               |               |
|                           | alb   | glair         | egg           | yolk          | anna          | gallco        | cas           | cap           | ov            | cow           | col           | mam           | fish          |
|                           | <b>(0.77)</b>                               | <b>(0.81)</b> | <b>(0.88)</b> | <b>(0.74)</b> | <b>(0.85)</b> | <b>(0.82)</b> | <b>(0.80)</b> | <b>(0.58)</b> | <b>(0.66)</b> | <b>(0.39)</b> | <b>(0.82)</b> | <b>(0.99)</b> | <b>(1.04)</b> |
| 1. Donkey milk            | 1,3   | 2,1           | 2,1           | 2,2           | 4,5           | 2,8           | 6,9           | 10,8          | 6,5           | 7,1           | 25            | 14            | 47            |
| 2. Human milk colostrum 1 | 1,2   | 2,3           | 1,9           | 2,1           | 4,4           | 3,6           | 6,7           | 12,5          | 7,4           | 9,3           | 24            | 17            | 44            |
| 3. Human milk colostrum 2 | 1,5   | 2,2           | 1,8           | 4,4           | 4,4           | 2,8           | 10,5          | 17,4          | 9,1           | 13,8          | 22            | 25            | 36            |
| 4. Cow milk               | 2,8   | 2,4           | 3,9           | 3,3           | 6,3           | 3,9           | <b>0,63</b>   | 2,0           | 1,1           | <u>0,43</u>   | 25            | 11            | 50            |
| 5. Skin glue              | 5,5   | 3,5           | 5,4           | 7,5           | 14            | 7,6           | 23            | 35            | 17            | 32            | <u>0,88</u>   | 1,98          | <b>0,61</b>   |
| 6. Strong glue            | 5,4   | 3,5           | 5,1           | 7,5           | 13            | 7,3           | 23            | 35            | 16            | 32            | <b>0,44</b>   | <b>0,57</b>   | 1,74          |
| 7. Hausenblasen Glue      | 5,4   | 3,3           | 5,3           | 7,4           | 13            | 7,5           | 23            | 34            | 17            | 32            | <u>0,94</u>   | 2,47          | <b>0,80</b>   |
| 8. Glue                   | 5,4   | 3,4           | 5,1           | 7,5           | 13            | 7,3           | 23            | 35            | 16            | 32            | <b>0,28</b>   | <b>0,54</b>   | 1,74          |
| 9. Rabbit glue            | 5,5   | 3,5           | 5,3           | 7,5           | 14            | 7,4           | 23            | 35            | 16            | 32            | <b>0,45</b>   | <b>0,64</b>   | 1,69          |
| 10. Skin pork glue        | 5,4   | 3,5           | 5,2           | 7,5           | 13            | 7,4           | 23            | 35            | 16            | 32            | <b>0,47</b>   | <b>0,73</b>   | 1,94          |
| 11. Pork bone glue        | 5,4   | 3,5           | 5,1           | 7,4           | 13            | 7,4           | 23            | 35            | 16            | 32            | <b>0,34</b>   | <b>0,55</b>   | 1,58          |
| 12. Cow bone glue 1       | 5,5   | 3,6           | 5,3           | 7,6           | 14            | 7,5           | 23            | 35            | 16            | 32            | <b>0,22</b>   | <b>0,77</b>   | 1,42          |
| 13. Cow bone glue 2       | 5,4   | 3,5           | 5,2           | 7,5           | 13            | 7,4           | 23            | 35            | 16            | 32            | <b>0,38</b>   | <b>0,58</b>   | 1,61          |
| 14. Artwork sample 1      | 3,9   | 4,8           | 7,1           | 9,7           | 12            | 7,8           | 11            | 12            | 12            | 17            | <b>0,79</b>   | 1,7           | <b>0,71</b>   |
| 15. Artwork sample 2      | 4,4   | 4,7           | 7,5           | 9,3           | 14            | 7,6           | 11            | 14            | 12            | 17            | <b>0,71</b>   | 1,4           | <b>0,93</b>   |
| 16. Artwork sample 3      | 4,4   | 4,9           | 7,7           | 9,6           | 13            | 7,7           | 12            | 14            | 12            | 18            | <b>0,65</b>   | 3,6           | <b>0,65</b>   |

640 Probability level of belonging to class or subclass: (1-5 %) underlined, > 5% in **bold**.

641 Albumin (alb), annatidae (anna), galliformes and columbiformes (gallco), casein (cas),  
 642 caprine (cap), ovine (ov), collagen (col), mammalian collagen (mam).

643

644

645

646

1  
2  
3  
4  
5  
6 647 **Figure captions**  
7

8 648 Figure 1. HPLC-UV chromatograms at 254 nm of 17 PTH-derivatives of amino acids  
9  
10 649 presents in hydrolysates of representative binders: a) albumin, b) casein and c) collagen  
11

12 650

13  
14  
15 651 Figure 2. Scores and loadings plot of the autoscaled chromatographic data of protein binder  
16 standards in the space of the first three principal components. Albumin (red), casein  
17 652 (green), and collagen (blue) classes.  
18  
19 653  
20

21 654

22  
23  
24 655 Figure 3. Scores plot of subclasses a) albumins: (Y) yolk, (W) whole egg, (G) glair; b)  
25 caseins: (G) goat, (C) cow, (S) sheep; and c) collagens, (M) mammals, (F) fish.  
26  
27 656  
28

29 657

30  
31 658 Figure 4. Dendrogram built with 17 variables from the albumin class. Cluster 1 contains  
32 amino acids used in modelling whereas cluster 2 shows amino acids not used in the  
33 659 polished model.  
34  
35 660

36 661

37  
38  
39 662 Figure 5. Cooman's plot of the squared SIMCA distances obtained from the data set for  
40 Anseriformes (A) and Galliformes-Columbiformes (G) glair samples.  
41  
42 663  
43

44 664

45 665

46 666

47 667

48 668

49 669

50 670  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

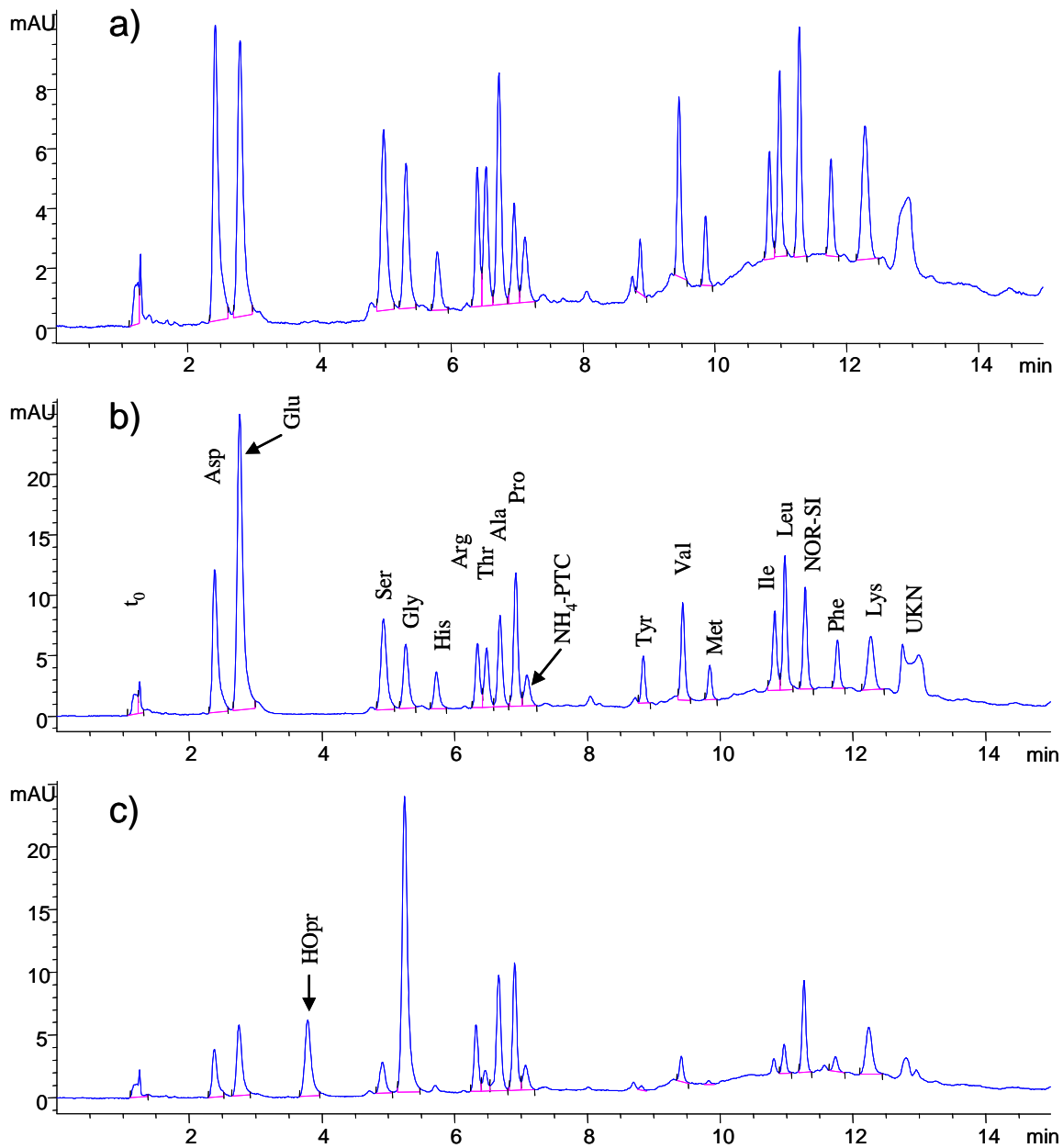


Figure 1

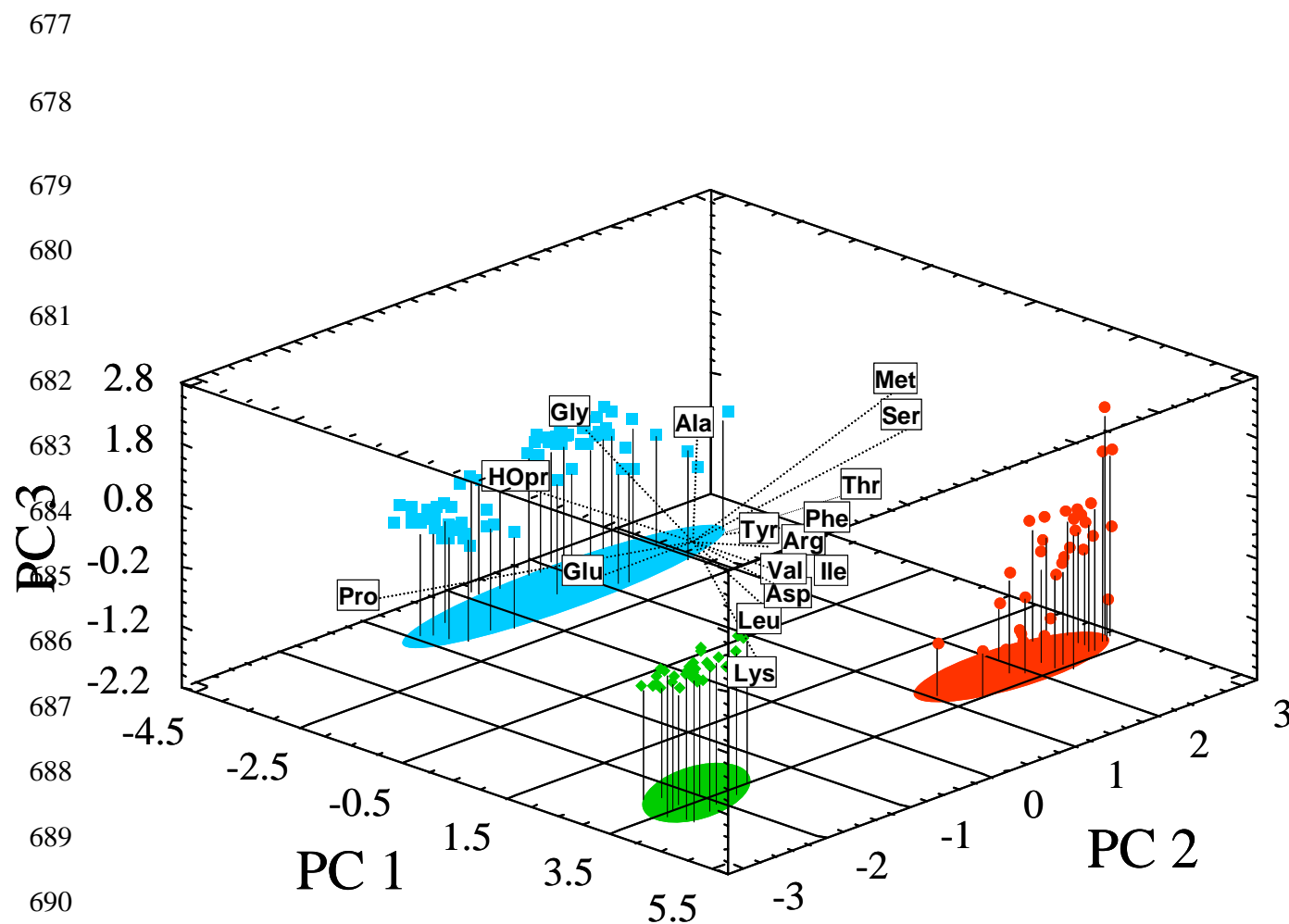


Figure 2

1  
2  
3  
4  
5  
6 696  
7  
8 697  
9  
10 698  
11  
12 699  
13  
14 700  
15  
16 701  
17  
18 702  
19  
20 703  
21  
22 704  
23  
24 705  
25  
26 706  
27  
28 707  
29  
30 708  
31  
32 709  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

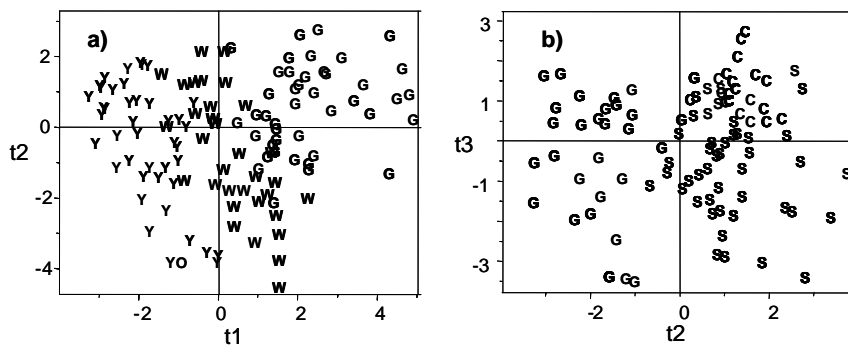
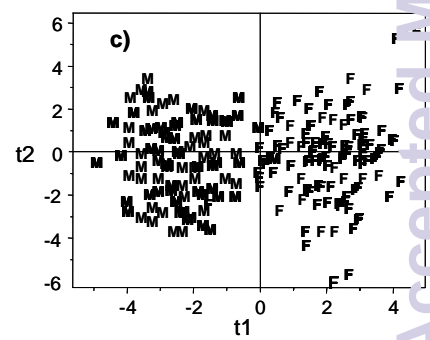


Figure 3



710

711

712

713

714

715

716

717

718

719

720

721

722

723

724

725

726

727

728

729

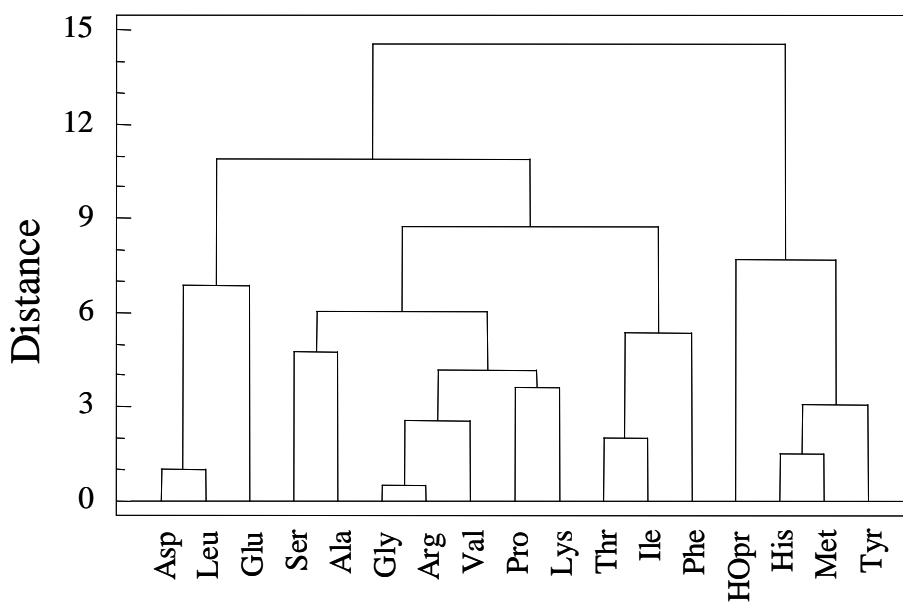


Figure 4

1  
2  
3  
4  
5  
6 730  
7  
8 731  
9  
10 732  
11  
12 733  
13  
14 734  
15  
16 735  
17  
18 736  
19  
20 737  
21  
22 738  
23  
24 739  
25  
26 740  
27  
28 741  
29  
30 742  
31  
32 743  
33  
34 744  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

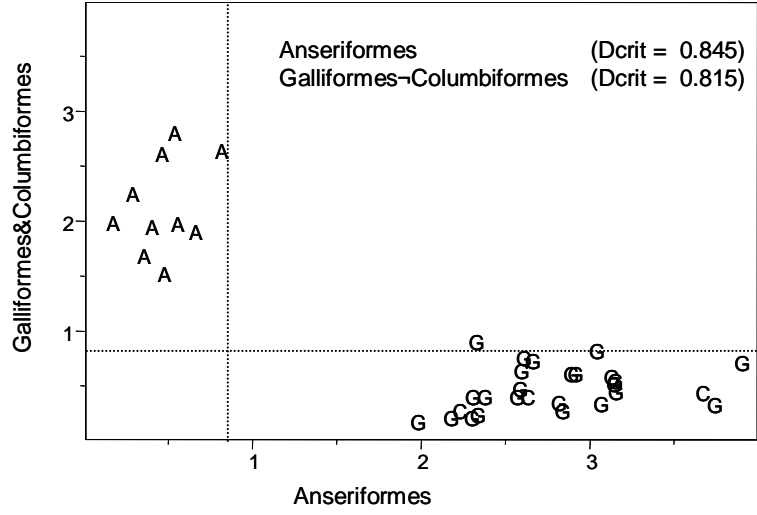


Figure 5