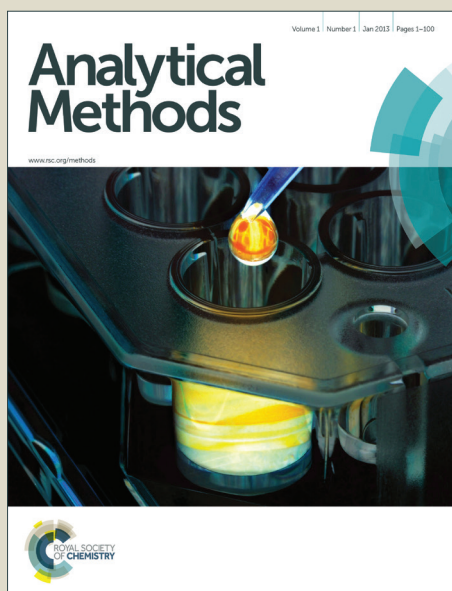


# Analytical Methods

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 1  
4  
5  
6  
7 2 **Assessing the discrimination potential of linear and non-linear**  
8  
9  
10 3 **supervised chemometrics methods on a filamentous fungi FTIR**  
11  
12  
13 4 **spectral database**  
14  
15  
16 5  
17  
18 6  
19

20 7 V. Gaydou<sup>1,2</sup>, A. Lecellier<sup>1,2</sup>, D. Toubas<sup>1,2,3</sup>, J. Mounier<sup>4</sup>, L. Castrec<sup>4</sup>, G. Barbier<sup>4</sup>, W.  
21  
22 8 Ablain<sup>5</sup>, M. Manfait<sup>1,2</sup>, G.D. Sockalingum<sup>1,2\*</sup>  
23  
24  
25 9  
26  
27  
28 10  
29

30 11 <sup>1</sup> Université de Reims Champagne-Ardenne, MÉDIAN-Biophotonique et Technologies pour  
31  
32 12 la Santé, UFR de Pharmacie, 51 rue Cognacq-Jay, 51096 REIMS cedex, France  
33

34 13 <sup>2</sup> CNRS UMR7369, Matrice Extracellulaire et Dynamique Cellulaire, MEDyC, Reims, France  
35

36 14 <sup>3</sup> Laboratoire de Parasitologie Mycologie, CHU de Reims, Hôpital Maison Blanche, 45 rue  
37  
38 15 Cognacq-Jay, 51092 Reims cedex, France  
39

40 16 <sup>4</sup> Laboratoire Universitaire de Biodiversité et Ecologie Microbienne (EA3882), SFR148  
41  
42 17 ScInBioS, Université Européenne de Bretagne, Université de Brest, ESIAB, Technopôle de  
43  
44 18 Brest Iroise, 29280 Plouzané, France  
45

46 19 <sup>5</sup> AES CHEMUNEX/BIOMERIEUX, Rue Maryse Bastié, CS17219 Ker Lann, 35172 Bruz  
47  
48 20 cedex, France  
49

50  
51  
52 21  
53  
54 22 \*Corresponding author:  
55

56 23 Ganesh D. Sockalingum  
57

58 24 Université de Reims Champagne-Ardenne  
59  
60

1  
2  
3 25 MéDIAN, Biophotonique et Technologies pour la Santé  
4

5 26 Unité MEDyC, CNRS FRE3481  
6

7 27 UFR Pharmacie, SFR CAP-Santé FED4231  
8

9 28 51 rue Cognacq-Jay, Reims, France.  
10

11 29 Tel: +33 3 26 91 35 53  
12

13 30 Fax: +33 3 26 91 35 50  
14

15 31 [ganesh.sockalingum@univ-reims.fr](mailto:ganesh.sockalingum@univ-reims.fr)  
16

17  
18  
19 32  
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

33

**Abstract**

This study proposes a comparative investigation of different linear and non-linear chemometrics methods applied to the same database of infrared spectra for filamentous fungi discrimination and identification. The database was comprised of 277 strains, (14 genus, 36 species), identified and validated by DNA sequencing, and analyzed by high-throughput Fourier Transform Infrared (FTIR) spectroscopy in the 4000-400  $\text{cm}^{-1}$  wavenumber range. A cascade of 20 supervised models based on taxonomic ranks was constructed to predict classes until the species taxonomic rank. The cascade modeling was used to test 11 algorithms (5 linear and 6 non-linear) of supervised classification methods. To assess these algorithms, indicators of classification rates and McNemar's tests were defined and applied in same way to each of them. For non-linear algorithms, the KNN (K Nearest Neighbors) method proved to be the best classifier (78%). Linear algorithms, PLS-DA (Partial Least Square - Discriminant Analysis) and SVM (Support Vector Machine) showed better performances than non-linear methods with the best classification potential (~93%). SVM and PLS-DA were comparable and a possible complementarity between these two algorithms was highlighted.

**Keywords:** Supervised classification, cascade models, infrared spectroscopy, fungi identification

52

## 1 Introduction

Spectrometric techniques play an important role in both research and industrial applications. The development of these techniques has continuously progressed in order to exploit at best their capacities. Among these, infrared spectroscopy has emerged as a promising approach for rapid analysis. In mid-infrared spectroscopy, the molecular fundamental vibrational modes are measured and involve wavelengths between 2.5 and 25 micrometers corresponding to the wavenumber range of 4000-400  $\text{cm}^{-1}$ . It relies on the absorption of mid-infrared light by vibrational transitions in covalent bonds. Fourier transform infrared (FTIR) spectroscopy has high molecular sensitivity and reveals numerous types and modes of vibrations. It is fast, label-free, cost-effective, easy to use, and applicable to various fields. However, it is perturbed by aqueous states and by the atmospheric water vapor and carbon dioxide. In near-infrared (NIR) spectroscopy, the sample receives wavelengths in the range of 800-2500 nm, whereby molecular overtone and combination vibrations are measured. The spectra are more complex and it can be difficult to assign specific features to specific chemical components. The molar absorptivity in the NIR region is typically quite small but NIR has the advantage that it can typically penetrate much farther into a sample than mid infrared radiation. NIR spectroscopy is, therefore, not a particularly sensitive technique, but it can be very useful in probing bulk material with little or no sample preparation. Numerous analytical attempts have been described in the literature with their advantages and disadvantages [1]. FTIR and NIR approaches have proved to be very effective in the characterization, differentiation, and identification of various fungi [2,3,4,5]. Another analytical method that has emerged as a new discovery tool for bacterial characterization (Lay, 2001) is matrix-assisted laser desorption/ionization time-of-flight mass spectroscopy (MALDI-TOF MS). It is an analytical tool sensitive to molecular composition and distinct mass signals can be observed in a mass-

1  
2  
3 78 to-charge ( $m/z$ ) ratio. It allows molecular profiling such as protein profiling. Its potential to  
4  
5 79 discriminate filamentous fungi of clinical origin at the species level has been demonstrated  
6  
7  
8 80 giving comparable results as with molecular identification methods (Cassagne et al., 2011; De  
9  
10 81 Carolis et al., 2012). However, MALDI-TOF spectrometric databases for filamentous fungi,  
11  
12 82 particularly from the food industry, are still under development (Santos et al., 2010b). A  
13  
14  
15 83 recent work, reports on the useful integration of different analytical imaging techniques,  
16  
17 84 including FTIR and MALDI-TOF, in a multimodality platform for a deeper characterization  
18  
19  
20 85 of the potential medicinal fungus *Hericium coralloides* [6].  
21

22 86  
23  
24 87 The progress in analytical spectroscopy and speed of data acquisition has also led to the  
25  
26  
27 88 construction of large and complex datasets. In order to exploit these large datasets  
28  
29 89 sophisticated statistical methods were developed [7,8,9]. The field of chemometrics has thus  
30  
31  
32 90 emerged as a powerful approach for data mining, interpretation, and understanding;  
33  
34 91 specifically for extracting relevant molecular information in different fields of spectroscopy.  
35  
36 92 Recent advances in computing and chemometric allow choosing a wide variety of statistical  
37  
38  
39 93 algorithms to analyze the same spectral data bank.

40  
41 94 The aim of this study was to compare the discriminating potential of 11 algorithms on the  
42  
43 95 same dataset of 5960 FTIR spectra of filamentous fungi collected from 277 fungal strains  
44  
45  
46 96 belonging to 14 genera and 36 species. Among these, 194 strains (4159 spectra) were used for  
47  
48  
49 97 the model optimization and calibration steps and 83 strains (1801 spectra) were used for the  
50  
51 98 external validation step. The assessed methods were all supervised discrimination methods  
52  
53 99 requiring a calibration step and grouped in two categories. The first category concerns the  
54  
55 100 linear methods with the Factorial Discriminant Analysis (FDA) method which is the most  
56  
57  
58 101 famous method and was introduced by Fisher in 1936 [10]. Then comes the Linear  
59  
60 102 Discriminant Analysis (LDA) method, with a discrimination rule quite equivalent to that of

1  
2  
3 103 FDA [11]. The Partial Least Square-Discriminant Analysis (PLS-DA) method is more recent,  
4  
5 104 ensues from the PLS algorithm, and was reported by Wold and Martens in 1983 [12]. The  
6  
7  
8 105 Soft Independent Modeling of Class Analogies (SIMCA) method was described by Wold and  
9  
10 106 Sjöström in 1977 [13], which is a less used method. The second category concerns the non-  
11  
12 107 linear methods comprised the Support Vector Machine (SVM) method was proposed by  
13  
14  
15 108 Vapnik et al. in 1963 [14]. The K-Nearest Neighbors (K-NN) methods, introduced by J. H.  
16  
17 109 Friedman in 1975 [15], was developed to answer discrimination challenge with various kind  
18  
19 110 of data. The Probabilistic Neural Network (PNN) method is relatively recent and was presented  
20  
21 111 by Specht in 1990 [16]. It is based on the analogy with the functioning of the brain of superior  
22  
23 112 organisms. Networks are formed by small units, called neurons connected them. Quadratic  
24  
25 113 Discriminant Analysis (QDA) method, proposed by Wold in 1976, is based on a quadratic  
26  
27 114 function to apply the discrimination law [9].  
28  
29  
30  
31  
32  
33

34 116 The chemometrics question raised in this study was to select the most appropriate statistical  
35  
36 117 method able to discriminate and identify an unknown strain of filamentous fungi from its  
37  
38 118 FTIR spectrum and using a pre-established and non-exhaustive spectral library recently  
39  
40  
41 119 generated in our group [2,3].  
42

43 120 In order to assess these supervised discrimination methods in terms of statistical significance,  
44  
45 121 indicators of classification rates and McNemar's tests were defined and applied in same way  
46  
47 122 to each of the studied algorithm.  
48  
49

50 123

51 124

## 52 125 **2 Materials and Methods**

53 126

### 54 127 **2.1 Sample preparation and FTIR analysis**

1  
2  
3 128  
4  
5 129 Two hundred and seventy-seven fungi strains belonging to 14 genera and 36 species and  
6  
7  
8 130 yielding 5960 spectra were used in this study. They were from the Université de Bretagne  
9  
10 131 Occidentale and Centraalbureau voor Schimmelcultures culture collections and were  
11  
12 132 identified by sequencing of specific DNA region like the rDNA internal transcribed spacer  
13  
14 133 (ITS) region.  
15  
16 134 Cryopreserved strains were first sub-cultured on Sabouraud agar slants (Becton Dickinson, Le  
17  
18 135 Pont de Claix, France) and incubated for 4 to 7 days at 25 °C depending on the strain. The  
19  
20 136 cultures were transferred to an M tube, each sample was dissociated using a cycle of 100  
21  
22 137 seconds. Two milliliters of dissociated mycelia suspension were then transferred into an  
23  
24 138 Eppendorf tube. The culture medium was then eliminated by centrifugation, the mycelia were  
25  
26 139 washed and suspended in 1 ml of 0.9 % physiological saline and the supernatant was  
27  
28 140 eliminated by another centrifugation. The mycelia pellets were recovered in 300 µl of  
29  
30 141 physiological saline. Finally the samples were deposited on an IR-transparent 384-well silicon  
31  
32 142 plate and dried into thin films. For reproducibility concerns, 3 independent cultures of each  
33  
34 143 strain prepared on 3 different days (biological replicates), were performed and for each  
35  
36 144 culture several instrumental replicates were recorded. The plate was then analyzed with a  
37  
38 145 high-throughput module (HTS-XT) coupled with a Tensor 27 FTIR spectrometer (Bruker  
39  
40 146 Optics, Ettlingen, Germany).  
41  
42 147 The FTIR acquisition parameters were 64 accumulations per well with a spectral resolution of  
43  
44 148 4 cm<sup>-1</sup>, in the spectral range of 4000-400 cm<sup>-1</sup>. The spectral recording and preprocessing  
45  
46 149 procedures were carried out by the OPUS 6.5. The preprocessing included baseline correction,  
47  
48 150 second derivative, and vector normalization. The wavenumber ranges selected for the data  
49  
50 151 bank were 800-800 cm<sup>-1</sup> and 2800-3200 cm<sup>-1</sup>. Details of the experimental protocol were  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60 152 reported recently [2,3].



153

## 2.2 Cascade modeling and building of the calibration and validation sets

155

156 For FTIR spectroscopic data, the establishment of a single model of discrimination,  
157 parameterized by more than around thirty clusters is quite challenging, particularly for linear  
158 algorithms. Such one-model procedure is difficult to implement since the zones of variance  
159 and co-variance overlap and become inconsistent with the number of clusters. For this reason,  
160 a modeling called “in cascade” has been developed [17] to circumvent the problem in this  
161 study (figure 2). The particularity of the cascade modeling is that it is parameterized from a  
162 reference arborescence and for the study presented here, it is the taxonomic classification of  
163 fungi that is used in this respect. At every taxonomic rank, samples were distributed in  
164 subphylum, class, order, family, genus, subgenus, section, serial, and species. In so doing,  
165 several “subgroups” were established at every rank and for each model the number of clusters  
166 was around 3 and so on, until the last rank called “species” rank is reached. The taxonomic  
167 tree is thus used to structure the data matrix in a subgroup and cluster cascade. We call  
168 “taxonomic nodes” the subgroups highlighted by the taxonomic tree. For every taxonomic  
169 node, a discrimination model was built. So, this technique allowed constructing the  
170 discrimination model in cascade including not less than 20 models with a maximum of 7  
171 models required to reach the species taxonomic rank as regards to *Camemberti* strains.  
172 The main advantage of the cascade modeling is that it allows obtaining a strong method of  
173 discrimination although the final number of clusters is high. On the other hand, this method is  
174 completely parameterized and thus totally dependent on the cascade reference to which it is  
175 associated. Yet, the fungal taxonomy is in constant evolution and consequently training  
176 variation on taxonomic nodes can influence the outcome in a significant way.

1  
2  
3 177 The data matrix was split into 2 sets; about 4159 spectra (194 strains) of samples were  
4  
5 178 attributed to the calibration set and the rest (1801 spectra, 83 strains) to the validation set. In  
6  
7  
8 179 fact, one third of strains of each of 36 species represented in spectral data bank was randomly  
9  
10 180 selected to constitute the validation spectral data set and to ensure that the relative variance of  
11  
12 181 the validation set is inferior to that of the calibration set [18].  
13  
14  
15 182

### 17 183 **2.3 Partial cross validation for parameter optimization**

18 184  
19  
20  
21  
22 185 Fundamentally, cross validation was developed for chemometrics experiments with a low  
23  
24 186 sample population [19]. Because of this low population it is impossible to split the data matrix  
25  
26 187 into calibration and validation sets while keeping a representative sample set. Thus, the cross  
27  
28 188 validation allows to estimate the accuracy and robustness with one sample set only. For the  
29  
30 189 present study, the calibration set was used with cross validation to optimize chemometrics  
31  
32 190 parameters of all the studied algorithms presented in table 1 [20].  
33  
34

35  
36 191 A large number of spectra are available in the calibration set. However, although the number  
37  
38 192 of samples is quite high, the proportion between the number of strains (194) and that of the  
39  
40 193 number of species (36) is close to five. But many species present in the calibration set is  
41  
42 194 represented by only 2 different strains. Therefore, the use of cross validation is justified. The  
43  
44 195 cascade structure of all models is complex and the calibration set is constituted of biological  
45  
46 196 and technical replicates. These two features must be taken into account in the implementation  
47  
48 197 of the cross validation.  
49  
50

51  
52  
53 198 The partial cross validation by culture was chosen and scripted such that in every cross  
54  
55 199 validation iteration, all replicate spectra associated to the same culture were removed, then  
56  
57 200 used as test in the validation phase. By applying this partial cross validation, it was able to test  
58  
59 201 all cultures of the calibration set. The partial cross validation by culture allows estimating  
60

1  
2  
3 202 (partially) the intra-species, intra-strains, and intra-cultures co-variances. Further, concerning  
4  
5 203 the species represented by only 2 strains, this cross validation algorithm was more stable less  
6  
7  
8 204 over-fitted than partial cross validation by strains.  
9

10 205

## 11 206 **2.4 Percentage of Good Prediction (PGP) and McNemar test**

12 207

13 208 The accuracy of each used algorithm was computed during the three following steps: the cross  
14  
15 209 validation step to explore the algorithm's parameters, the computing model step to build  
16  
17 210 discrimination models with optimized parameters, and the validation step to evaluate the final  
18  
19 211 expected accuracy.  
20

21 212 The statistical indices here called the Percentage of Good Prediction (PGP) were calculated at  
22  
23 213 the end of these three steps. This index was calculated by dividing the number of well  
24  
25 214 identified spectra by the total number of spectra to predict. They allowed estimating the  
26  
27 215 accuracy of the discrimination models at each step. It is possible to calculate for these three  
28  
29 216 steps the PGP for each model independently. However, during the validation step and for  
30  
31 217 better presentation, the PGP was calculated only by taxonomic rank.  
32

33 218 McNemar's test is a statistical procedure that allows estimating if the prediction powers of  
34  
35 219 two methods are significantly different. This test is based on a  $\chi^2$  with one degree of freedom  
36  
37 220 because the sample number of each model is always higher than twenty ( $\alpha$ : type I error). The  
38  
39 221  $\chi^2$  critical value with a 95% level of significance, written  $\chi^2_{(1,0.95)}$  is equal to 3.8414.  
40

41 222 The McNemar's test was chosen because a unique training set was used for each model and  
42  
43 223 each algorithm. In this condition, the McNemar's test allows low probability of Type I error  
44  
45 224 and presents a powerful ability to differentiate between two algorithms. [21].  
46  
47

48 225

## 49 226 **2.5 Linear and non-linear chemometrics methods**

1  
2  
3 227  
4  
5 228 The methodological rules of these two categories are entirely different and the data are not  
6  
7  
8 229 visualized in the same way. For linear methods, the variance of the explanatory variables is  
9  
10 230 considered as linear and a proportionality relationship between them and the variables to  
11  
12 231 explain is assumed. Non-linear methods take into account two types of variances, the global  
13  
14 232 variance of the explanatory variables and variance of variables to explain, and then try to  
15  
16  
17 233 correlate these by means of a non-linear function such as the polynomial Kernel function for  
18  
19 234 SVM algorithm. Also, for these two categories of algorithms, chemometrics models were not  
20  
21 235 built around the same statistical rules. For supervised discrimination studies, the variety of  
22  
23 236 chemometrics methods available is quite diverse. The linear methods are generally the most  
24  
25  
26 237 used with spectroscopic data. Indeed, the linearity relationship put forward by the Beer-  
27  
28 238 Lambert expression, linking concentration and absorption, implies that the linear approach  
29  
30 239 appears better [22]. However, the evolution of non-linear methods has allowed the elaboration  
31  
32 240 of effective approaches such as SVM or Neural Network, which have been successfully  
33  
34 241 applied to numerous experimental cases, including complex biological spectral data [23].  
35  
36 242 In order to optimize data mining and improve the understanding of biological phenomena  
37  
38 243 from spectral results, it becomes essential to evaluate both linear and non-linear methods.  
39  
40 244 Many of these linear and non-linear algorithms were declined in various specific algorithms,  
41  
42 245 e.g., for the PLS algorithm, it was declined in robust or double PLS, quadratic PLS, splines  
43  
44 246 function PLS or GIFI-PLS and many algorithms were combined such as the neural networks  
45  
46 247 PLS or the least squares SVM [24]. For this study, only the “classical” (not “declined”)  
47  
48 248 algorithms were used in order to assess the fundamental computing methodology of each of  
49  
50 249 the following described algorithms.  
51  
52  
53  
54  
55  
56  
57  
58 250  
59  
60 251

### 2.5.1 Linear algorithms

1  
2  
3 252  
4  
5 253 LDA is a linear method of supervised discrimination that can improve the spreading of the  
6  
7  
8 254 sample distribution. The aims of this method are to maximize the ratio of the inter- to intra-  
9  
10 255 class distances and to find a linear transformation allowing to achieve the maximum class  
11  
12 256 discrimination. However, for the classical LDA the scatter matrices must be non-singular,  
13  
14  
15 257 which is well-known as the under sampling problem. To get round this problem many  
16  
17 258 solutions exist. One of them is to precede LDA by a Principal Component Analysis (PCA) in  
18  
19 259 order to extract the discriminant information. For this study PCA-LDA was tested although  
20  
21 260 this algorithm may lead to a loss of discriminant information during the PCA step [25].  
22  
23  
24 261 FDA aims at finding the subspace of the original variable space that best separates clusters by  
25  
26 262 maximizing the inter-class variance with regard to the total variance [26]. This descriptive  
27  
28 263 analysis builds a discriminant model to determine which cluster a new sample belongs to.  
29  
30  
31 264 This is simply done by projecting this sample onto the eigenvectors space and by selecting the  
32  
33 265 nearest cluster. Several distances can be used for this decision, the Euclidean distance was  
34  
35 266 preferred.  
36  
37  
38 267 Wold and Sjöström were the first to describe the SIMCA chemometrics method [13]. It is a  
39  
40 268 supervised classification method which considers every “cluster of samples” or “groups”  
41  
42 269 separately. This method is very useful for classifying high-dimensional observations because  
43  
44 270 it incorporates PCA for dimension reduction. So for every cluster, decomposition into  
45  
46 271 principal components (PC) is carried out providing a matrix of scores and loadings for each.  
47  
48 272 The most practical interest of this analysis is that each cluster can be reduced to a set of PCs  
49  
50 273 [27] and during the calibration step, the optimal PCs are determined by means of their  
51  
52 274 explained variance. After PCA steps, the discrimination models are built using Euclidean  
53  
54 275 distance between clusters and PCA subspaces, taking into account the information and  
55  
56 276 properties of clusters.  
57  
58  
59  
60

1  
2  
3 277 PLS-DA is a supervised classification method based on the multivariate PLS regression  
4  
5 278 algorithm [28]. This algorithm allows to mathematically maximize the variance-covariance  
6  
7  
8 279 between the explanatory variable matrix and the property variable matrix. PLS-DA applies the  
9  
10 280 multivariate PLS algorithm to establish discrimination rules by means of a binary matrix. The  
11  
12 281 validation samples were attributed by means of the predicted binary code. The highest  
13  
14  
15 282 predicted binary code variable gives the predicted cluster [29].  
16  
17  
18 283

### 19 284 **2.5.2 Non-linear chemometrics methods**

20 285  
21  
22 286 QDA is non-linear algorithm because it is based on a quadratic function but it is not very  
23  
24  
25 287 much different from LDA except that it is assumed that the covariance matrix can be different  
26  
27 288 for each cluster, where it is estimated separately as a Gaussian distribution. The Gaussian  
28  
29 289 parameters for each cluster are computed from training points with maximum likelihood  
30  
31 290 estimation [30]. For this study, this method was applied on the PCA scores of the data matrix.  
32  
33 291 KNN techniques were developed to answer challenges about density estimation and pattern  
34  
35 292 classification [31]. Processing of this algorithm consists of basically ordering the training  
36  
37 293 samples in a d-dimensional unit hypercube by means of a metrics distance measure. Then, for  
38  
39 294 each tested sample, the training matrix is examined in the order of their projected distance  
40  
41 295 from the tested sample on the sorted coordinate. The prediction of the unknown sample is  
42  
43 296 determined by the most representative cluster of the k nearest neighbors [32]. To optimize the  
44  
45 297 training model, the k integer and the metrics of distance can be adjusted.  
46  
47 298 Neural networks were successfully used to solve complicated pattern recognition and  
48  
49 299 classification problems in different domains. The probabilistic neural networks (PNN) method  
50  
51 300 presents a few advantages over the conventional neural network [33]. It provides a robust  
52  
53 301 classification with noisy data. PNN combines different concepts: neural computing, Bayes

1  
2  
3 302 classification rule, and non parametric estimation of the probability density function. In this  
4  
5 303 study, the PNN method was employed on the eigenvalues of the data matrix, after PCA  
6  
7  
8 304 preprocessing, and the Mahalanobis method was used for distance computing.  
9  
10 305 SVM is a supervised method originally proposed by Vapnik et al. in 1963. Fifty years later,  
11  
12 306 many publications reporting on SVM and its extensions as a multiclass classification method  
13  
14  
15 307 can be found in literature [34]. The SVM algorithm classifies data by finding the best  
16  
17 308 hyperplane that separates all data points of one class from the others classes. The best  
18  
19 309 hyperplane for an SVM corresponds to the one with the largest margin between the two  
20  
21 310 classes. In this study, the nu-SVM algorithm was employed and this algorithm could be used  
22  
23 311 with many Kernel functions (table 1). When the SVM was used with a linear Kernel function,  
24  
25 312 this algorithm was considered as a linear algorithm, and for this study, the results of SVM  
26  
27 313 with a linear Kernel function, also called linear SVM, were associated with the results  
28  
29 314 obtained by other linear methods.  
30  
31  
32  
33

34 315

## 36 316 **2.6 Computing**

37 317

38  
39  
40 318 All the chemometrics analyses were performed with Matlab R2013a (32-bit) (Mathwork,  
41  
42 319 USA) and was used to classify the samples using their explanatory variables. The algorithms  
43  
44 320 used for LDA, QDA and KNN were available in the pure Matlab. The algorithms used for  
45  
46 321 SIMCA were developed by Cleiton A. Nunes, Brazil (available on Mathwork/matlabcentral).  
47  
48 322 The algorithms used for SVM called lib-SVM were developed by Chih-Chung Chang and  
49  
50 323 Chih-Jen Lin, China [34]. The algorithms used for FDA, PLS-DA and PNN were developed  
51  
52 324 by Dominique Bertrand and Christophe Cordella, INRA, France [35].  
53  
54  
55 325 The computing was realized on a personal computer with 2Go RAM, an Intel Core2 Duo  
56  
57 326 2.66GHz as processor and Microsoft Vista (32 bit). The required time for models computing  
58  
59  
60

1  
2  
3 327 step and validation step was negligible (only few decades of seconds) compared to that of the  
4  
5 328 optimization step. The number of parameters to optimize was the factor which had the most  
6  
7  
8 329 influence on the total required time. For example, the PLS-DA algorithm required to optimize  
9  
10 330 only one parameter and took twelve hours to explore this parameter from 1 to 35,  
11  
12 331 corresponding to 1 845 900 computing models (20 models x 35 range parameter x 879 total  
13  
14 332 cross validated strains x 3 cultures per strain). On the other hand, for polynomial SVM, four  
15  
16 333 days were needed to optimize four parameters (113 127 300 computing models).  
17  
18  
19  
20 334  
21  
22 335  
23

### 24 336 **3 Results and discussion**

25  
26 337  
27  
28  
29 338 For all linear and non-linear methods described previously, the results obtained took into  
30  
31 339 account the three following steps: optimization of the chemometric parameters, model  
32  
33 340 computing, and validation in cascade. The details of the models are presented in table 2 and  
34  
35 341 each model was validated taking into account the taxonomic tree as described in figure 2.  
36  
37 342 Twenty models were built to complete the cascade and one cascade was built for each tested  
38  
39 343 algorithm (a total of 220 optimized models). During the optimization and computing steps, the  
40  
41 344 models were built independently by means of the calibration set (or part of the calibration  
42  
43 345 set), but during the validation step, the models were tested by the validation set and were  
44  
45 346 interlocked into each other.  
46  
47  
48  
49

#### 50 347 51 52 348 **3.1 Optimization of chemometric parameters and model computing**

53 349  
54  
55 350 The various methods of discrimination compared in this work required an optimization of  
56  
57 351 parameters. These parameters were different from each other and were directly associated to  
58  
59  
60



1  
2  
3 352 the chemometrics method used (table 1). They naturally have a strong influence on the final  
4  
5 353 results and it was thus essential to optimize these parameters in the most rigorous way.  
6  
7  
8 354 The parameters were optimized by means of the calibration set only and for each algorithm,  
9  
10 355 the influence of the variability of the various parameters or combination of parameters (e. g.,  
11  
12 356 SVM with polynomial Kernel function) was explored culture wise by partial cross validation.  
13  
14  
15 357 In the scope of this study, it is not possible to describe all the chemometric parameters for all  
16  
17 358 algorithms and all models. Thus, we have therefore chosen to illustrate with the PLS-DA  
18  
19 359 algorithm.  
20  
21  
22 360 Concerning this algorithm, the parameter to optimize was the used number of latent variables  
23  
24 361 (LV). It corresponds to the number of computed regression vectors. In this study, the LV  
25  
26 362 parameter varied from 1 to 35. The limit of 35 was chosen principally for computing reasons.  
27  
28  
29 363 These parameters were optimized by partial cross validation and each LV was tested with  
30  
31 364 each culture. The average of Percentage of Good Prediction (PGP) as a function of the LV  
32  
33 365 was computed and plotted highlighting a maximum of PGP (figure 3). The LV corresponding  
34  
35 366 to this maximum was taken as the best parameter. In fact, the LV optimization was important  
36  
37 367 in order to minimize model under- and over-fitting [9]. All of the LV optimal number were  
38  
39 368 defined for each model of the cascade and each model was built with its appropriate LV as  
40  
41 369 presented in table 3 (see the LV column). It was possible to observe that the LV parameters  
42  
43 370 highlighted the complexity of the model because the more the model took into account a high  
44  
45 371 number of clusters, the higher was the LV parameter.  
46  
47  
48 372 After parameter optimization, the models were computed taking into account these optimized  
49  
50 373 parameters. During the model computing step, the discrimination abilities of each model and  
51  
52 374 each algorithm was evaluated by means of the PGP of calibration. The calibration results  
53  
54 375 concerning PLS-DA and linear SVM algorithms are presented in table 3.  
55  
56  
57  
58  
59  
60 376

### 377 3.2 Validation step

378

379 The prediction capacity of all the classification models was evaluated by means of the sample  
380 validation set. This step allowed observing, in real conditions, the behavior of the various  
381 models tested in this investigation. Figures 4.1 and 4.2 show the broken curves corresponding  
382 to the PGP of validation spectra of each tested algorithm versus the taxonomic rank.

383 Concerning the linear methods, the best results were obtained with the PLS-DA method  
384 (Figure 4.1). This method allowed reaching a PGP of 98.9% for the genus taxonomic rank and  
385 93.2% for the species taxonomic rank. The LDA and FDA methods respectively gave a PGP  
386 around 3% and 6% less than the PLS-DA method, with 96.4% and 95% for the genus  
387 taxonomic rank and 89.6% and 85.8% for the species taxonomic rank. The broken curve of  
388 SIMCA is not shown so as to preserve the best scale for PGP. This method showed the worst  
389 results with PGP of 66.5% for the genus taxonomic rank and less than 50% for the species  
390 taxonomic rank. The linear SVM gave very good results with a PGP of 99.8% for the genus  
391 taxonomic rank and 91.3% for the species taxonomic rank. This algorithm was able to reach  
392 the second best PGP for the last taxonomic rank and it appears also adapted to the  
393 identification of fungi described in this study.

394 The PLS-DA and linear SVM algorithm were equivalent and showed a superior ability for  
395 correct identification than the other linear methods. This is because LDA and FDA methods  
396 use a mathematical algorithm based mainly on variance of the spectral data matrix, whereas  
397 the PLS-DA and SVM algorithms are based on combination of the spectral and the reference  
398 matrices. Concerning LDA and FDA algorithms, the validation results showed a relative  
399 similarity for both methods. This particularity could be explained by the methodology  
400 employed, which is based on minimization of the Euclidian distance between validation  
401 samples and the cluster barycenters. On the opposite, the SIMCA method uses the internal

1  
2  
3 402 variance of each cluster separately. The pertinent variance searched for each model becomes  
4  
5 403 finer when the modeling is near to the taxonomic species rank. When the SIMCA method  
6  
7  
8 404 estimates the internal variance of each cluster, the variance of interest is probably occluded by  
9  
10 405 other internal cluster variances.  
11

12 406  
13  
14  
15 407 In figure 4.2, the 3 non-linear SVM (RBF, sigmoid, and polynomial) showed PGP values near  
16  
17 408 to 100% only down to the family taxonomic rank. For the following ranks, these PGP  
18  
19 409 decreased strongly, at the genus taxonomic rank (92%, 82% and 43%) and at species  
20  
21 410 taxonomic rank (42%, 51% and 25%). Concerning the other non-linear algorithms, the best  
22  
23 411 result was obtained with the KNN algorithm and gave a PGP of 90.4% and 78.2%  
24  
25 412 respectively for genus and species taxonomic ranks. The PGP of this algorithm was close to  
26  
27 413 100% down to the family taxonomic rank and was about 15% less than the PLS-DA algorithm  
28  
29 414 from the genus to the species rank. The second best non-linear algorithm was the QDA  
30  
31 415 algorithm. This algorithm gave PGP values close to the KNN algorithm, nearly 5% less, with  
32  
33 416 a PGP of 71.5% for the species taxonomic rank. Finally, the PNN algorithm gave the worst  
34  
35 417 results, comparable to SVM with the polynomial Kernel function, with PGP values around  
36  
37 418 50% for the class taxonomic rank, which then decreased substantially down to the species  
38  
39 419 taxonomic rank.  
40  
41  
42  
43  
44

45 420 Non linear SVM did not perform very well probably because the variance of interest was not  
46  
47 421 adapted to these parameters. The optimization step could also induce an over-fitting in the  
48  
49 422 model. In addition, the interlocked cascade models could most certainly exacerbate this effect.  
50  
51 423 Concerning the other non-linear algorithms, a plausible explanation is that the variance of  
52  
53 424 interest could not be efficiently extracted by these non-linear algorithms probably due to the  
54  
55 425 linearity rules associating the spectral data matrix with the taxonomy.  
56  
57  
58  
59  
60 426

1  
2  
3 427 In order to evaluate if the prediction power of two methods was significantly different, the  
4  
5 428 McNemar's test was applied to the investigated methods pairwise. The results are displayed in  
6  
7  
8 429 table 4 in the form of a two-dimension correlation matrix. The tests were computed for the  
9  
10 430 species taxonomic rank. This test showed that all these algorithms were significantly different  
11  
12 431 except for LDA versus linear SVM and PNN versus SVM (polynomial Kernel function).  
13  
14  
15 432 For the linear and non-linear methods all the curves presented in figure 4.1 and 4.2 showed a  
16  
17 433 decreasing tendency. This decrease could be correlated with the variance sought at each  
18  
19 434 taxonomic rank. We noticed that the change from the genus to the subgenus and from the  
20  
21  
22 435 section to the taxonomic species rank induced complications. These may be associated firstly,  
23  
24 436 to the morphological proximity at the subgenus rank and to the closeness of the biochemical  
25  
26 437 structures at the species rank. These observations seemed to converge with difficulties  
27  
28 438 encountered during morphological identification, in particular for *Aspergillus 2*, *Camemberti*,  
29  
30 439 *Chrysogena* and *Roquefortorum* species models [36, 37, 38].  
31  
32  
33  
34 440

### 36 441 **3.3 Combined cascade**

37  
38 442  
39  
40  
41 443 Table 3 underlines some problematical models for both PLS-DA and linear SVM algorithms.  
42  
43 444 The concerned models were principally *Aspergillus 2*, *Camemberti*, *Chrysogena*, and  
44  
45 445 particularly *Roquefortorum*, presenting respectively validation PGP values of 85.7, 83.3, 86.5  
46  
47 446 and 61% for PLS-DA and 0, 37.5, 75 and 59% for SVM. These 4 models gave, during the  
48  
49 447 calibration step, a PGP of 100% (or close to 100%) but their validation was low and  
50  
51 448 illustrated a lack of accuracy. These inconveniencies could be explained for *Aspergillus 2*,  
52  
53 449 *Chrysogena* and *Camemberti* models, by the low number of strains and by the high number of  
54  
55 450 needed taxonomic nodes. For the *Roquefortorum* model, the taxonomy of strains concerned  
56  
57 451 by this model is still in evolution [36] and the difficulties to discriminate by sequencing and  
58  
59  
60

1  
2  
3 452 the genetic proximity of these concerned strains were real and correlated with the outcome of  
4  
5 453 our chemometrics models.

7  
8 454 On other hand, a complementarity between PLS-DA and linear SVM algorithms was  
9  
10 455 suggested through the data presented in figure 4.1. At the subgenus rank, these two algorithms  
11  
12 456 gave similar accuracy with respectively PGP of 98.9 and 99.0%. At the genus rank, the PGP  
13  
14 457 of these two algorithms were 98.9 and 99.9% respectively, and at the species rank, 93.2 and  
15  
16 458 90.1% respectively. The remarkable intersection point between these two broken curves is  
17  
18 459 pointed by a circle in figure 4.1. Due to this singularity between these two algorithms, a  
19  
20 460 “combined cascade” model was built. The linear SVM algorithm was used to elaborate the  
21  
22 461 first 8 models from the subphylum to the genus taxonomic ranks while the PLS-DA algorithm  
23  
24 462 was used to build 12 models from the sub-genus to species taxonomic ranks. The validation  
25  
26 463 results of this combined cascade are presented in figure 5.

27  
28 464 The "combined cascade" model showed the best performances compared to all the previous  
29  
30 465 "regular cascades", with 94.2% of PGP at species taxonomic rank. The spectra that were  
31  
32 466 wrongly predicted by PLS-DA were then correctly predicted by linear SVM at genus  
33  
34 467 taxonomic rank, and were almost all correctly identified until the species taxonomic rank. The  
35  
36 468 gain of one percent (illustrated by  $\epsilon$  in figure 5), due to the combined cascade, was maintained  
37  
38 469 from the genus down to the species taxonomic ranks. The “combined cascade” model  
39  
40 470 revealed that linear SVM appeared to be the most pertinent algorithm to discriminate fungi  
41  
42 471 strains until the genus taxonomic rank. This suggested that linear SVM could be better  
43  
44 472 adapted than PLS-DA algorithm for voluminous sample sets. On the other hand, PLS-DA  
45  
46 473 appeared as the most pertinent method to identify fungi at the species taxonomic rank,  
47  
48 474 suggesting that PLS-DA could be adapted for reduced and clustered sample sets. Thus, the  
49  
50 475 combination of both methods indicated an improvement of the identification capacity.  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60 476

477

**4 Conclusion**

479

480 This is the first study that compares eleven linear and non linear supervised classification  
481 algorithms on such a large dataset of food-related fungi FTIR spectra. The results obtained  
482 highlight the suitability of the linear classification methods, in particular the PLS-DA and  
483 linear SVM algorithms, for discriminating and identifying from the family to the species  
484 taxonomic ranks. These findings are promising but also pointed out the dependence due to the  
485 taxonomic references and consequently the limits of the supervised cascade computing for the  
486 application to spectral data. These observations seem to corroborate with difficulties  
487 associated with the morphological and biochemical identification. The “combined cascade”  
488 modeling including the two well suited models, PLS-DA and linear SVM, gave an  
489 improvement of the identification accuracy from the subphylum to the species taxonomic  
490 ranks. This study also highlights the interest of the concept of cascade modeling based on  
491 taxonomy because of the size and nature of the data set. Indeed, extending a complex  
492 discrimination problem into several steps allowed to distribute the studied variance on several  
493 models and thus to target the adequate variance on every taxonomic node. Further, the  
494 supervised cascade model amplifies the discrimination capacity of each tested algorithm by  
495 means of the interlocked models. The McNemar’s results pinpoint that the choice of the  
496 supervised cascade to develop chemometrics discrimination methods was appropriate to  
497 assess many discrimination algorithms. To perform knowledge about abilities of these  
498 supervised discrimination algorithms, linear and non-linear algorithms need to be assessed  
499 using other types of data, such as bio-morphological data or using other study cases. Also, the  
500 fungi spectral data bank could be used in a non-supervised way to define new clustering or  
501 new cascade of classification. In addition, by means of PLS-DA regression vector, rand

1  
2  
3 502 feature or ANOVA (ANalysis Of Variance) algorithms, it would very interesting to study  
4  
5 503 spectroscopic markers for each model in order to link spectral, biological, and chemical  
6  
7  
8 504 properties of fungi.  
9

10 505

11  
12 506

## 15 507 **6 Acknowledgements**

16  
17 508

18  
19  
20 509 The technological platform PICT-IBiSA “Imagerie Cellulaire et Tissulaire” and the “Pôle de  
21  
22 510 compétitivité” Valorial, La Région Bretagne, La Région Champagne-Ardenne are gratefully  
23  
24 511 acknowledged. Financial supports under the MOLDID project, project "Mycotech" of the  
25  
26 512 European Union, the Région Bretagne, and the Conseil Général du Finistère are also  
27  
28 513 gratefully acknowledged. The authors thank Marie-Anne Le Bras and Valérie Vasseur for  
29  
30 514 their expertise and help in fungal identification, Amélie Weill and Olivia Le Bourhis for their  
31  
32 515 excellent technical assistance, Cyril Gobinet for his reading of the chemometrics section, and  
33  
34  
35 516 Dominic Bertrand for his pertinent and constructive review of the manuscript.  
36  
37

38 517

39  
40 518  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

**Table 1:** Optimized parameters used for the different chemometrics methods

Chemometrics methods	Used parameters	Chemometrics methods	Used parameters
<i>Linear methods</i>		<i>Non-linear methods</i>	
<b>LDA</b>	<b>Kdim</b> (positive integer included in 1 to 35): size of eigenvalues matrix	<b>QDA</b>	<b>Kdim</b> (positive integer included in 1 to 35): size of eigenvalues matrix
<b>FDA</b>	<b>maxscore</b> (integer included in 1 to 35): size of PCA-score matrix allowed to the model (PCA step) <b>Kdim</b> (positive integer included in 1 to 35): size of eigenvalues matrix	<b>KNN</b>	<b>NumNeighbors</b> (positive integer included in 1 to 30) : specifying the number of nearest neighbors in calibration data to find for classifying each point when predicting <b>Metric choice:</b> function use to specify the distance metric between neighbors (among 11 distances metric process)
<b>SIMCA</b>	<b>maxscores</b> (positive integer included in 1 to 35): size of PCA-score matrix allowed to each clusters (PCA step)	<b>PNN</b>	<b>FN</b> (positive integer included in 1 to 35): the number of computed iterations <b><math>\sigma^2</math></b> (positive real included in 0 to $+\infty$ ): "smoothing parameter" of the probability function estimator
<b>PLS-DA</b>	<b>LV</b> (positive integer included in 1 to 35) is the numbers of computed regression vector defined by the Latent Variables	<b>SVM</b>	<b>Kernel function choice</b> (among 3 K-functions: <b>RBF</b> , <b>Sigmoid</b> and <b>polynomial</b> ) <b>v</b> (positive real included in 0 to 1): "level of detail" or hyperplan resolution <b><math>\gamma</math></b> (positive real included in 0 to $+\infty$ ): selected value of $\gamma$ in Kernel function (RBF, sigmoid and polynomial choice) <b>coef0</b> (positive real included in 0 to $+\infty$ ): selected value of coef0 in Kernel function (RBF and sigmoid choice) <b>d</b> (positive integer included in 1 to 5): selected degree in kernel function (polynomial choice)
<b>SVM</b>	<b>Linear Kernel function choice</b> <b>v</b> (positive real included in 0 to 1): "level of detail" or hyperplan resolution		



**Table 2:** Details of the 20 discrimination models tested.

Taxonomic rank	Model name	Corresponding clusters names	Number of spectra (number of strains)	
			Optimization & calibration	External validation
Subphylum	<i>Micromycetes</i>	<i>Pezizomycotina</i> <i>Mucoromycotina</i>	4159 (194)	1801 (83)
Class	<i>Pezizomycotina</i>	<i>Eurotiomycetes</i> <i>Sordariomycetes</i> <i>Saccharomycetes</i> <i>Dothideomycetes</i>	3302 (154)	1430 (66)
Order	<i>Dothideomycetes</i>	<i>Dothideales</i> <i>Pleosporales</i>	144 (7)	39 (2)
	<i>Sordariomycetes</i>	<i>Xylariales</i> <i>Hypocreales</i>	1087 (51)	433 (19)
Family	<i>Mucorales</i>	<i>Mucoraceae</i> <i>Lichtheimiaceae</i>	857 (42)	371 (17)
	<i>Hypocreales</i>	<i>Cordycipitaceae</i> <i>Nectriaceae</i>	1046 (45)	412 (18)
Genus	<i>Mucoraceae</i>	<i>Rhizopus</i> <i>Mucor</i> <i>Actinomucor</i>	822 (40)	355 (16)
	<i>Trichocomaceae</i>	<i>Paecilomyces</i> <i>Penicillium 1</i> <i>Aspergillus 1</i>	1948 (91)	886 (42)
Subgenus	<i>Penicillium 1</i>	<i>Penicillium 2</i> <i>Aspergilloides</i>	1043 (49)	441 (22)
Section	<i>Aspergillus 1</i>	<i>Flavi</i> <i>Fumigati</i> <i>Nidulantes</i> <i>Nigri</i> <i>Aspergillus 2</i>	851 (39)	427 (19)
	<i>Penicillium 2</i>	<i>Brevicompacta</i> <i>Fasciculata</i> <i>Chrysogena</i> <i>Roquefortorum</i> <i>Penicillium 3</i>	808 (37)	337 (17)
Serial	<i>Fasciculata</i>	<i>Camemberti</i> <i>Verrucosa</i>	190 (9)	45 (2)
Species	<i>Nidulantes</i>	<i>E. nidulans</i> <i>A. versicolor</i>	147 (7)	57 (3)
	<i>Aspergillus 2</i>	<i>E. chevalieri</i> <i>E. amstelodami</i>	110 (5)	21 (1)
	<i>Fusarium</i>	<i>F. equiseti</i> <i>F. graminearum</i> <i>F. verticillioides</i> <i>F. oxysporum</i>	1001 (43)	392 (17)
	<i>Mucor</i>	<i>M. circinelloides</i> <i>M. racemosus</i> <i>M. spinosus</i>	647 (31)	320 (13)
	<i>Camemberti</i>	<i>P. biforme</i> <i>P. camemberti</i>	108 (5)	24 (1)
	<i>Chrysogena</i>	<i>P. nalgiovense</i> <i>P. chrysogenum</i>	141 (7)	84 (4)
	<i>Roquefortorum</i>	<i>P. roqueforti</i> <i>P. carneum</i> <i>P. paneum</i>	250 (12)	119 (6)
	<i>Aspergilloides</i>	<i>P. corylophilum</i> <i>P. glabrum</i> <i>P. oxalicum</i>	235 (11)	104 (5)

**Table 3:** Comparison of calibration and validation PGP between PLS-DA and linear SVM models

Taxonomic rank	Model name	PLS-DA models			linear SVM models		
		LV	PGP Calibration	PGP Validation	$\gamma$	PGP Calibration	PGP Validation
<b>Subphylum</b>	<i>Micromycetes</i>	7	100	100	0,001	100	100
<b>Class</b>	<i>Peizizomycotina</i>	12	99,9	99,2	0,01	99,8	100
<b>Order</b>	<i>Dothideomycetes</i>	5	100	100	0,01	100	100
	<i>Sordariomycetes</i>	10	100	99,6	0,005	99,3	100
<b>Family</b>	<i>Mucorales</i>	16	100	100	0,005	98,4	100
	<i>Hypocreales</i>	5	100	99,5	0,001	100	100
<b>Genus</b>	<i>Mucoraceae</i>	20	100	100	0,005	99,1	100
	<i>Trichocomaceae</i>	25	99,9	98,8	0,05	95,7	99,9
<b>Subgenus</b>	<i>Penicillium 1</i>	13	100	99,6	0,2	98,7	95,7
<b>Section</b>	<i>Aspergillus 1</i>	20	100	97,7	0,05	99,2	99,5
	<i>Penicillium 2</i>	15	100	95,4	0,2	94,7	84,7
<b>Serial</b>	<i>Fasciculata</i>	5	100	100	0,2	100	100
<b>Species</b>	<i>Nidulantes</i>	6	100	100	0,01	95,8	100
	<i>Aspergillus 2</i>	5	100	85,7	0,8	100	0,0
	<i>Fusarium</i>	20	99,1	95,5	0,2	96,1	94,5
	<i>Mucor</i>	14	100	99,0	0,2	96,6	95,3
	<i>Camemberti</i>	10	100	83,3	0,1	93,3	37,5
	<i>Chrysogena</i>	5	100	86,5	0,1	95,0	75,0
	<i>Roquefortorum</i>	20	100	66,0	0,01	100	59,0
	<i>Aspergilloides</i>	5	100	100	0,01	100	100

**Table 4:** Correlation matrix of McNemar's test, presenting McNemar's value for each pair of chemometrics methods at the species level

Misclassified sample's number (species level)	Tested chemometrics algorithms	<i>linear methods</i>					<i>non-linear methods</i>					
		LDA	FDA	SIMCA	PLS-DA	SVM Linear	QDA	PNN	KNN	SVM RBF	SVM Sigmoid	SVM Polynomial
188	LDA	0	10	703	13	0,2	150	845	71	591	454	882
255	FDA	10	0	581	45	13	87	715	29	475	349	751
1164	SIMCA	703	581	0	840	723	251	9,2	382	6,8	37	14
123	PLS-DA	13	45	840	0	10	239	987	139	722	577	1026
178	SVM Linear	0,2	13	723	10	0	162	865	80	610	472	903
514	QDA	150	87	251	239	162	0	351	16	178	99	378
1316	PNN	845	715	9,2	987	865	351	0	499	32	83	0,6
392	KNN	71	29	382	139	80	16	499	0	293	191	531
1041	SVM RBF	591	475	6,8	722	610	178	32	293	0	12	41
888	SVM Sigmoid	454	349	37	577	472	99	83	191	12	0	97
1356	SVM Polynomial	882	751	14	1026	903	378	0,6	531	41	97	0

Figure 1: Raw and preprocessed FT-IR spectra of an *Alternaria alternata* strain with a tentative band assignment of major macromolecules.

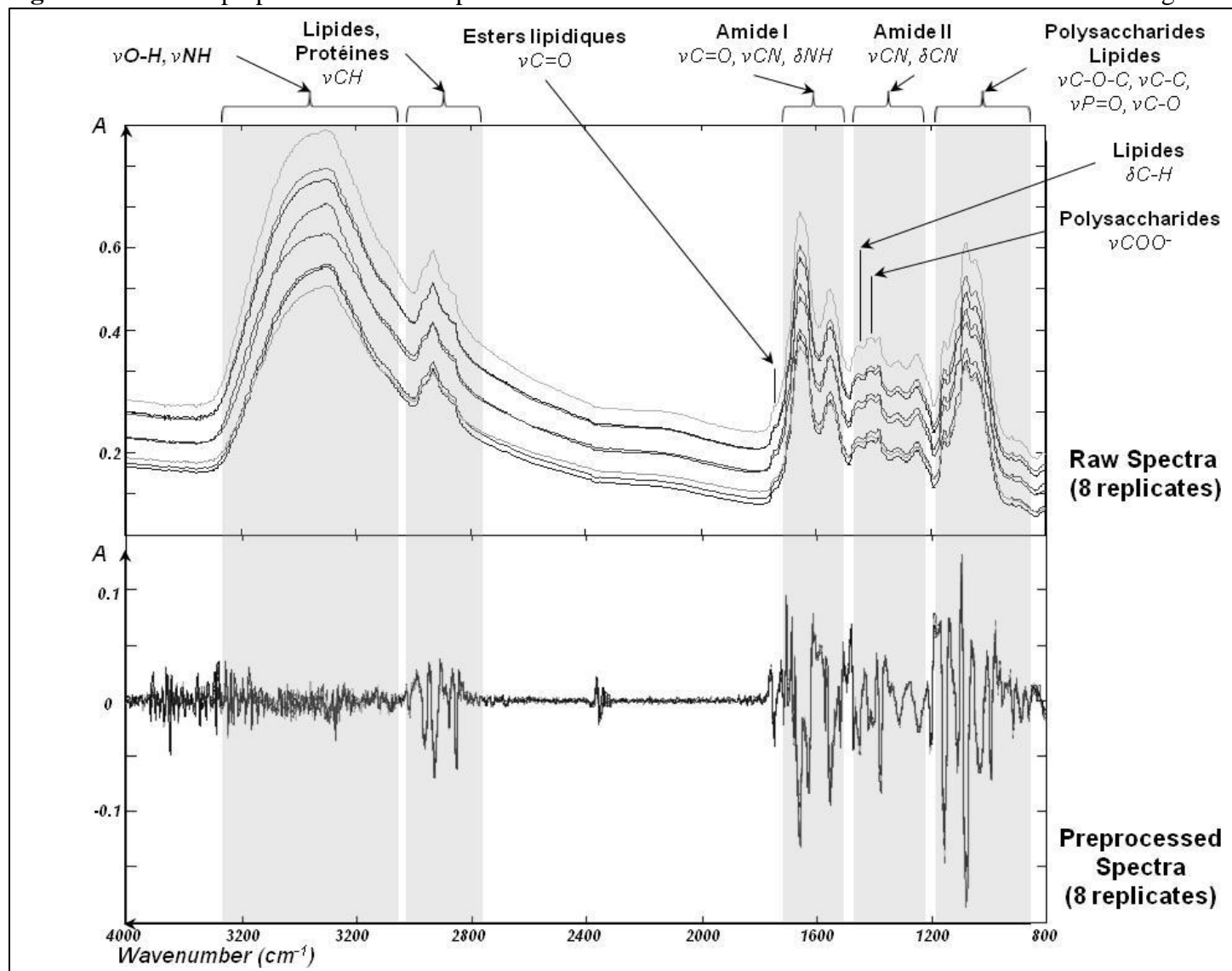
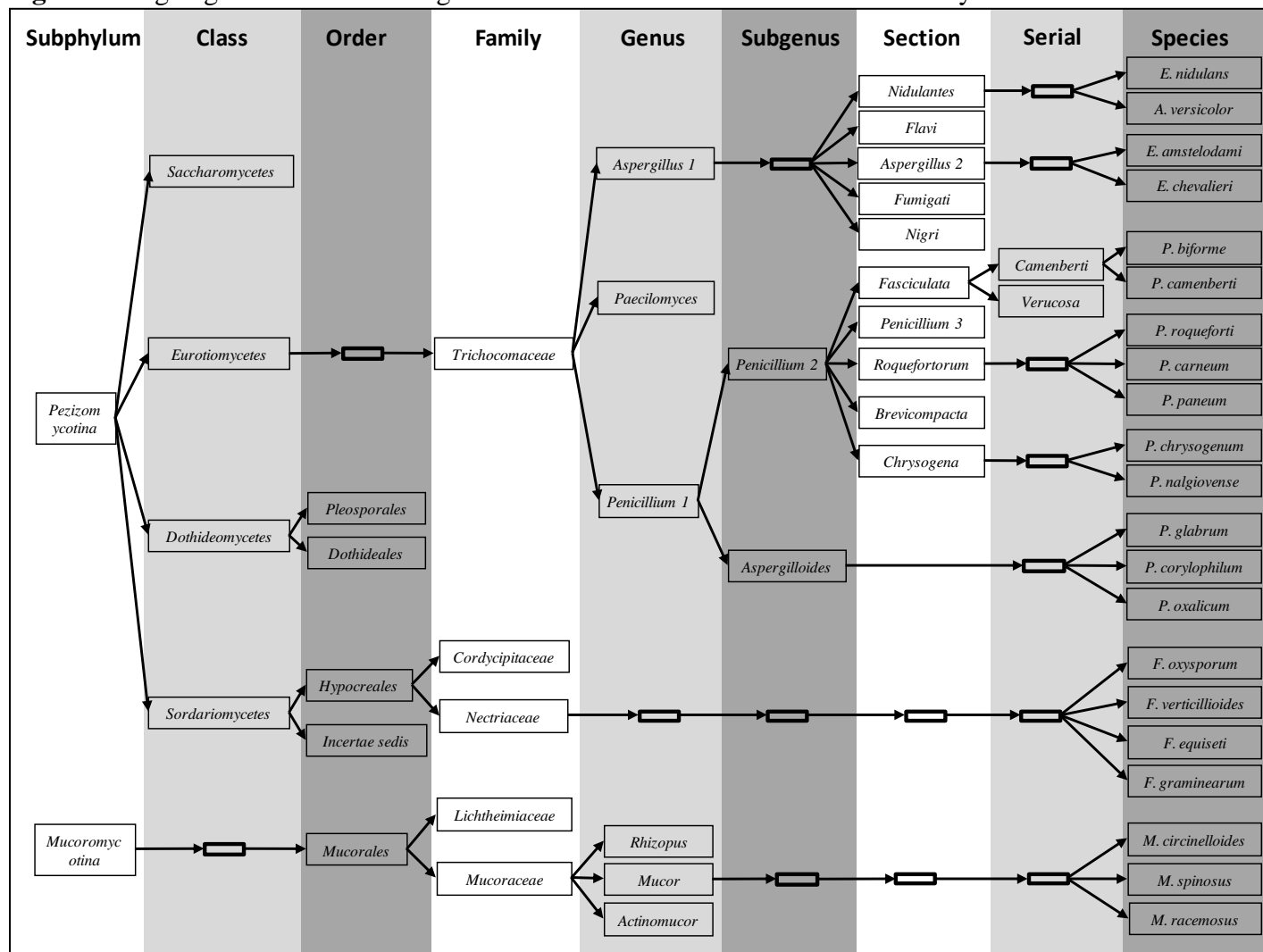
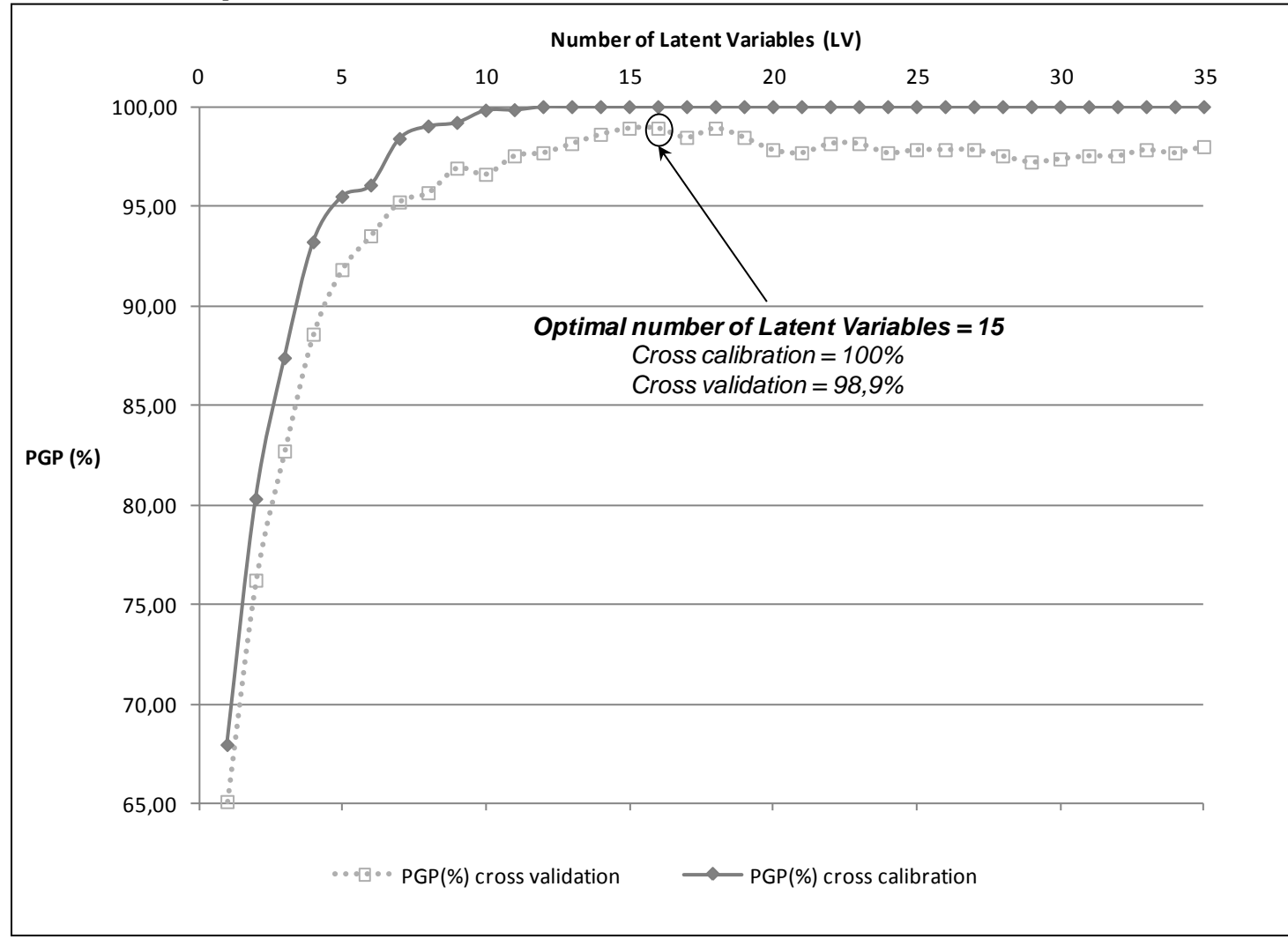


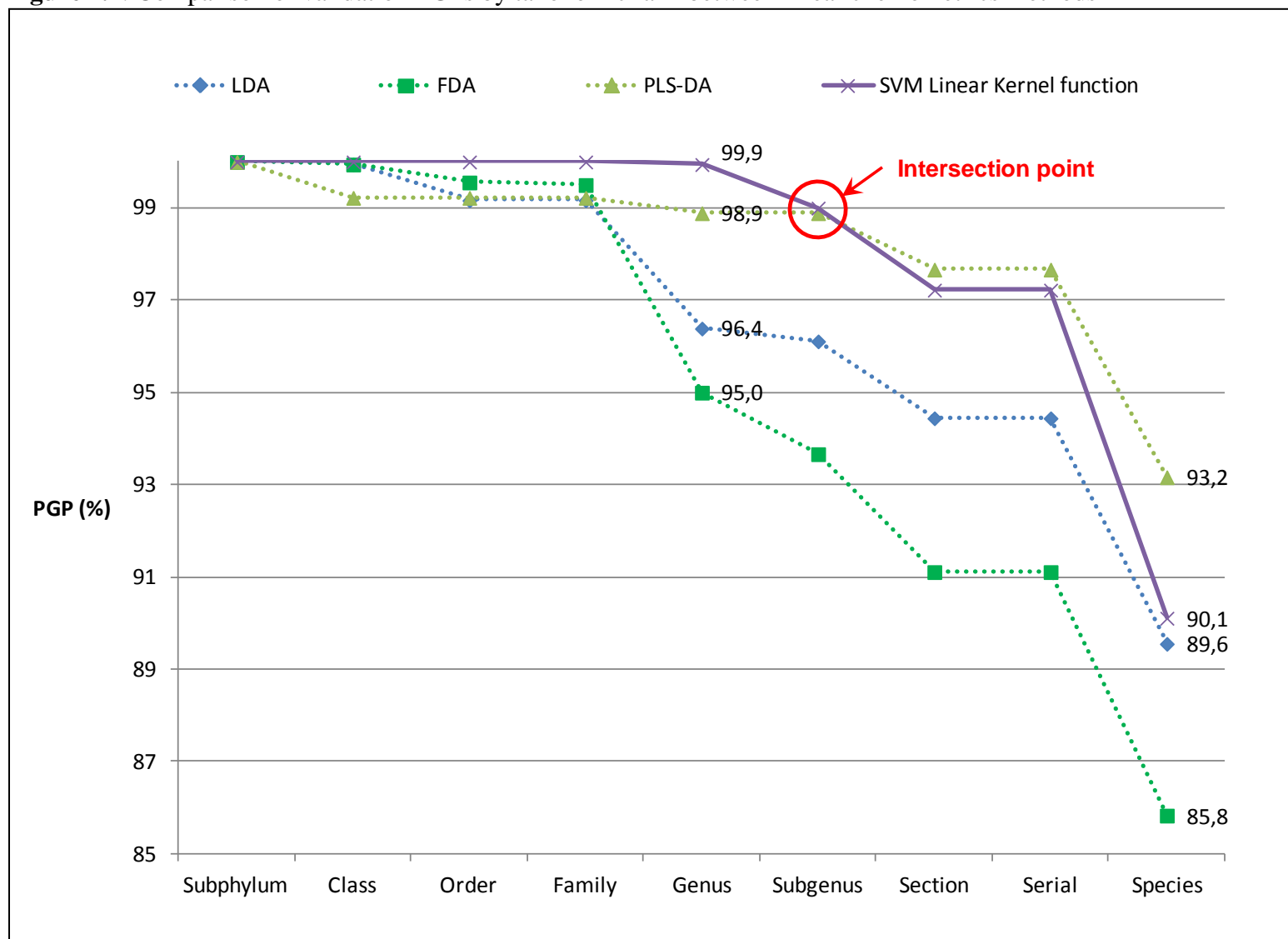
Figure 2: Organigram of the modeling cascade based on the current mold taxonomy



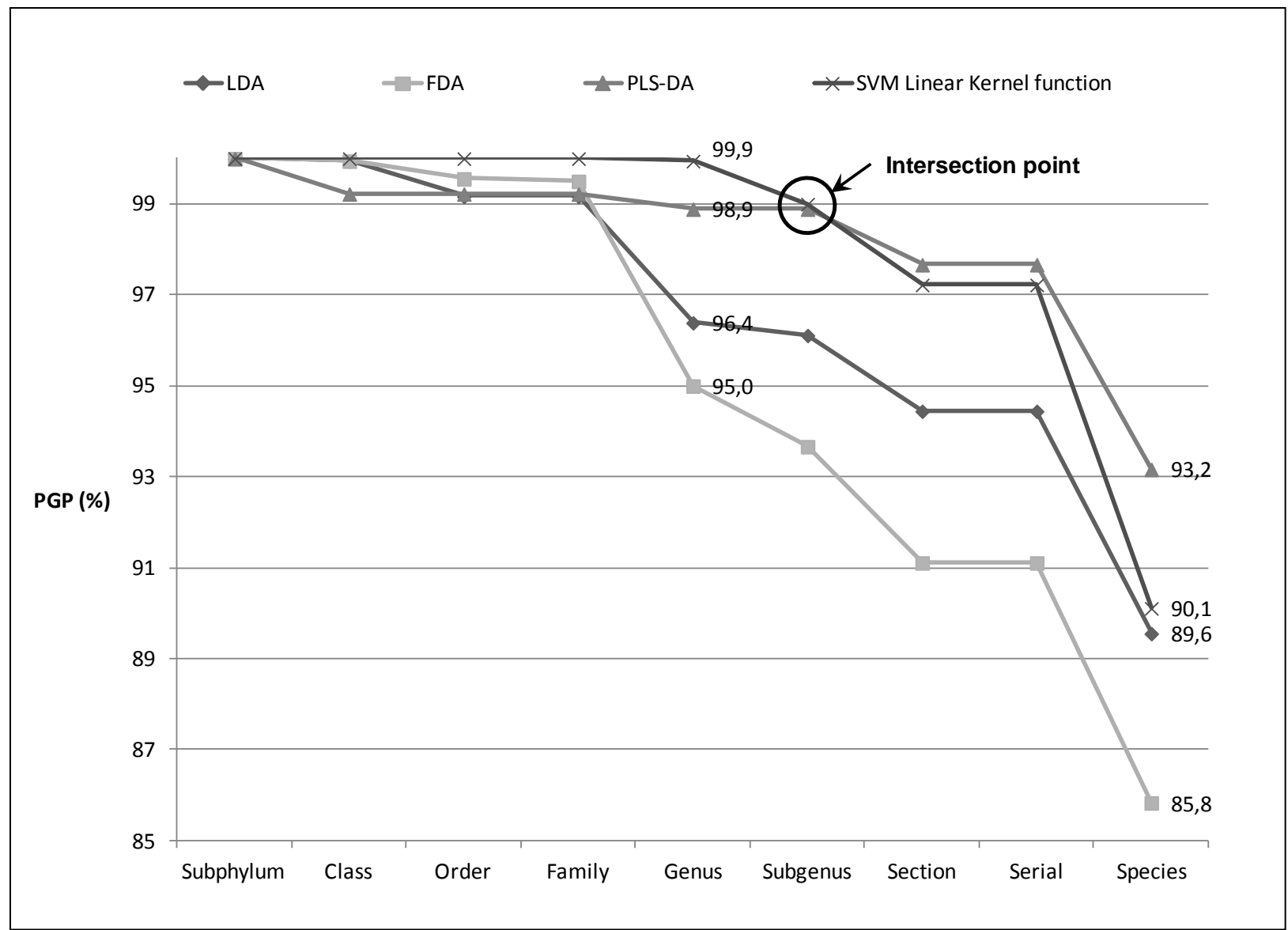
**Figure 3:** Number of Latent Variables (LV) optimization (PLS-DA methods) for the *Mucor* species model with 3 clusters (*M. racemosus*, *M. circinelloides*, *M. spinosus*)



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

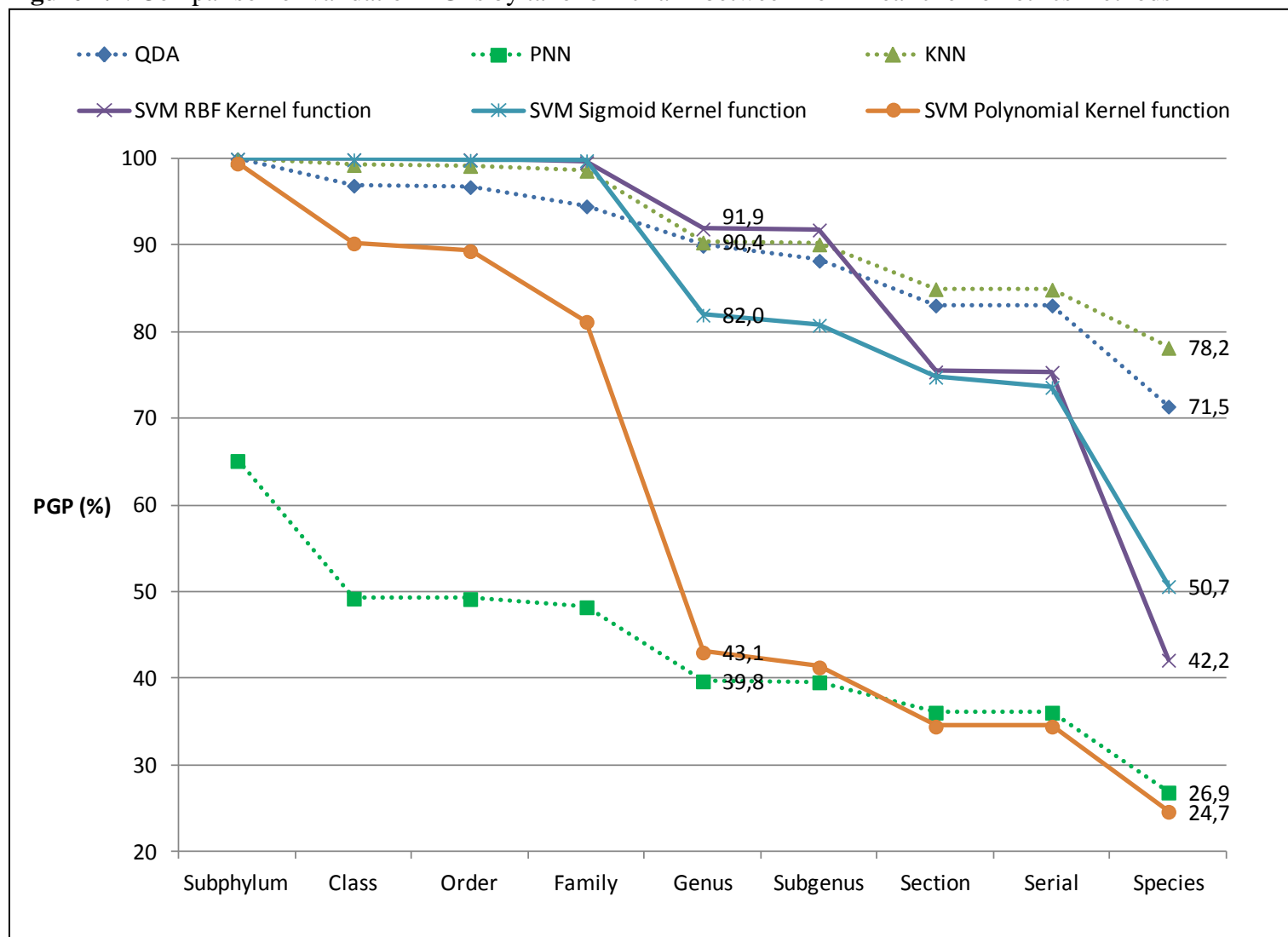
**Figure 4.1:** Comparison of validation PGPs by taxonomic rank between linear chemometrics methods

Reproduced in color on the Web and in black-and-white in print

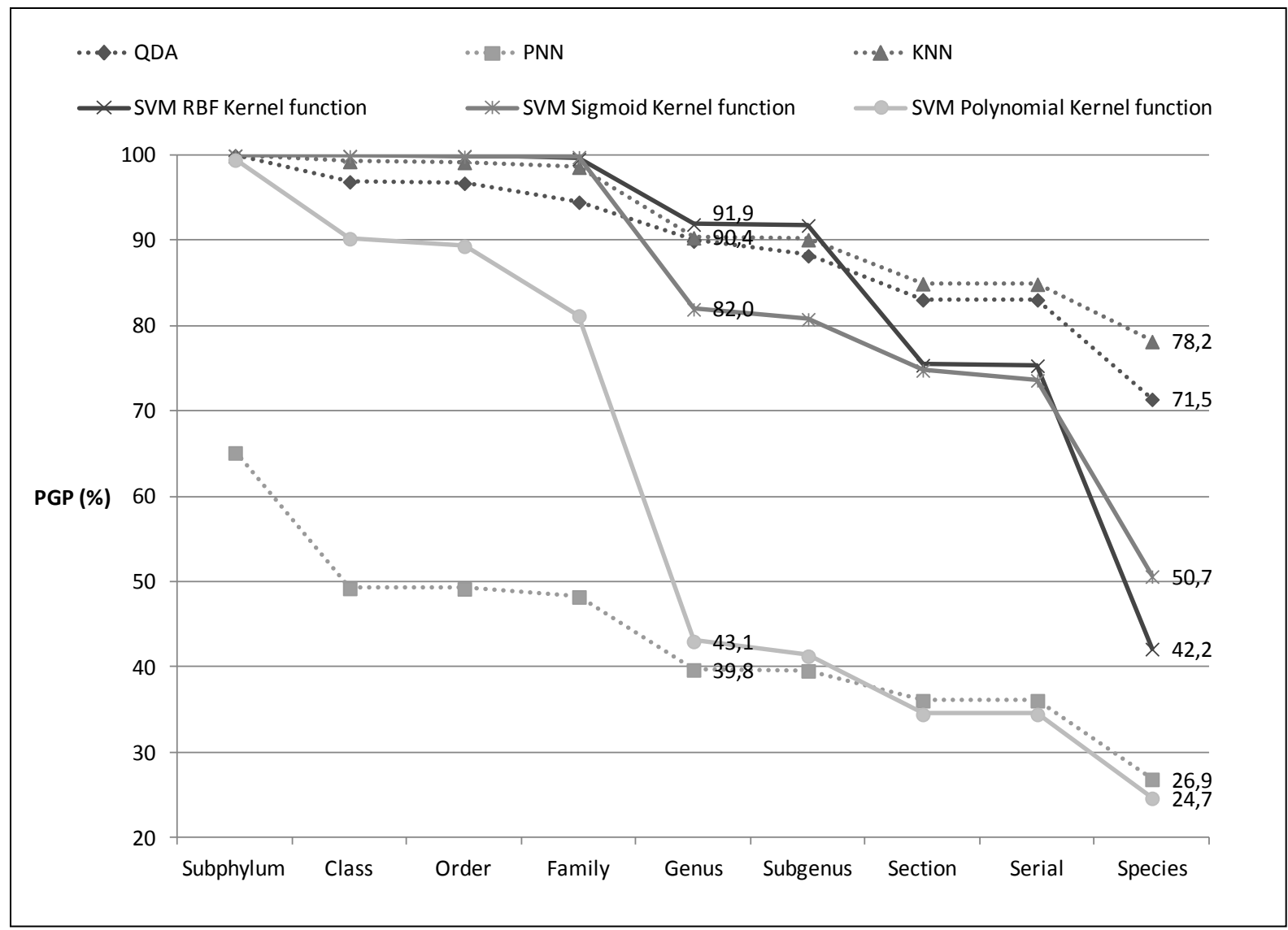


1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47



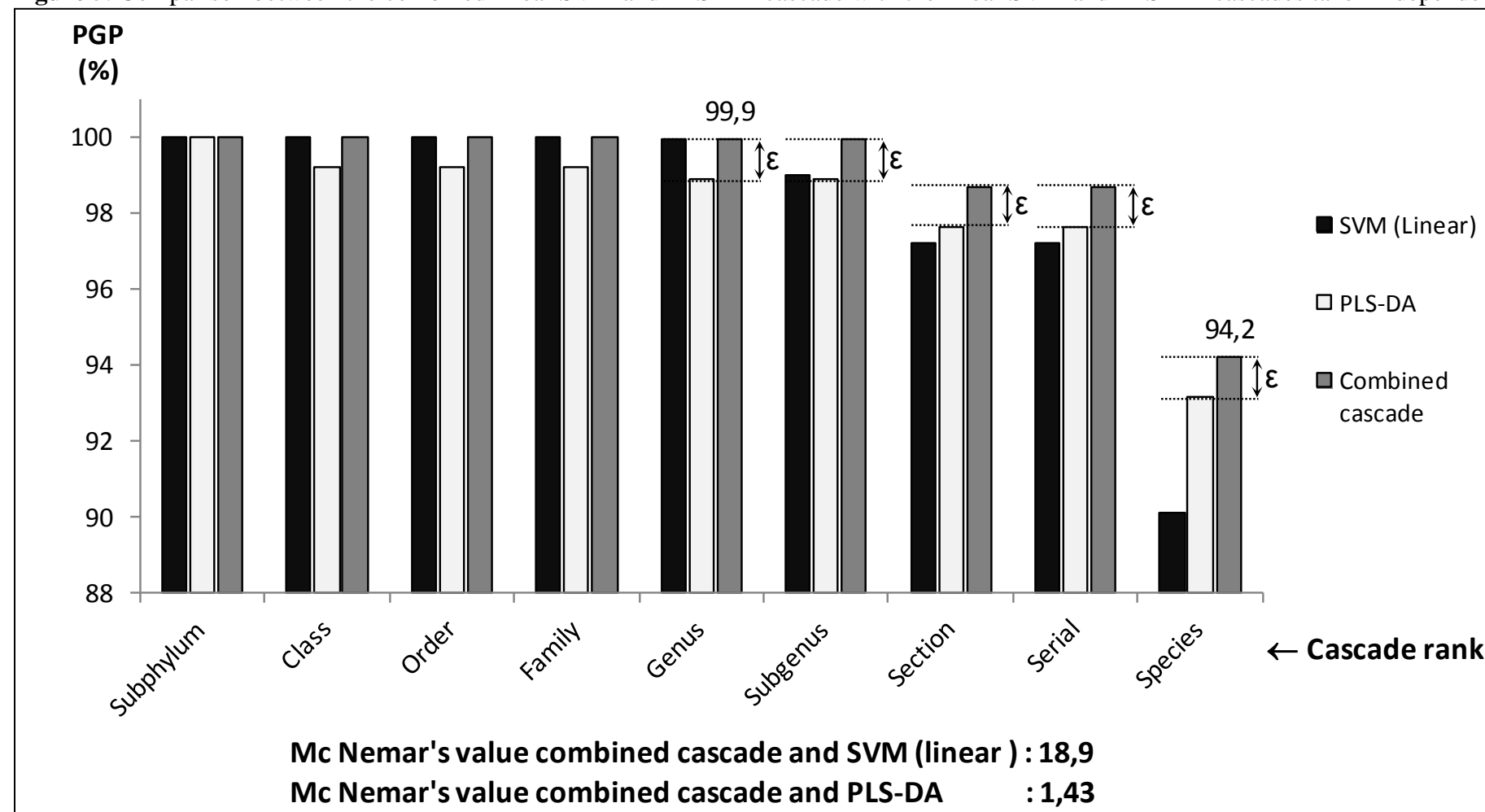
**Figure 4.2:** Comparison of validation PGPs by taxonomic rank between non-linear chemometrics methods

Reproduced in color on the Web and in black-and-white in print



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

**Figure 5:** Comparison between the combined linear SVM and PLS-DA cascade with the linear SVM and PLS-DA cascades taken independently.



## 7 References

- 
- [1] J. Laane, *Frontiers of Molecular Spectroscopy*, 1st ed. Elsevier Science, Texas, USA, **2008**.
- [2] A. Lecellier, J. Mounier, V. Gaydou, L. Castrec, G. Barbier, W. Ablain, M. Manfait, D. Toubas, G.D. Sockalingum, Differentiation and identification of filamentous fungi by high-throughput FTIR spectroscopic analysis of mycelia, *J Food Microbiol.*, 32 (**2014**) 168-169.
- [3] A. Lecellier, V. Gaydou, J. Mounier, L. Castrec, G. Barbier, W. Ablain, M. Manfait, D. Toubas, G.D. Sockalingum, Implementation of an FTIR spectral library of 486 filamentous fungi strains for rapid identification of molds, *J Food Microbiol.*, DOI information: 10.1016/j.fm.2014.01.002. In press (**2014**).
- [4] V. Shapaval, J. Schmitt, T. Moretro, H.P. Suso, I. Skaar, A.W. Asli, D. Lillehaug, A. Kohler, Characterization of food spoilage fungi by FTIR spectroscopy, *J. Appl. Microbiol.* 114 (**2013**) 788-796.
- [5] M. Decker, P.V. Nielsen, H. Martens, Near-infrared spectra of *Penicillium camemberti* strains separated by extended multiplicative signal correction improved prediction of physical and chemical variations, *Appl Spectrosc.* 59 (**2005**) 56-68.
- [6] J.D. Pallua, W. Recheis, R. Poeder and al., Morphological and Tissue Characterization of the Medicinal Fungus *Herichium corraloides* by a Structural and Molecular Imaging Platform. *Analyst* 137 (**2012**) 1584-1595.
- [7] B. R. Kowalski, *Chemometrics : view and proposition*, *J. Chem Inf. Comput. Sci.* 15 (**1975**) 201-203.
- [8] A. Höskuldsson, *Prediction methods in science and technology, Basic Theory vol.1;* Thor Publishing: Copenhagen, Denmark, **1996**, pp. 245.

- 1  
2  
3  
4  
5 [9] D. Bertrand, E. Dufour, Chimiométrie appliquée à la spectroscopie infrarouge, La  
6 spectroscopie infrarouge et ses applications analytiques, 2nd ed. Lavoisier : Paris, **2006**,  
7 pp. 309-401.  
8  
9  
10  
11 [10] R. Fisher, The use of multiple measurements in taxonomic problems, Ann. Eugenics 7  
12 (**1936**) 179-188.  
13  
14  
15 [11] D. H. Moore, Combining linear and quadratic discriminants, Comput. Biomed. Res. 6  
16 (**1973**) 422-429.  
17  
18  
19 [12] Wold, S.; Martens, H.; Wold, H. The multivariate calibration problem in chemistry  
20 solved by the PLS method, In: Ruhe, A., Kastrom, B. (Eds.) Springer, Heidelberg  
21 (**1983**), 286-293.  
22  
23  
24 [13] S. Wold, M. Sjostrom, SIMCA: A method for analyzing chemical data in terms of  
25 similarity and analogy, Chemometrics, 52 (**1977**) 243.  
26  
27  
28 [14] V. Vapnik, A. Lerner, Pattern recognition using generalized portrait method, Automat.  
29 Rem. Contr. 24 (**1963**) 774-780.  
30  
31  
32 [15] J.H. Friedman, F. Baskett, L.J. Shustek, An Algorithm for Finding Nearest Neighbors,  
33 Trans. Comput. 24 (**1975**) 1000-1006.  
34  
35  
36 [16] D. Specht, Probabilistic neural networks, Neural networks 3 (**1990**) 110-118.  
37  
38  
39 [17] D. Bertrand, E. Dufour, Identification et caractérisation des microorganismes, La  
40 spectroscopie infrarouge et ses applications analytiques, 2nd ed. Lavoisier : Paris, **2006**,  
41 pp. 561-581.  
42  
43  
44 [18] L.J. Tashman, Out-of-samples tests of forecasting accuracy : an analysis and review,  
45 Int. J. Forecasting 16 (**2000**) 437-450.  
46  
47  
48 [19] S. Arlot, A. Celisse, A survey of cross-validation procedures for model selection,  
49 Statistics Surveys 4 (**2010**) 40-79.  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

- 1  
2  
3  
4  
5 [20] M. Stone, Cross-Validatory Choice and Assessment of Statistical Predictions, J R Stat  
6 Soc Series B Stat Methodol. 36 (1974) 111-147.  
7  
8  
9 [21] T.G. Dietterich, Approximate Statistical Tests for Comparing Supervised Classification  
10 Learning Algorithms, Department of Computer Science, Oregon State University, 10  
11 (1998) 1895-1923.  
12  
13  
14 [22] J. Workman, Review of Chemometrics Applied to Spectroscopy: Quantitative and  
15 Qualitative Analysis, The Handbook of Organic Compounds, NIR, IR, Raman, and UV-  
16 Vis Spectra Featuring Polymers and Surfactants 1 (2001) 301-326.  
17  
18  
19 [23] F. Chauchard, R. Cogdill, S. Roussel, J.M. Roger, V. Bellon-Maurel, Application of LS-  
20 SVM to non-linear phenomena in NIR spectroscopy : Development of a robust and  
21 portable sensor for acidity prediction in grapes, Chemometr. Intell. Lab. 71 (2004) 141-  
22 150.  
23  
24 [24] M. Mörtzell, M. Gulliksson, An overview of some non-linear techniques in  
25 Chemometrics, Rapportserie FSCN - ISSN 1650-5387 2001:6, Mid-Sweden University,  
26 2001.  
27  
28  
29 [25] D.L. Swets, J. Weng, Using discriminant eigenfeatures for image retrieval, Pattern Anal  
30 and Machine Intelligence, 18 (1996) 831-836.  
31  
32  
33 [26] Romeder, J. M. Méthodes et Programmes d'Analyse Discriminante, Dunod, Paris,  
34 France, 1973.  
35  
36  
37 [27] Sharaf, M. A.; Illman, D. L.; Kowalski B.R. Chemometrics, Wiley, New York, 1986.  
38  
39  
40 [28] P.H. Garthwaite, An interpretation of partial least squares, J. Amer. Statist. Assoc. 89  
41 (1994) 122-127.  
42  
43  
44 [29] M. Tenenhaus, L'algorithme de régression PLS1, In Tenenhaus M. (ed.), Paris, France  
45 1998 pp. 75-77.  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

- 1  
2  
3  
4  
5 [30] T. Hastie, R. Tibshirani, J. Friedman, The Elements of Statistical Learning, Springer-  
6 Verlag, New York, **2001**.  
7  
8  
9 [31] D.O. Loftsgaarden, C.P. Quesenberry, A Nonparametric Estimate of a Multivariate  
10 Density Function, Ann. Math Statist. 36 (**1965**) 1049-1051.  
11  
12 [32] L. Labart, A. Morineau, N. Tabart, Technique de la description statistique, Méthodes et  
13 logiciels pour l'analyse des grands tableaux, Ed. Dunod, Paris, France, **1987**.  
14  
15 [33] P. Wasserman, Advanced methods in neural networks, Van Nostrand Reinhold, New  
16 York, USA, **1993**.  
17  
18 [34] C.C. Chang, C.J. Lin, T. ACM, LIBSVM : A Library for Support Vector Machines, Int.  
19 Sys. Techn, 2(27) (**2011**) 1-27, Software available at  
20 <http://www.csie.ntu.edu.tw/~cjlin/libsvm>.  
21  
22 [35] [http://www.chimietrie.fr/saisir\\_conceptors.html](http://www.chimietrie.fr/saisir_conceptors.html).  
23  
24 [36] M. Boysen, P. Skouboe, J. Frisvad, L. Rossen, Reclassification of the Penicillium  
25 roqueforti group into three species on the basis of molecular genetic and biochemical  
26 profiles, Microbiology, 142 (**1996**) 541-9.  
27  
28 [37] F. Giraud, T. Giraud, G. Aguilera, E. Fournier, R. Samson, C. Cruaud, S. Lacoste, J.  
29 Ropars, A. Tellier, J. Dupont, Microsatellite loci to recognize species for the cheese  
30 starter and contaminating strains associated with cheese manufacturing, Int. J. Food  
31 Microbiol. 137 (**2010**) 204-13.  
32  
33 [38] V. Hubka, M. Kolarik, A. Kubatova, S.W. Peterson, Taxonomic revision of Eurotium  
34 and transfer of species to Aspergillus, Mycologia 105 (**2013**) 912-37.  
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