



HAL
open science

Combination of multivariate curve resolution with factorial discriminant analysis for the detection of grapevine diseases using hyperspectral imaging. A case study: flavescence dorée

Sílvia Mas Garcia, Maxime Ryckewaert, Florent Abdelghafour, Maxime Metz, Daniel Moura, Carole Feilhes, Fanny Prezman, Ryad Bendoula

► To cite this version:

Sílvia Mas Garcia, Maxime Ryckewaert, Florent Abdelghafour, Maxime Metz, Daniel Moura, et al.. Combination of multivariate curve resolution with factorial discriminant analysis for the detection of grapevine diseases using hyperspectral imaging. A case study: flavescence dorée. *Analyst*, 2021, 146 (24), pp.7730-7739. 10.1039/D1AN01735G . hal-03538618

HAL Id: hal-03538618

<https://hal.inrae.fr/hal-03538618>

Submitted on 21 Jan 2022

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

1 **Combination of Multivariate Curve Resolution with Factorial Discriminant Analysis for**
2 **the detection of grapevine diseases using Hyperspectral imaging. A case study:**
3 **Flavescence Dorée.**

4 Silvia Mas Garcia^{1,2*}, Maxime Ryckewaert^{1,2}, Florent Abdelghafour¹, Maxime Metz^{1,2}, Daniel
5 Moura¹, Carole Feilhes³, Fanny Prezman³, Ryad Bendoula¹.

6 ¹ITAP, INRAE, Institut Agro, University Montpellier, 34196 Montpellier, France

7 ²ChemHouse Research Group, 34196 Montpellier, France

8 ³IFV, 1920 Route de Lisle-sur-Tarn, 81310 Peyrole, France

9

10 * Author for correspondence:

11 E-mail: silvia.mas-garcia@inrae.fr

12

13

14

15

16

17

18

19

20

21

22

23 **Abstract**

24 Hyperspectral imaging is an emergent technique in viticulture that can potentially detect
25 bacterial diseases in a non-destructive manner. However, the main problem is to handle the
26 substantial amount of information obtained from this type of data, for which reliable data
27 analysis tools are necessary. In this work, combination of multivariate curve resolution-
28 alternating least squares (MCR-ALS) and factorial discriminant analysis (FDA) is proposed to
29 detect the Flavescence dorée grapevine disease from hyperspectral imaging.

30 The main purpose of MCR-ALS in this work was providing chemically meaningful basic
31 spectral signatures and distribution maps of the constituents needed to describe both healthy
32 and infected images by Flavescence dorée. MCR scores (distribution maps) were used as
33 starting information for FDA to distinguish between healthy and infected pixels/images. Such
34 an approach is presumably more powerful than the direct use of FDA on the raw imaging data,
35 since MCR scores are compressed and noise-filtered information on pixel properties, which
36 makes them more suitable for discrimination analysis. High levels of correct pixels
37 discrimination rates (CR=85,1%) for the MCR-ALS/FDA discrimination model were obtained.
38 The model present a lesser ability to determine infected leaves than healthy leaves.
39 Nevertheless, only two images were misclassified. Therefore, proposed strategy constitutes a
40 good approach for the detection of the Flavescence dorée that could be potentially used to detect
41 other phytopathologies

42

43 **Keywords:** Multivariate Curve Resolution-Alternating Least Squares (MCR-ALS), Factorial
44 Discriminant Analysis (FDA), Hyperspectral imaging, vineyard diseases, Flavescence dorée

45

46

47 **1. Introduction**

48 Epidemiological surveillance is a crucial issue in agriculture and especially in viticulture. As a
49 matter of fact, the grapevine (*Vitis vinifera*) is sensitive to a wide range of biopests. To cope
50 with these threats, preventive chemical control is required. To reduce the use of chemical inputs
51 while ensuring the protection of the vineyard, it is necessary to implement more parsimonious
52 spraying practices. The development of sustainable crop protection systems is closely related
53 to the knowledge regarding the physiological state and the health status of the vineyard ¹.

54 To date, the evaluation of sanitary risks is conducted by visual and tactile inspection, which is
55 time consuming and labour intensive. The analysis of light-matter interaction can provide
56 information related to physiological properties such as hydric status, nitrogen content,
57 pigmentation or even cellular structure ². Therefore, optical instruments and especially
58 multispectral (MSI) and hyperspectral imaging (HSI) are relevant tools for the automated and
59 non-invasive detection of phytopathology ³⁻⁶. In this context, conventional analysis of
60 hyperspectral and multispectral images, such as determination of spectral vegetation indices
61 (SVIs) ⁷, performs only a limited use of the substantial amount of information available with
62 this type of data. Therefore, in order to successfully interpret these images, the application of
63 advanced data processing tools is necessary. In this work, we will focus on the application of
64 HSI to discriminate an important vine disease: the Flavescence dorée” (FD, also known as
65 “yellowing”).

66 FD is a phytopathology caused by the bacteria *Candidatus phytoplasma vitis* that can spread
67 fast through a leafhopper (*Scaphoideus titanus*). It represents a very serious threat, since
68 without proper management; it can lead to the complete loss of the harvest or even the death of
69 the vine stocks. Recently, spectral imaging have been used to detect FD. Albetis et al. ⁶
70 evaluated the potentiality of Unmanned Aerial Vehicle (UAV) multispectral imagery for the
71 airborne detection of FD symptoms under field conditions. For this purpose, they analysed

72 several spectral bands, vegetation indices, and biophysical parameters. However, the specific
73 detection of FD appears to be limited. Al-Saddik et al. ⁵ used a portable spectroradiometer (350–
74 2500 nm) to collect hyperspectral reflectance data of healthy and symptomatic leaves. The aim
75 of this study was to develop specific spectral disease indices (SDIs) for the detection of FD
76 disease in grapevines, thereby, reaching discrimination accuracies of more than 90%. However,
77 the SDIs were dependent on the disease infestation state and the grapevine variety considered;
78 the best wavelengths selected were different from one case to another, and hence no single best
79 index for FD in all situations was identified. To deal with these limitations, this work aims to
80 propose a new general methodology (i.e. not depending of the variety) to discriminate between
81 healthy and infected leaves based on HSI measurements and data analysis methods.

82 The data analysis workflow proposed in this work relies on two steps:

83 a) Multivariate curve resolution-Alternating Least Squares (MCR-ALS) ^{8,9} model to
84 provide chemically meaningful spectral signatures and related distribution maps of the
85 image constituents. This unmixing method allows a global differentiation between
86 infected and healthy images. However, some components related to the two different
87 class types (infected and healthy pixels) may overlap, and hence, a supervised
88 discrimination method is necessary to achieve a harder separation between them.

89 b) Factorial discriminant analysis (FDA) ¹⁰ model on MCR scores (distribution maps).
90 This supervised classifier will help to discriminate between infected and healthy images,
91 using previous pixel labelling (classes infected or healthy).

92 Previous studies have already shown the capability of the application of MCR-ALS combined
93 with supervised classification methods to the analysis of imaging data ^{11,12}. These works
94 demonstrated that the use of MCR outputs as starting information for classification methods
95 allows a compound-wise selection and preprocessing of the input information to be submitted
96 to the classification algorithm. This is due to the fact that MCR results are chemically

97 meaningful and express concentrations or spectra of the pixel constituents in the images. Such
98 a specificity of the method allows discarding components related to background signal
99 contributions in the classification task. In this work, FDA is chosen as the supervised classifier
100 because it is one of the simplest and fastest approach for discrimination that has proven its
101 efficiency for various analytical chemistry applications ^{13,14}. However, to the best of our
102 knowledge this is the first time that MCR-ALS combined with such a classification method is
103 used as phytopathology detection model. This study demonstrates that the proposed
104 methodology has the potential to improve disease detection in agriculture applications.

105 **2. Material and Methods**

106 **2.1 Samples**

107 Leaves were collected during September 2020 on previously identified plots with Flavescence
108 Dorée. All cultivars were sampled with a similar proportion of red and white varieties. In total
109 109 leaves were collected on the field. The number of leaves from the different varieties
110 selected for this study are summarised in Table 1.

111

112 **TABLE 1**

113

114 Infected leaves were chosen in order to represent at best the variability of the available
115 symptoms in terms of severity and stages of infections. Leaves were selected when foliar
116 symptoms were undoubtedly caused by FD from vines exhibiting clear symptoms of FD on
117 other organs. Each leaf and each vine from which they were extracted were diagnosed by a
118 phytopathology expert. Leaves were extracted from the front face, in the middle of the canopy
119 so that to avoid the younger and older organs which can present different physiological

120 behaviour. Regarding the healthy leaves, they were selected in the same regions and they were
121 asserted absent of symptoms of FD or any other visible pathology. However, some of the
122 healthy sample can exhibit light forms of mechanical or chemical wounds (due to protection,
123 management operations) and some slight damage caused by insects.

124 **2.2 Image acquisition**

125 Acquisitions of leaf images were performed with a hyperspectral camera (IQ, Specim, Finland).
126 Imaging of grapevine leaves was carried out in the spectral range of 400-900 nm, with a spectral
127 resolution of 7 nm. Images in RGB were also registered. Illumination was provided by a
128 halogen lamp (Arrilite 750 Plus ARRI, Munich, Germany). Constant angles of -50° and 50°
129 were maintained between the halogen lamp and the hyperspectral camera. These angles were
130 chosen to optimise the intensity of the reflected beam and to reduce specular reflection.

131 For each sample image, the intensity of the reflected light ($I(\lambda)$) was measured. The Dark current
132 ($I_n(\lambda)$) *i.e.* signal without light, was recorded from all measured spectra and then subtracted. A
133 white reference (SRS99, Spectralon®) ($I_0(\lambda)$) was measured to standardise spectra and prevent
134 nonlinearities of all the instrumentation components (light source, lens, fibbers and
135 spectrometer). From these measurements, a reflectance image ($R(\lambda)$) was calculated for each
136 sample, as follows:

$$137 \quad R(\lambda) = \frac{I(\lambda) - I_n(\lambda)}{I_0(\lambda) - I_n(\lambda)} \quad \text{Equation 1}$$

138 **3. Data analysis**

139 The proposed workflow for data analysis follows the three following steps:

140 a) **Image preprocessing**

141 b) **MCR** to recover basic spectral signatures and distribution maps of pure compounds
142 contributions, allowing differentiation between infected and healthy images.

143 c) **FDA model using the MCR scores (concentration profiles) resulting from the MCR**
144 **results** to predict the class (infected or healthy) of the images.

145 These steps are described in detail in the following subsections

146 **3.1 Image preprocessing**

147 In HSI (Hyperspectral Imaging), the generated data can be arranged into a data cube in which
148 the x-and y-axis correspond to the pixel coordinates and the z-axis corresponds to the
149 wavelengths values registered in each pixel. Data preprocessing is required to improve the
150 signal quality and to compress the acquired raw data for further analysis.

151 Firstly, the pixels in the images were binned by a factor of 4 in x and y . This spatial binning
152 produced an image of 128x128 pixels from an original image of 512x512. Afterwards, a mask
153 was create for each image to extract only the vegetation pixels. The Spectral Angle Mapper
154 (SAM) ¹⁵ was used for this purpose. To identify vegetation pixels, SAM compare image spectra
155 to a reference spectrum by calculating the spectral angle between them. Smaller angles
156 represent closer matches to the reference spectrum, and hence the corresponding pixels are
157 classified as vegetation pixels, whereas pixels further away than the specified maximum angle
158 threshold are not classified.

159 Finally, A matrix \mathbf{Di} (n,m) of dimension n equal to $(x \times y)$ pixels by 175 wavelengths was
160 generated per each image.

161 **3.2 Multivariate Curve Resolution/ Factorial Discrimination Analysis (MCR/FDA) model**

162 Before using MCR-ALS and FDA methods, the dataset was divided into two sets of samples:
163 training and independent test sets. The training sets were used to build the models. The test sets
164 were left for external validation and are not used to build the models. Healthy and infected
165 images were both divided with the same split ratio of 2/3 and 1/3 respectively for training and

166 test, as detailed in Table 1. This division was made randomly and assuring a similar distribution
167 of all classes in both training and test sets.

168 **3.2.1 Multivariate Curve Resolution- Alternating Least Squares (MCR-ALS)**

169 The goal of the MCR-ALS algorithm is the decomposition of the image data **D** into distribution
170 maps (relative amounts or concentration) and pure spectra of the constituents present in the
171 imaged sample^{8,9,16}. In matrix form, the hyperspectral images can be described by a bilinear
172 model based on the Beer-Lambert law (Equation 1). Where the matrix **D** contains the pixel
173 spectra obtained after the preprocessing described in section XX. Each spectra is then
174 decomposed into a set of concentration profiles (**C** matrix) corresponding to pure spectra
175 (**S^T** matrix) of the constituents present in the image. **E** is the matrix associated with noise or
176 experimental error (residuals).

$$177 \mathbf{D} = \mathbf{CS}^T + \mathbf{E} \quad \text{Equation 2}$$

178 Figure 1 shows the application of MCR-ALS to an individual image data **D**. It can be observed
179 that every row of the resolved **S^T** matrix corresponds to the pure spectrum of an image
180 constituent, while every column of the resolved **C** matrix of concentration profiles corresponds
181 to the related pixel-to-pixel variation of its chemical concentration. It is worth mentioning that
182 each column of the resolved **C** matrix can be refolded appropriately in order to recover the
183 original two-dimensional spatial image structure and then pure distribution maps are obtained.

184

185 **FIGURE 1**

186

187 In order to recover the bilinear model expressed in Equation 1, MCR-ALS begins with
188 determining the number of signal contributions in the original data set **D** by Singular Value

189 Decomposition (SVD) ¹⁷. Afterwards, an initial **C** or **S^T** matrix with as many profiles as the
 190 number of components estimated for **D** is constructed to initiate the iterative resolution process.
 191 In this work, the initial **S^T** was generated by a pure variable selection method based on Simple-
 192 to-use Interactive Self-modelling Mixture Analysis (SIMPLISMA) ¹⁸. Such estimate **S^T** and
 193 the matrix **D** are used to initialise the least squares alternating optimisation of the profiles in
 194 matrices **C** and **S^T** of the bilinear model under the constraints until convergence is achieved.
 195 The convergence criterion can be a maximum number of iterations or a value related to the
 196 difference in fit improvement between consecutive iterations.

197 The quality of the MCR results are described by the explained variance (% r^2), which are
 198 calculated according to the following expressions:

$$199 \quad \% r^2 = 100 \times \left(1 - \frac{\sum e_{ij}^2}{\sum d_{ij}^2} \right) \quad \text{Equation 3}$$

200 where e_{ij} is equal to $d_{ij} - d_{ij}^*$, d_{ij}^* are the values of the data set reproduced by the bilinear model
 201 and d_{ij} the original values in the original data set **D**. In order to consider that MCR results of
 202 an analysis are adequate, the variance explained must be sufficiently high and the concentration
 203 profiles and spectra obtained must be chemically meaningful and show shapes consistent with
 204 the variation in the original data sets.

205 MCR-ALS can also be used to analyse simultaneously several images in a single multiset
 206 structure to provide more reliable results ^{9,19}. Resolved features would define much better
 207 general traits analysed together than if they were analysed individually. In this study, the
 208 multiset structures were obtained by setting different images **Di** one on top of each other to
 209 form a column-wise augmented matrix **Daug**. The bilinear model in Equation 1 is now extended
 210 to the augmented data set as shown in Equation 4:

$$211 \quad \mathbf{Daug} = [\mathbf{D}_1; \mathbf{D}_2; \dots; \mathbf{D}_n] = [\mathbf{C}_1; \mathbf{C}_2; \dots; \mathbf{C}_n] \mathbf{S}^T + [\mathbf{E}_1; \mathbf{E}_2; \dots; \mathbf{E}_n] = \mathbf{Caug} \mathbf{S}^T + \mathbf{Eaug} \quad \text{Equation 4}$$

212 where \mathbf{C}_{aug} is a column-wise augmented matrix formed by as many submatrices \mathbf{C}_i as images
213 in the multiset, and \mathbf{S}^T is a single data matrix of pure spectra, assumed to be common and valid
214 for all the images in the multiset. The concentration profiles in each of these submatrices can
215 be also refolded conveniently to recover the related distribution maps of each image (see Figure
216 1b).

217 The MCR-ALS analysis of a single image or an image multiset takes the benefit of the use of
218 constraints on \mathbf{C} or/and \mathbf{S}^T to obtain chemically meaningful and more accurate spectral
219 signatures and distribution maps. In this study, the most common constraints in image
220 resolution, such as non-negativity and normalisation, were used. Moreover, the constraint of
221 correspondence among species to encode the information related to the presence/absence of
222 some components in the different \mathbf{C}_i submatrices in the multiset structures was also applied^{9,16}.

223 MCR-ALS distribution maps (\mathbf{C} matrix) and pure spectra (\mathbf{S}^T matrix) are excellent low
224 dimension, noise-filtered meaningful basis of the pixel and the spectral space of the image,
225 which may be further used to obtain additional information. In this work, the MCR scores
226 (distribution maps) were fed into the FDA to predict the type-class (healthy or infected) of the
227 images.

228 It is worth mentioning that a multiset structure containing all the training dataset from both
229 healthy and infected images (\mathbf{D}_{trFH}) was used for the MCR approach. Then, the distribution
230 maps related to the multiset structure containing all the test dataset (\mathbf{D}_{testFH}) were calculated
231 by a single non-negative least-squares step taking MCR pure spectra obtained in the training
232 stage (\mathbf{S}^{Ttr}).

233 3.2.2 Factorial Discrimination Analysis (FDA)

234 The aim of FDA¹⁰ is to predict the membership of an individual to a group of samples according
235 to pre-defined groups. This method searches for relationships between a qualitative variable

236 (healthy or infected) and a group of quantitative explanatory variables (wavelengths,
237 intensities...). The use of the qualitative variable within a population allows the division of this
238 population into different groups, with each individual assigned to one group. Discrimination
239 of the groups consists of maximising the variance between their gravity center. For each group,
240 the distance from the different gravity center of the groups is calculated and then, the sample is
241 assigned to the group where its distance between the centre of gravity is the nearest. Comparison
242 of the assigned group to the real group is an indicator of the quality of the model, and hence,
243 discrimination rate (CR) is taken as a criterion of goodness for the developed model

244 In this work, FDA was performed to determine the affiliation of each pixel/image whether to
245 the healthy or to the infected class. High correlations can occurred among the wavelengths or
246 intensities of the pixels/images, therefore, MCR scores (distribution maps) coming from the
247 augmented **CtrFH** matrix obtained by MCR-ALS have been used as the pixel input information
248 for FDA. Therefore, no variable reduction algorithm such as PCA or ICA *need to be done due*
249 *to the fact that MCR scores (concentration profiles) are compressed and noise-filtered*
250 *information on pixel properties*. The gravity centre of each sample type in the model was
251 calculated from these training sample scores. The Mahalanobis distance ²⁰ from each to each
252 level of the gravity centres was measured. Finally, test samples were assigned to the group
253 with the nearest gravity centre.

254 **3.3 Software**

255 All data processing has been performed in MATLAB platform (Version 2015b, MathWorks
256 Inc., Natick, MA, USA). The application of MCR-ALS has been performed using the MCR
257 GUI (multivariate curve resolution graphical user interface) developed by the chemometrics
258 group of Universitat de Barcelona and IDAEA-CSIC ²¹, which is can be downloaded from the
259 MCR webpage <http://www.mcrals.info/>. FDA analysis method has been applied using in-house
260 routines, partly based on the PLS Toolbox (Eigenvector Research Inc., Manson, WA, USA).

261 **4. Results and discussion**

262 **4.1. MCR. Global differentiation between infected and healthy images**

263 The first MCR-ALS analysis was focused on identifying significant contributions with a
264 specific reflectance signature for each leaf type (healthy and infected). For that purpose, two
265 multisets were built, one formed by the 47 training images corresponding to the infected leaves
266 from all varieties (**DaugtrF**), and the other multiset formed by the 25 training images
267 corresponding to the healthy leaves from all varieties (**DaugtrH**). MCR-ALS was applied
268 separately to each of these multiset structures using non-negativity constraints in concentration
269 and spectra profiles and spectra normalisation.

270 Table 2 summarises the number of resolved components and the explained variance obtained
271 from the MCR-ALS analyses of both multisets. Resolution of three contributions was necessary
272 in both cases. The inclusion of a different number of contributions gave solutions worse
273 mathematically or unreliable spectra or distribution maps.

274

275 **TABLE 2**

276

277 Figure 2a and b show the MCR-ALS resolved distribution maps (with their corresponding RGB
278 images) and pure spectra of each analysed multiset, respectively. To simplify, resolved
279 distribution maps of only one image per variety is shown. It can be seen that the blue and red
280 contributions present resolved pure spectra rather similar in both multisets, with a Pearson
281 correlation coefficient higher than 0.90. The blue contribution shows a low intensity plateau in
282 visible region (from 400 to 700 nm) and then an increment of the intensity in the near infrared
283 region that ends rather stable from 750 to 900 nm. This contribution seems to present the typical

284 profile related to the cell structure of the leaf. The red contribution presents a peak at 550 nm
285 and a low intensity between 600 and 640 nm, which could correspond to the pigment content,
286 especially anthocyanins and chlorophyll. The green contribution presents a greater spectral
287 dissimilarity between the two multisets. Remarkably, the component from the infected leaf
288 multiset (Figure 2a) has a characteristic peak located at 700 nm, a second peak located at 650
289 nm and lower intensity values at 400 nm and 500 nm. This green contribution could be
290 attributed to a difference in slope level in the red-edge region, an imbalance between
291 chlorophyll a and chlorophyll b, and an appearance of carotenoids. For the healthy multiset
292 (Figure 2b), the green present values in the visible region that oscillate between 0.6 and 0.4 and
293 the slope in the near infrared region increases from 750 to 800 nm. Therefore, this component
294 seems to reflect an intensity level in the pigment region.

295 The distribution maps use a graduated colour scale per column, where the blue colour
296 corresponds to small concentration values and the red colour to large values. Differences
297 between the scores of white and red wine varieties can be observed for the infected multiset
298 (Figure 2a). Unlikely, there is no visible differences between the white and red grape varieties
299 on the healthy multiset (Figure 2b). This seems to show that the spectra obtained vary according
300 to the grape variety at the onset of the disease. For example, the white varieties (Chardonnay,
301 Colombars and Loin de l'oeil) have abnormally high values for the red component (figure 2a)
302 might due to the fact that these varieties have low anthocyanin levels but still retain the
303 Chlorophyll pigments. Therefore, it could explain why these leaves retain their green colouring
304 in contrast to the red grape varieties (see RGB images in Figure 2). Indeed, very low scores for
305 the third component will translate into a redder and greener colouration of the leaves.

306

307 **FIGURE 2**

308

309 Once, the basic spectral signatures that differentiate between infected and healthy images are
310 resolved, MCR-ALS analysis of the multiset formed by both infected and healthy training
311 images (**DaugtrFH**) was performed. In this case, the correspondence among species constraint
312 was also used since the presence/absence of constituents in each sample was known. From this
313 information, a matrix containing 72 blocks, (representing the 47 infected and the 25 healthy
314 training images analysed simultaneously) and 4 columns (representing the number of
315 constituents: both common blue and red contributions and the specific contributions for each
316 multiset) coding the presence (1) or absence (0) of each constituent in each image was
317 introduced as information in the resolution process. The absent constituents in the image were
318 then forced to have null concentration profiles.

	1	2	3	4
47 blocks Flavescence Dorée	1	1	1	0
25 blocks Healthy	1	0	1	1

319

320 Figure 3 shows the MCR-ALS resolved distribution maps (corresponding to the same images
321 in Figure 2) and pure spectra of the **DaugtrFH** multiset. The resolved spectra in Figure 3 are
322 rather similar to the pure spectra obtained from the MCR-ALS analyses of both infected and
323 healthy multisets (see Figure 2). Blue distribution maps refer to absent constituents in images
324 and the rest of the maps are consistent with those obtained in Figure 2, matching the relative
325 concentration of the different constituents in the images. A rather similar fit to previous MCR-
326 ALS analysis (see Table 2) was obtained ($r^2\% = 99.96$), strongly supporting the MCR results.
327 The introduction of the correspondence among species constraint does not perturb the natural

328 behaviour of the dataset. On the contrary, it improves the accuracy of the resolved profiles and
329 reduces ambiguity.

330

331 **FIGURE 3**

332

333 In order to validate this model, distribution maps related to the multiset structure containing all
334 the test dataset (**DtestFH**) were calculated by a single non-negative least-squares step using the
335 pure spectra obtained in Figure 3 (**S^TtrFH**). Satisfactory results (calculated $r^2 = 99.95$) with
336 consistent distributions maps (data not shown) are obtained, validating the MCR results.

337 Now, the basic spectral signatures and distribution maps of pure compounds contributions of
338 both infected and healthy images can be estimated. However, some components related to the
339 two different class types (infected and healthy pixels) may overlap. Thus, this unsupervised
340 method is not sufficient to distinguish between these two class-type, and hence an appropriate
341 method for discrimination is required. Therefore, the MCR scores of both training (**CaugtrFH**)
342 and test (**CaugtestFH**) sets were used for discriminant analysis. In practice, the FDA enables
343 to determine the relation between these scores and the most probable class of the samples.

344 **4.2. FDA model. Class assignment at pixel and leaf scale.**

345 The FDA model is calibrated based on the MCR scores of training dataset (**CaugtrFH**). Once
346 the MCR/FDA model is estimated, it was used to predict the class (infected or healthy) of each
347 pixel in the test dataset (**CaugtestFH**). At pixel level, the discrimination rate (CR) of test set is
348 equal to 85.1%. Both infected and healthy test pixels were correctly classified into their
349 corresponding class with more than of 75% and 95% in accuracy, respectively. For infected
350 pixels, a lower CR value is obtained, consequently, the model present a lesser ability to
351 determine infected than healthy pixels. This can be attributed to the labelling process. Indeed,

352 leaves affected by Flavescence dorée are entirely labelled as infected, *i.e.* every pixel of the leaf
353 is labelled the same. Indeed, at early stages, infected leaf images most likely include healthy
354 pixels or pixels presenting slight symptoms that were labelled in a single infected class and then
355 used in the calibration. Therefore, the model for FD could be depreciated by the presence of
356 healthy samples, hence the lesser accuracy in discrimination.

357 In order to better evaluate the capacity of the MCR/FDA strategy to discriminate between
358 infected and healthy leaves, pixel-wise decisions are summarised at the scale of the leaf. Indeed
359 the chemical information is relevant/ consistent at the scale of the spectrum/pixel. However, on
360 a pythopathological view, it is more sensible to consider at the scale of an organ (*i.e.* the leaf in
361 this case). Considering the characteristics of the symptoms and the development of the disease,
362 it is proposed to consider that a leaf is infected if more than 50% of its pixels are classified as
363 such. Therefore the CRs for images are calculated as the percentage of correct predictions to
364 the total number of pixels for each image (see Table 3). On the other hand, a healthy leaf should
365 exhibit in total very few abnormal pixels. Therefore, it is considered that a leaf is healthy if 75%
366 of its pixels are healthy (to take into account the fact that some part of the leaf could be less
367 vigorous but still unaltered by any disease).

368

369 **TABLE 3**

370

371 From Table 3, a satisfactory CR higher than 74 % for all healthy leaves can be observed except
372 for the image *DtestgH2* (63.3 %). Similar results are obtained for infected leaves. Only the
373 image *DtestcoF2* presents a CR lower than 50 %. However, *DtestgF7*, *DtestfF1*, *DtestcF1*,
374 *DtestIF2* and *DtestF3* show also CR lower than 60 %. This lesser accuracy for infected images
375 could be explained by the greater variability induced by the diversity of severity stages of the
376 pathology. In addition, the model was calibrated using a single multiset, that included diverse

377 red and white grape varieties that exhibit different visible symptoms. Nevertheless, since
378 discrimination results obtained for almost all images are satisfactory, the MCR/FDA strategy
379 could be considered adequate and future leaves are expected to be properly classified into their
380 corresponding class.

381 For a better evaluation of the lowest CR results of both **DtestcoF2** and **DtestgH2** images, Figure
382 4 a and b shows their predicted distribution maps (**CtestcoF2** and **CtestgH2**, respectively)
383 alongside their corresponding RGB images. It can be seen that **DtestcoF2** which presents only
384 25.2 % of infected pixels exhibits early and slight symptoms (as shown by a general low-level
385 green hue), and hence its uncertain state could explain its low accuracy. Likewise, **DtestgH2**
386 image does not have the appearance of a healthy leaf due to the presence of some stains
387 (possibly confounding factors such as stresses or fungal diseases). This suggests that possibly
388 a binary discrimination assignment without an external class for confounding factors is
389 insufficient for this application, which could also explain the lower CR for this particular leaf.
390 Moreover, as way of example Figure 4 c and d shows two examples of good predictions for
391 both **DtestgF1** and **DtestgH1** images.

392

393

394 **FIGURE 4**

395

396 In summary, it can be said that the combination of HSI and the method MCR-ALS with FDA
397 model proved to be efficient to distinguish between infected and healthy images. However, to
398 evaluate the discriminant potential of the proposed approach, larger data sets showing a greater
399 variability of symptoms and infection stages is required. Moreover, confounding factors such
400 as abiotic stresses or other phytopathology exhibiting similar symptoms should be also tested.

401 Ultimately, the processing of images representing canopies rather than isolated single leaves
402 should be taken into consideration to guarantee its feasibility in field conditions.

403

404 **5. Conclusions**

405 The strategy of combining MCR-ALS and FDA proved its interest for the discrimination
406 between healthy and infected leaves by Flavescence dorée based on the use of hyperspectral
407 images. For the first time, this strategy was applied as a phytopathology detection approach.

408 MCR-ALS enables to extract some relevant signatures that can discriminate healthy leaves from
409 leaves infected by Flavescence dorée. The pure component resulting from this model can be
410 interpreted concerning the visible symptoms of FD and to some associated physiochemical
411 disruptions. The relative abundances of these components within the leaves (MCR scores) can
412 be processed with FDA and provide an efficient discrimination of the leaves.

413 To improve the proposed strategy and reach a practical application in viticulture, some aspects
414 such as confounding factors, progressive infection stages and feasibility in the field should be
415 taken into account. Another development to improve these results, would be to upgrade the
416 labelling process, e.g. by selecting areas of the leaves clearly identified as infected rather than
417 assigning a class to the whole leaf. Nonetheless, Hyperspectral imaging combined with the
418 proposed data processing approach has the potential to be a valuable strategy to detect grapevine
419 diseases.

420 **6. Acknowledgements**

421 This work was supported by the INTERREG SUDOE SOE3/P2/E0911 Viniot projet

422 **7. Conflict of interest**

423 The authors report there are not conflicts of interest

424 **8. Bibliography**

- 425 1 E. Tona, A. Calcante and R. Oberti, *Precis. Agric.*, 2018, **19**, 606–629.
- 426 2 G. A. Carter and A. L. A. N. K. K. Napp, 2001, **88**, 677–684.
- 427 3 A. Mahlein, U. Steiner, C. Hillnhütter, H. Dehne and E. Oerke, 2012, 1–13.
- 428 4 S. Sankaran, A. Mishra, R. Ehsani and C. Davis, *Comput. Electron. Agric.*, 2010, **72**,
- 429 1–13.
- 430 5 F. Dor, G. Disease and H. Al-saddik, , DOI:10.3390/s17122772.
- 431 6 J. Albetis, A. Jacquin, M. Goulard, H. Poilvé, J. Rousseau, H. Clenet, G. Dedieu and S.
- 432 Duthoit, , DOI:10.3390/rs11010023.
- 433 7 C. F. Jordan, *Ecology*, 1969, **50**, 663–666.
- 434 8 R. Tauler, *Chemom. Intell. Lab. Syst.*, 1995, **30**, 133–146.
- 435 9 R. de Juan, Anna; Rutan, S. and Tauler, in *Comprehensive Chemometrics*, ed. B.
- 436 Brown, S. D. ; Tauler, R. and Walczak, Elsevier, 2010, vol. 2, pp. 325–344.
- 437 10 D. . Hand, *Biometrical J.*, 1985, **27**, 148.
- 438 11 V. Olmos, M. Marro, P. Loza-Alvarez, D. Raldúa, E. Prats, F. Padrós, B. Piña, R.
- 439 Tauler and A. de Juan, *J. Biophotonics*, 2018, **11**, e201700089.
- 440 12 S. Mas, A. Torro, L. Fernández, N. Bec, C. Gongora, C. Larroque, P. Martineau, A. De
- 441 Juan and S. Marco, *Talanta*, 2020, 208, 120455
- 442 13 S. Hennessy, G. Downey and C. P. O'Donnell, *J. Agric. Food Chem.*, 2009, **57**, 1735–
- 443 1741.
- 444 14 R. Karoui, M. Hammami, H. Rouissi and C. Blecker, *Food Chem.*, 2011, **127**, 743–
- 445 748.
- 446 15 F. A. Kruse, A. B. Lefkoff, J. W. Boardman, K. B. Heidebrecht, A. T. Shapiro, P. J.
- 447 Barloon and A. F. H. Goetz, *Remote Sens. Environ.*, , DOI:10.1016/0034-
- 448 4257(93)90013-N.

449 16 A. de Juan, in *Data Handling in Science and Technology*, 2020.

450 17 N. Lord, G. H. Golub and C. F. Van Loan, *Math. Gaz.*, , DOI:10.2307/3621013.

451 18 W. Windig, C. E. Heckler, F. A. Agblevor and R. J. Evans, *Chemom. Intell. Lab. Syst.*,
452 , DOI:10.1016/0169-7439(92)80104-C.

453 19 R. Tauler, A. Smilde and B. Kowalski, *J. Chemom.*, 1995, **9**, 31–58.

454 20 P. . Mahalanobis, in *On the generalized distance in statistics*, Proceedings of the
455 National Institute of Sciences (Calcutta) (1936), 1936.

456 21 J. Jaumot, A. de Juan and R. Tauler, *Chemom. Intell. Lab. Syst.*, 2015, **140**, 1–12.

457

458

459

460

461

462

463

464

465

466

467

468

469

470

471

472

473

474 **9. Figure captions**

475 **Figure 1.** MCR application to a) an individual hyperspectral image, b) an image multiset
 476 structure.

477 **Figure 2.** MCR-ALS results for a) the multiset of training infected dataset (**DaugtrF**) and b)
 478 the multiset of training healthy dataset (**DaugtrH**). Left plots: related MCR-ALS distribution
 479 maps with their corresponding RGB images. Right plots: resolved pure MS spectra. Varieties
 480 in italics correspond to *white varieties*.

481 **Figure 3.** MCR-ALS results for the multiset of both infected and healthy training (**DaugtrFH**).
 482 Left plots: related MCR-ALS distribution maps with their corresponding RGB images. Right
 483 plots: resolved pure MS spectra. Varieties in italics correspond to *white varieties*.

484 **Figure 4** Predicted distribution maps of: a) **DtestcoF2**, b) **DtestgH2**, c) **DtestgF1** and d)
 485 **DtestgH1** images with their corresponding RGB images.

486

487

488

489

490

491 **Table 1.** Total number of leaves images selected from the different varieties and the number
 492 of images both in the training and in the independent test set for the MCR-ALS/FDA models
 493 (see section 3.2 for more information).

494

Varieties	Flavesc. Dorée			Healthy		
	Total	Training	Test	Total	Training	Test
Gamay (g)	19	12	7	10	7	3
Fer (f)	10	7	3	5	3	2
Duras (d)	9	6	3	3	2	1

<i>Chardonnay (c)</i>	12	8	4	11	7	4
<i>Colombard (co)</i>	10	7	3	5	3	2
<i>Loin de l'œil (l)</i>	10	7	3	5	3	2

Varieties in italics correspond to *white varieties*

495

496

497

498

499

500

501

502

503

504

505

506

507

508

509

510

511

512 **Table 2.** Number of resolved components and variance explained by MCR-ALS analysis of
513 **DaugtrF** and **DaugtrH** multiset structures.

514

Multiset	Explained variance	Resolved components
DaugtrF	99.91 %	3
DaugtrH	99.97 %	3

515

516

517
518
519
520
521
522
523
524
525
526
527
528
529
530
531

532 **Table 3.** Discrimination Rate (CR) of infected and healthy images of the test dataset.

Infected Dabcd*	CR	Healthy Dbacd*	CR
DtestgF1	95.1	DtestgH1	97.7
DtestgF2	94.2	DtestgH2	63.3
DtestgF3	93.2	DtestgH3	74.7
DtestgF4	90.8	DtestfH1	94.6
DtestgF5	73.5	DtestfH2	99.9
DtestgF6	78.6	DtestdH1	99.3
DtestgF7	53.8	DtestcH1	99.9

DtestfF1	53.1	<i>DtestcH2</i>	99.9
DtestfF2	88.4	<i>DtestcH3</i>	99.4
DtestfF3	83.9	<i>DtestcH4</i>	99.9
DtestdF1	93.9	<i>DtestcoH1</i>	97.2
DtestdF2	90.4	<i>DtestcoH2</i>	98.9
DtestdF3	94.0	<i>DtestlH1</i>	100
DtestcF1	53.5	<i>DtestlH2</i>	99.8
<i>DtestcF2</i>	73.8		
<i>DtestcF3</i>	86.4		
<i>DtestcF4</i>	91.3		
DtestcoF1	51.7		
DtestcoF2	25.2		
<i>DtestcoF3</i>	94.2		
<i>DtestlF1</i>	81.7		
DtestlF2	55.6		
DtestlF3	50.1		

*Image code Dabcd; a= training (tr) or test ; b=variety; c=flavescence doré (F) or healthy (H) and d=sample number
 Varieties in italics correspond to *white variety*. Images in red present low CR results.

533

534

535

536

537