



**HAL**  
open science

## Massive spectral data analysis for plant breeding using parSketch-PLSDA method: Discrimination of sunflower genotypes

Maxime Ryckewaert, Maxime Metz, Daphné Héran, Pierre George, Bruno Grèzes-Besset, Reza Akbarinia, Jean Michel Roger, Ryad Bendoula

### ► To cite this version:

Maxime Ryckewaert, Maxime Metz, Daphné Héran, Pierre George, Bruno Grèzes-Besset, et al.. Massive spectral data analysis for plant breeding using parSketch-PLSDA method: Discrimination of sunflower genotypes. Biosystems Engineering, 2021, 210, pp.69-77. 10.1016/j.biosystemseng.2021.08.005 . hal-03329674

**HAL Id: hal-03329674**

**<https://hal.inrae.fr/hal-03329674>**

Submitted on 7 Jun 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.



Distributed under a Creative Commons Attribution - NonCommercial - NoDerivatives 4.0 International License

1 Massive spectral data analysis for plant breeding using  
2 parSketch-PLSDA method: discrimination of sunflower  
3 genotypes

4 Maxime Ryckewaert<sup>a,b</sup>, Maxime Metz<sup>a,b</sup>, Daphné Héran<sup>a</sup>, Pierre George<sup>c</sup>,  
5 Bruno Grèzes-Besset<sup>c</sup>, Reza Akbarinia<sup>d</sup>, Jean-Michel Roger<sup>a,b</sup>, Ryad  
6 Bendoula<sup>a</sup>

7 <sup>a</sup>*ITAP, Univ Montpellier, INRAE, Institut Agro, Montpellier, France*

8 <sup>b</sup>*ChemHouse Research Group, Montpellier, France*

9 <sup>c</sup>*Innolea, 6 chemin des Panedautes, 31700 Mondonville, France*

10 <sup>d</sup>*Inria & LIRMM, Univ Montpellier, France*

---

11 **Abstract**

12 In precision agriculture and plant breeding, the amount of data tends to  
13 increase. This massive data is becoming more and more complex, leading  
14 to difficulties in managing and **analysing** it. Optical instruments such as  
15 NIR Spectroscopy or hyperspectral imaging are gradually expanding directly  
16 in the field, increasing the amount of spectral database. **Using these tools**  
17 **allows access to non-destructive and rapid measurements to classify new va-**  
18 **rieties according to breeding objectives.** Processing this massive amount of  
19 spectral data is challenging. In a context of genotype discrimination, we pro-  
20 pose to apply a method called parSketch-PLSDA to **analyse** such a massive  
21 amount of spectral data. **ParSketch-PLSDA is a combination of an index-**  
22 **ing strategy (parSketch) and the reference method (PLSDA) for predicting**  
23 **classes from multivariate data.** For this purpose, a spectral database was  
24 formed by collecting 1,300,000 spectra **generated from hyperspectral images**  
25 of leaves of four different sunflower genotypes. ParSketch-PLSDA is com-

26 pared to a PLSDA. Both methods use the same set of calibration and test.  
27 The prediction model obtained by PLSDA has a classification error close  
28 to 23% on average across all genotypes. ParSketch-PLSDA method outper-  
29 forms PLSDA by greatly improving prediction qualities by 10%. **Indeed, the**  
30 **model built with ParSketch-PLSDA has the ability to take into account non-**  
31 **linearities among data sets.** These results are encouraging and allow us to  
32 anticipate the future bottleneck related to the generation of a large amount  
33 of data from phenotyping.

34 *Keywords:* Spectroscopy, Massive data, Digital Agriculture, Precision  
35 Agriculture, Chemometrics

---

## 36 **1. Introduction**

37 In recent years, precision agriculture and plant breeding have tended to  
38 increase the quantity and complexity of phenotyping related data (Mahlein,  
39 2016; Tripodi et al., 2018; Awada et al., 2018). Managing and **analysing**  
40 huge amounts of data are identified as a future bottleneck in phenotyping  
41 (Tripodi et al., 2018). Indeed, over the last few years, high throughput  
42 phenotyping (HTP) platforms in the laboratory or directly in the field have  
43 been flourishing (Chawade et al., 2019; Shakoor et al., 2017). These platforms  
44 provide a monitoring of one or more phenotypic traits of the vegetation. This  
45 information can be obtained at different spatial, spectral or temporal scales  
46 depending on the studied level, which could be vegetation organ, individual  
47 or even population (Dhondt et al., 2013; Mahlein, 2016; Mutka et al., 2016).  
48 The higher the spatial, spectral or temporal resolution, the larger the amount  
49 of data.

50 Spectroscopy in the visible and near-infrared range (VIS-NIR) has proven  
51 to be relevant for providing useful information for vegetation monitoring.  
52 Several plant phenotyping issues can be tackled with high spectral resolution  
53 measurements such as biochemical variable access (Vigneau et al., 2011; Jay  
54 et al., 2017), disease (Lu et al., 2018; Yu et al., 2018) or stress detection  
55 (Behmann et al., 2014; Christensen et al., 2005).

56 From a technological point of view, spectral acquisitions directly in the  
57 field have been made possible thanks to spectrometer **miniaturisation** (Yan  
58 and Siesler, 2018; Beć et al., 2020) or hyperspectral imager evolution (Mishra  
59 et al., 2020; Fiorani and Schurr, 2013). Associated with mobile vectors (such  
60 as UAV, tractor, pedestrian), these tools become **HTP** instruments and gen-  
61 erate a large amount of spectral data. However, simple computations on this  
62 amount of data such as outlier detection or the use of pre-processing become  
63 difficult to perform and very time consuming (Szymańska, 2018). Processing  
64 this massive amount of spectral data is challenging.

65 In chemometrics, most popular methods as Partial Least Square (PLS)  
66 (Wold et al., 2001) are based on an assumption of linear relationship be-  
67 tween spectral data and specific variables (Mark and Workman, 2007). These  
68 methods are popular because of their good predictive performances and low  
69 computation time. Conversely, using these methods may not provide good  
70 prediction models when relationships between spectra and variable of interest  
71 are non-linear.

72 When dealing with a large amount of spectral data, complex structures  
73 and non-linear relationships can arise which may compromise linear regres-  
74 sion approaches (Dardenne et al., 2000). In practice, using a linear classifi-

75 cation or regression to predict a complex database would lead to degraded  
76 results (Bertran et al., 1999; Ni et al., 2014). As a consequence, linear meth-  
77 ods are challenged on large amounts of data. Furthermore, some methods,  
78 called local methods, exploit data based on a restricted **neighbourhood** of  
79 individuals which greatly improves prediction quality (Dardenne et al., 2000;  
80 Pérez-Marín et al., 2007; Davrieux et al., 2016; Naes et al., 1990). These  
81 methods can be used to overcome non-linearity problems under the assump-  
82 tion that with a restricted **neighbourhood**, the relationship between spectra  
83 and variables becomes linear. The **parSketch-PLSDA** method has recently  
84 been proposed to implement a local approach to a large volume of data  
85 (Metz et al., 2020). Therefore, **parSketch-PLSDA** can be used to address the  
86 complex analysis of large amount of spectral data from phenotyping.

87 In this paper, we propose to study the use of the **parSketch-PLSDA**  
88 method to exploit a large amount of spectral data, **generated from hyper-**  
89 **spectral images of leaves of four different sunflower genotypes**. Additionally,  
90 we compare this method with a reference method in an application of dis-  
91 crimination of different sunflower varieties.

## 92 **2. Materials and methods**

### 93 *2.1. Biological material*

94 Four sunflower genotypes (called A, B, C and D) were grown in a green-  
95 house at INRAE France in **well-watered** conditions **by using a flood and**  
96 **drain system refilling water every 48 hours**. **All pots used the same pot-**  
97 **ting soil (Pot Clay coarse, Floradur, Floragard)**. Water and lighting con-  
98 ditions were similar for each pot with a day-night cycle of 16h/8h. **The**

99 temperature was 25°C with a relative humidity in the greenhouse ranging  
100 from 50% to 60%. The greenhouse was equipped with multispectral light-  
101 ing (450 nm, 560 nm, 660 nm, 730 nm and 6000°K) controlled by Herbro  
102 automaton (GreenHouseKeeper) with Photosynthetically Active Radiation  
103 (PAR) set at 400  $\mu\text{mol}/\text{m}^2/\text{s}$ .

104 For the four selected genotypes, two potted plants (called P1 and P2)  
105 of each were grown. Four leaves were collected at the upper and middle  
106 parts of each plant, except for the genotype D where only two leaves of each  
107 plant were collected. Leaf petioles were immediately wrapped with water  
108 soaked paper before measurements. In total, 28 leaves were then collected  
109 and measured.

## 110 2.2. Spectral acquisitions

111 Spectral data of the prepared leaf samples were acquired in the reflectance  
112 mode by using a laboratory-based line scanning Hyperspectral Imaging Sys-  
113 tem (HIS). The HIS system was composed of a linear halogen light (Haloline,  
114 Osram, 150 W), a translation rail (Linear Unit LES 4, Iselautomation, Ger-  
115 many), and a detection system. The sample was placed on a translation  
116 rail, synchronised with the acquisition software (NEO Hypspec, Norsk Elek-  
117 tro Optikk AS) which can record images when sample was scanned under the  
118 hyperspectral camera (NEO Hypspec VNIR-1600 with 30cm-objective, Norsk  
119 Elektro Optikk AS, Skedsmokorest, Norway). Spectral data were acquired  
120 in the 400 – 1000 nm wavelength range with 3.7 nm intervals.

121 For each sample, the reflected light intensity ( $I_s(\lambda)$ ) was measured at  
122 each wavelength . Dark current image ( $I_b(\lambda)$ ) was also recorded for each  
123 measure. A white reference (SRS99, Spectralon  $\text{\textcircled{R}}$ ) was used as a reference

124 ( $I_o(\lambda)$ ) to standardize images from non-uniformities of all components of the  
125 instrumentation (light source, lens, detector). From these measurements,  
126 reflectance ( $R_s(\lambda)$ ) was calculated for each sample:

$$R_s(\lambda) = \frac{I_s(\lambda) - I_b(\lambda)}{I_0(\lambda) - I_b(\lambda)} \quad (1)$$

127 For all hyperspectral images, vegetation pixels were selected to form a  
128 spectral data set. This selection was made by a threshold procedure (Fig. 1).  
129 Indeed, vegetation and background pixels were easily identified by comparing  
130 their reflectance value at 800 nm to a threshold defined here at 30%. Leaf  
131 spectra were collected from the 28 hyperspectral images, representing more  
132 than 1,300,000 spectra.



Figure 1: Pixel selection from a hyperspectral image of a sunflower leaf (a) image, (b) mask based on threshold values

### 133 2.3. Data analysis

134 Two methods were used to compare their ability to discriminate sunflower  
135 genotypes. Both methods were applied to a similar data set, called test set,  
136 built out of the spectra database. Calculations were performed with the  
137 R software (version 3.6.1 (Core Team, 2013)) and rnirs package (<https://>

138 [github.com/mlesnoff/rnirs](https://github.com/mlesnoff/rnirs)) was used for classical discrimination methods  
139 (PLSDA).

### 140 2.3.1. *PLSDA method*

141 The Partial Least Squares for Discrimination Analysis (PLSDA) ([Barker](#)  
142 [and Rayens, 2003](#)) was used as reference method for classification. This  
143 method consisted of building models between multivariate data and a vector  
144 coding different classes (here, the four genotypes).

145 Multivariate data was represented by a matrix  $\mathbf{X}$  of size  $(n, p)$  where  $n$   
146 was the observation number and  $p$  the variable number. The  $n$  observations  
147 were identified by their corresponding class in the vector  $y$  of size  $(n,1)$  where  
148 values ranged from 1 to  $q$ , where  $q$  was the class number. The first step was  
149 to transform  $y$  into a dummy matrix  $\mathbf{Y}$  of size  $(n, q)$  also called disjunctive  
150 table.

$$y = \begin{bmatrix} 1 \\ 1 \\ 1 \\ 2 \\ 2 \\ 2 \\ 3 \\ 3 \\ 3 \end{bmatrix} \rightarrow \mathbf{Y} = \begin{bmatrix} 1 & 0 & 0 \\ 1 & 0 & 0 \\ 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 1 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \\ 0 & 0 & 1 \\ 0 & 0 & 1 \end{bmatrix} \quad (2)$$

151 An example of a dummy matrix is given in equation 2 with nine observa-  
152 tions belonging to three classes. The matrix  $\mathbf{Y}$  contains binary values (0,1)



153 where each column corresponded to a class. For a given observation, the  
154 class-corresponding column has a value of 1 while other columns were equal  
155 to 0.

156 Then, a Partial Least Square (PLS) model (Wold et al., 2001) was applied  
157 between  $\mathbf{X}$  and  $\mathbf{Y}$ .  $\mathbf{Y}$  being multidimensional, the algorithm PLS2 adapted  
158 to the prediction of several responses was used. Finally, a linear discriminant  
159 analysis (LDA) (Fisher, 1936) was applied between the PLS2 scores and  $\mathbf{Y}$ .

### 160 2.3.2. ParSketch-PLSDA method

161 The other strategy was to apply the parSketch-PLSDA method (Metz  
162 et al., 2020), an extension of the K-Nearest Neighbours (KNN)-PLSDA  
163 method for massive data processing (Lesnoff et al., 2020). ParSketch-PLSDA  
164 was used to combine an indexation strategy (parSketch) and the PLSDA.  
165 An approximation of the neighbourhood was defined for each spectrum to be  
166 classified. This neighbourhood was then used to compute a PLSDA model  
167 and to predict which class belong new spectra.

168 ParSketch was performed in three steps: dimension reduction, grid cre-  
169 ation, neighbourhood approximation. Three method parameters ( $v$ ,  $s$ ,  $m$ )  
170 were defined, corresponding to these three steps, and are described below.

171 First, a dimension reduction was achieved by calculating the matrix  $\mathbf{T}$   
172 corresponding to the sketch of the matrix  $\mathbf{X}$  as follows:

$$\mathbf{T} = \mathbf{XP} \tag{3}$$

173 Where  $\mathbf{P}$  was a matrix of size  $(p,v)$  containing values of -1 or 1 according  
174 to a random selection. The first parameter of ParSketch ( $v$ ), corresponding  
175 to the column number of  $\mathbf{P}$  was then defined. The higher the value of  $v$

176 the better the approximation of the **neighbourhood**. However, the larger the  
 177 value of  $v$  the longer the parSketch method computation time.

178 The second step corresponding to the grid creation process (see Fig. 2)  
 179 was to segment the space (2d) formed by adjacent pairs of  $\mathbf{T}$  columns. The  
 180 number of segments ( $s$ ) is the second parSketch parameter. The higher the  
 181 value of  $s$  the better the approximation of the nearest **neighbours**. However,  
 182 the greater the value of  $s$ , the smaller the number of **neighbours**.

183

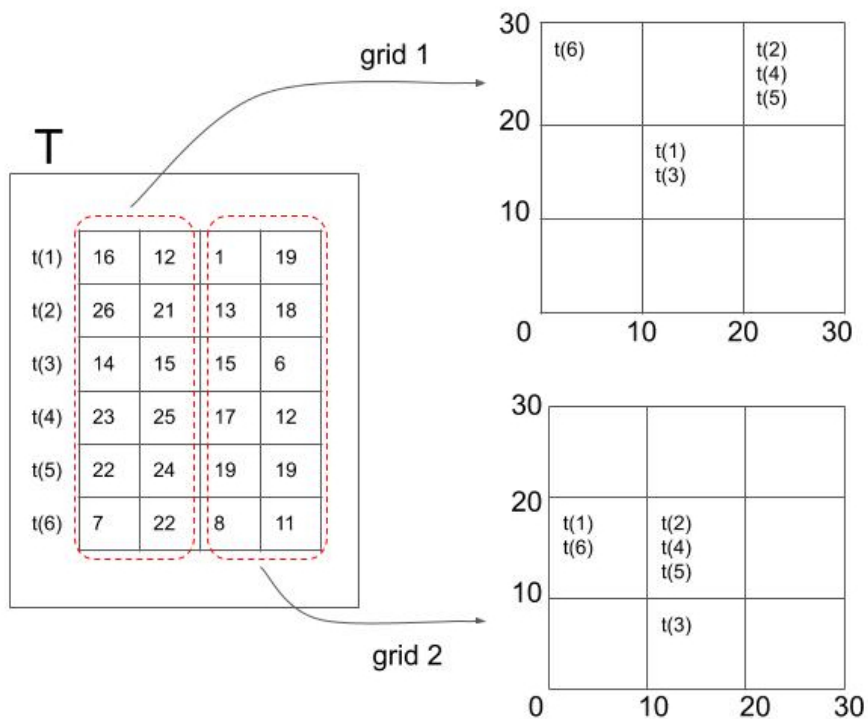


Figure 2: An illustrated example of grid creation with a segment number  $s = 3$

184 The last step to configure parSketch was to define the minimal number  
 185  $m$  of grids returned in the **neighbour's** search. This step corresponded to the

186 **neighbourhood** approximation for the grid search process (see Fig. 3). The  
 187 higher the value of  $m$  the better the approximation of the nearest **neighbours**.  
 188 Observations present in the same cell for at least  $m$  grid number are selected  
 189 as **neighbours** of the individual to be predicted. However, the greater the  
 190 value of  $m$ , the smaller the number of **neighbours** returned by parSketch  
 191 method.

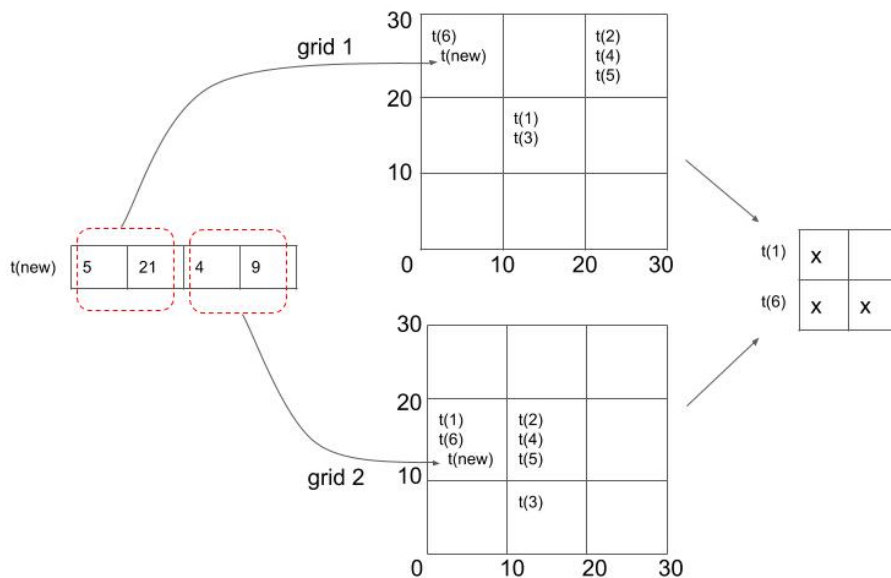


Figure 3: An illustrated example of grid search. For a new measure  $t(new)$  and a  $m$  value equals to 2,  $t(6)$  will be returned as a **neighbour** because it is present in two grids next to  $t(new)$ , whereas  $t(1)$  will not be considered as a **neighbour**

192 *2.4. Evaluation strategies and method parameterization*

193 The data set was divided into two independent data sets: a calibration  
 194 set and an independent test set . The calibration set was formed with the  
 195 14 images acquired on P1 plants and corresponding to about 650,000 spec-

196 tra. The PLSDA model and parSketch-PLSDA method were both calibrated  
197 using all spectra of this calibration set.

198 The test set was formed with the other 14 images acquired on P2 plants  
199 (independent from P1). 1000 spectra were randomly selected in each image  
200 totalling 14,000 spectra for the test set. Spectra that could not be predicted  
201 by parSketch due to lack of **neighbours** were removed from the test set. In  
202 the end, the same test set were used for parSketch-PLSDA and for PLSDA.

203 For both methods, validation steps were performed on the calibration set  
204 in order to minimize overfitting.

205 To build the PLSDA model, the cross-validation step consisted of splitting  
206 the calibration data set into different blocks in order to calculate calibration  
207 and validation errors. This approach, also called k-fold validation ([Wold,](#)  
208 [1978](#); [Camacho and Ferrer, 2012](#)) was carried out with five blocks repeated  
209 three times. Validation errors were then computed and led to the number of  
210 latent variables to be retained.

211 For parSketch-PLSDA, a parametrisation step was performed to config-  
212 ure the three parSketch parameters. This step was performed by analyzing  
213 distributions of returned **neighbours** according to two parSketch parameters:  
214 number of segments  $s$  and the common minimum grids  $m$ . Here, the number  
215 of random vectors  $v$  was set to a value of 20. Afterwards, a PLSDA model  
216 was established. The number of latent variables was optimized for a subset of  
217 the calibration set, called the validation set. This validation set was formed  
218 with four images of the calibration set by randomly selecting 1000 spectra in  
219 each image.

220 In order to compare both methods, confusion matrices were obtained and

221 percentages of precision and recall were calculated according to the following  
222 equations:

$$\text{Precision} = \frac{tp}{tp + fp} \quad (4)$$

$$\text{Recall} = \frac{tp}{tp + fn} \quad (5)$$

223 Where  $tp$ ,  $fp$  and  $fn$  corresponded to true positives, false positives and  
224 false negatives respectively. On the one hand, for a given class, precision  
225 value assessed the predictive quality of the model based on the proportion  
226 of well-classified observations among all observations that were classified in  
227 the same corresponding class. On the other hand, recall, also called sensi-  
228 tivity, evaluated the number of well-classified observations compared to the  
229 total number of observations of the given class. These two criteria are com-  
230plementary to evaluate the model performances. These two figures of merit  
231 were expressed as percentages. The higher the values, the better the model  
232 performance.

233 **3. Results and discussion**

234 *3.1. Data visualization*

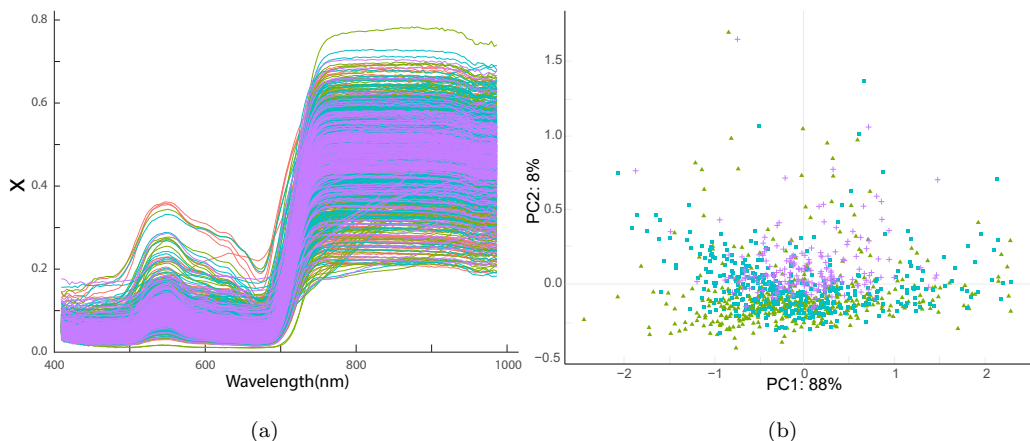


Figure 4: All data set including Genotype A (red), B (green), C (cyan) and D (violet): (a) spectra, (b) score plot of the first two principal components

235 Reflectance spectra shown in Fig. 4a correspond to 1000 spectra per  
236 class randomly selected among all the data set. These spectra correspond  
237 to vegetation spectra (Xu et al., 2019) : specific hollows at 450 nm and  
238 650 nm related to chlorophyll content, anthocyanin content at 550 nm; the  
239 red-edge towards 780 nm and a plateau in the near-infrared between 780 nm  
240 and 1000 nm. Besides, the main observed variability in the spectrum plot  
241 corresponds to an additive effect due to the scattering effect of the structure  
242 of the leaves. However, the number of spectra is too large to be able to  
243 describe difference between classes.

244 A principal component analysis was applied to these spectra. Figure 4b  
245 shows the score plot of the two first components. The first component rep-  
246 represents 88% of the spectra variability and 8% for the second component. On

247 these two components, scores are uniformly distributed without any evident  
248 distinction between genotypes. The exploratory study of the spectra shows  
249 that there are no outliers and that there is no distinct group on the first two  
250 components.

### 251 3.2. Model calibration

#### 252 3.2.1. PLSDA

253 Figure 5 shows the cross-validated error rate curve for PLSDA applied to  
254 all spectra of the calibration data set. The behaviour of the curve decreases  
255 continuously according to the latent variable (LV) number.

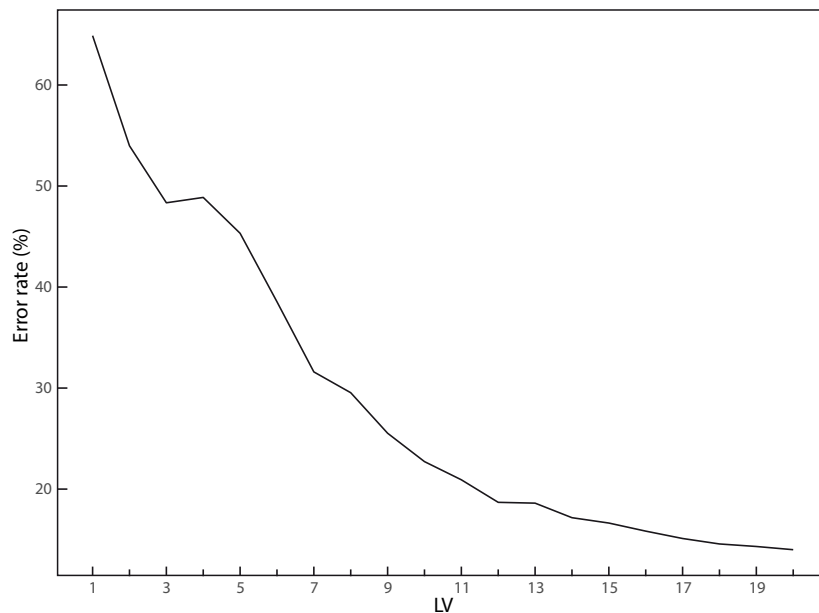


Figure 5: Evolution of the cross-validated error rate as a function of latent variables (LV) for the PLSDA applied to all spectra of the calibration data set

256 A high value of LV number generally shows the complex structure of a  
257 data set. This is expected with spectral measurements on vegetation (Metz

258 [et al., 2020](#)). With 16 LVs, an error rate with a value close to 12% is ob-  
 259 tained. However, after 16 LVs the predictive performance gain is very small.  
 260 Consequently, the PLSDA model is set to 16 LVs.

261 *3.2.2. ParSketch-PLSDA*

262 ParSketch parameters  $s$  and  $m$  are studied according to the statistical  
 263 distribution of the number of returned **neighbours** (Fig. 6).

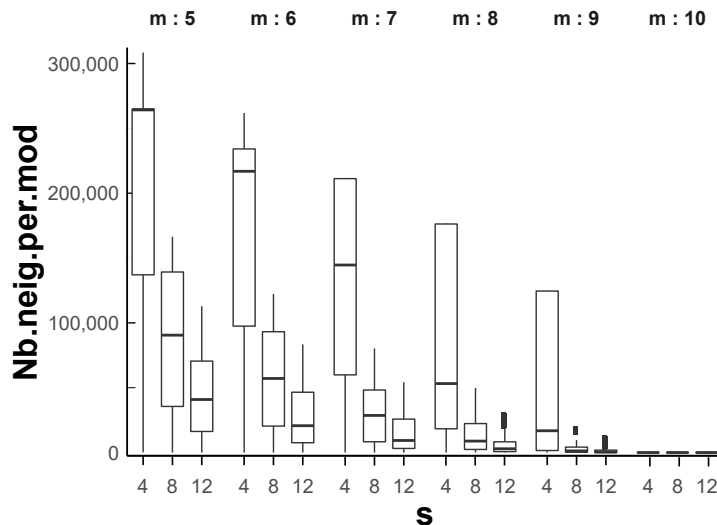


Figure 6: Distribution of the returned **neighbours** according to parSketch parameters  $s$  and  $m$  (number of segments and common minimum grids). Here,  $v$  (number of random vectors) parameter is fixed to 20

264 The number of **neighbours** decreases to a value close to zero when param-  
 265 eter values  $(s, m)$  increase. Indeed, on the one hand, an increase of number  
 266 of segments  $s$ , the number of returned **neighbours** will be lower. And on the  
 267 other hand, by increasing the minimum number of common grids  $m$ , the risk  
 268 of not having **neighbours** is high. This global trend is expected ([Metz et al.](#),



269 2020).

270 By contrast, when parameter values are low, the number of returned  
271 neighbours is high (close to 300 000 neighbours by individual to be pre-  
272 dicted). This situation is not desirable, as it may cause problems related  
273 to computation time constraints. As a result, parameters  $m$  and  $s$  must  
274 be chosen to have a sufficient number of neighbours, neither too much nor  
275 too little. As several values are possible, four combinations of the parSketch  
276 parameters are selected (Table 1) to compare their model performances.

Table 1: Combinations of the selected parSketch parameters and the corresponding median number of returned neighbours

Combination	$m$	$v$	$s$	Median neighbour number
(a)	9	20	8	1246
(b)	8	20	12	2903
(c)	7	20	12	9303
(d)	6	20	12	27030

277 Table 1 shows the retained values of the three parSketch parameters for  
278 these four combinations. The combination (a) was selected because the  
279 median number of neighbours is 1246. This low number of neighbours enables  
280 to quickly calibrate PLSDA models but it could be insufficient to have a  
281 good predictive quality. Indeed, a low median number of neighbours means  
282 that a large amount of observations do not have neighbour at all. For the  
283 combinations (b) and (c), higher numbers of neighbours are returned, with  
284 median values of 2903 and 9303 respectively. Finally, the highest median  
285 number of neighbours returned by parSketch is chosen with the combination

286 **(d)** with a value equal to 27030. In this case, constraints in computation  
287 time might appear. Moreover, the linear relationship between spectra and  
288 class variable of a small **neighbourhood** might be lost.

289 Validation error curves for the four retained combinations for parSketch-  
290 PLSDA are shown in the figure 7.

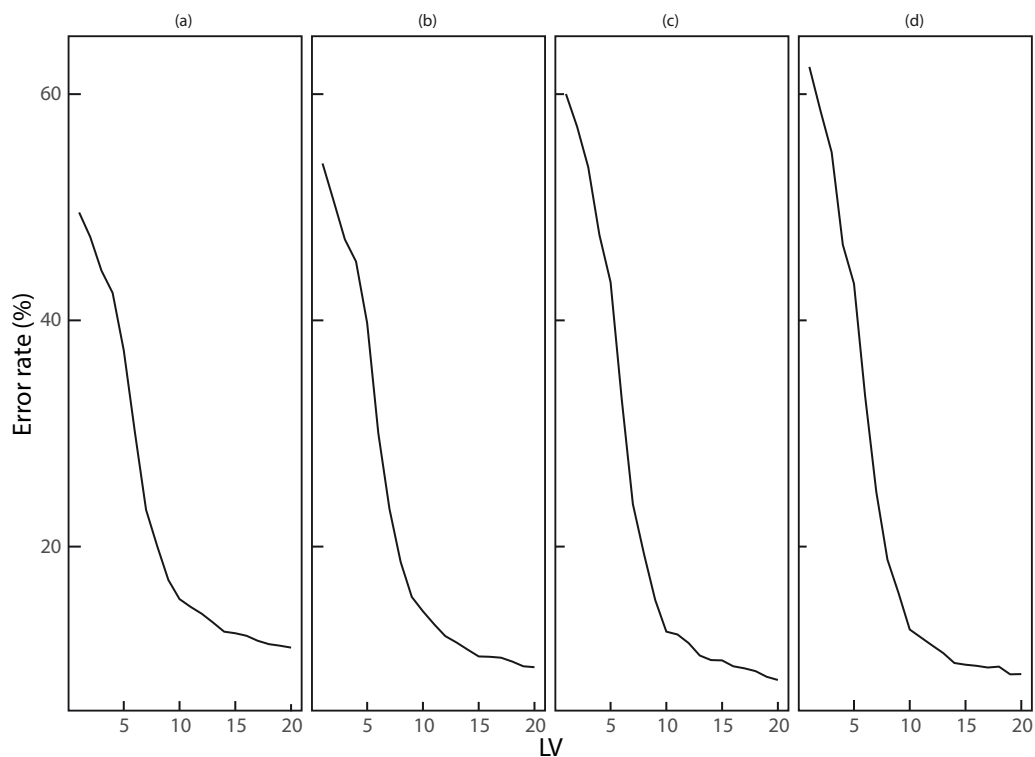


Figure 7: Evolution of the validation error rate as a function of latent variables for the four parameter combinations of parSketch-PLSDA

291 With higher error values, the combination **(a)** is less predictive than  
292 other parameter combinations. As expected, this combination having the  
293 smallest number of **neighbours**, the resultant model has poorer predictive  
294 performances than the three other ones.

295 The error curve obtained with the combination **(b)** reaches lower rates  
296 than the combination **(a)** curve. This means that the predictive capabilities  
297 of the model can be improved by slightly increasing the number of **neighbours**.  
298 The combination **(c)** has best predictive performance for the validation set  
299 with lowest values of classification error. The combination **(d)** has lower  
300 prediction quality than the combination **(c)** for a larger median number of  
301 **neighbours** per sample to be predicted (cf. Table 1).

302 The model with the lowest predictive quality has a high number of **neigh-**  
303 **bours**. This degradation reflects a non-linear aspect of the data set. By  
304 further increasing the number of returned **neighbours**, prediction qualities of  
305 parSketch will be close to the PLSDA method on the whole data set.

306 Finally, for this validation set the optimal parameter combination is the  
307 combination **(c)**. In this case, the number of latent variables is not easy to  
308 define. The number of latent variables is defined by a trade-off between the  
309 size of the model and the benefit of adding an extra dimension to the model.  
310 The number of latent variables chosen is therefore 16.

### 311 *3.3. Model testing*

312 The PLSDA model has been calibrated with all the spectra of the cali-  
313 bration set. Then this model is applied to the test set defined previously and  
314 its prediction performances are assessed in Table 2.

Table 2: Confusion matrix for PLSDA (16 LVs)

	<b>A</b>	<b>B</b>	<b>C</b>	<b>D</b>	<b>Recall (%)</b>
<b>A</b>	3120	121	339	250	<b>81</b>
<b>B</b>	238	2812	619	154	<b>74</b>
<b>C</b>	127	241	3463	75	<b>89</b>
<b>D</b>	283	173	270	1213	<b>63</b>
<b>Precision(%)</b>	<b>83</b>	<b>84</b>	<b>74</b>	<b>72</b>	

315 Precision and recall values are high for all classes A, B, C and D with  
316 values ranging from 72% to 84% for precision and from 63% to 89% for  
317 recall. Genotypes A and B have the highest precision values with values of  
318 83% and 84%, respectively. This means that 83% of spectra classified in  
319 genotype A, actually belong to genotype A. Few other genotypes are found  
320 in this class. The same argumentation holds true for 84% of genotype B  
321 spectra. For recall, genotypes A and C have the best values with 81% and  
322 89%, respectively. For these genotypes, spectra are mainly well-classified  
323 that is infrequently assigned to other classes.

324 The percentage missings from recall values correspond to the prediction  
325 error for each class. The prediction error of the whole data set, corresponding  
326 to the average error, is close to 23%. It is expected to have value for the test  
327 error slightly higher than the 12% observed during calibration (see Fig. 5).  
328 This means that the calibration set samples are representative of the test set  
329 despite their independence (as mentioned above, the test set corresponds to  
330 other plants of the same genotype).

Table 3: Confusion matrix for parSketch-PLSDA with combination(c) ( $m = 7$ ,  $v = 20$ ,  $s = 12$ )

	<b>A</b>	<b>B</b>	<b>C</b>	<b>D</b>	<b>Recall (%)</b>
<b>A</b>	3547	114	115	54	<b>93</b>
<b>B</b>	40	3256	473	54	<b>85</b>
<b>C</b>	58	211	3590	47	<b>92</b>
<b>D</b>	51	112	260	1516	<b>78</b>
<b>Precision(%)</b>	<b>96</b>	<b>88</b>	<b>81</b>	<b>91</b>	

331 Table 3 shows the parSketch-PLSDA prediction performance by giving  
 332 percentages of precision and recall for each genotype. Genotypes A and D  
 333 have the highest precision values with values of 96% and 91% respectively.  
 334 Besides, genotypes A and C have the highest recall values with values of  
 335 93% and 92% respectively. Genotype D has low recall values with both  
 336 methods (63% for PLSDA and 78% for parSketch-PLSDA). This is probably  
 337 due to the under-representation in the data set which may degrade the model  
 338 calibration. Indeed, only two images were acquired per plant of genotype D  
 339 compared to four images for the other genotypes.

340 Finally, overall recall and precision values have increased by almost 10%  
 341 with parSketch-PLSDA (83% and 87% respectively) compared to PLSDA  
 342 (76% and 77% respectively). Consequently, the model prediction error de-  
 343 creases to a value of 13%. This implies that parSketch-PLSDA model per-  
 344 forms better than the reference discriminant strategy. This improvement  
 345 in the classification results demonstrates the advantage of using a limited  
 346 number of neighbours to create a model.

347 As the methods used are locally linear, this improvement confirms the hy-  
348 pothesis that with a limited number of **neighbours**, the problem becomes lin-  
349 ear. The prediction improvement obtained with parSketch-PLSDA method  
350 highlights the presence of non-linear relationships between spectra and a class  
351 variable in the whole data set. which can be encountered when building a  
352 large spectral database.

#### 353 4. Conclusion

354 In this study, we compared the two classification strategies on the same  
355 calibration and test data sets.

356 For both methods, classification results are encouraging and confirm the  
357 interest of VIS-NIR spectroscopy for variety discrimination. Results showed  
358 that parSketch-PLSDA method outperforms PLSDA by improving predic-  
359 tion qualities by 10%. The use of the parSketch-PLSDA procedure in the  
360 exploitation of massive spectral data is confirmed and shows the interest of  
361 using a close **neighbourhood** of the spectra to be predicted.

362 It would be interesting to test such methods on a larger number of geno-  
363 types. This increase in the spectral database can potentially lead to an  
364 increase in complexity hence reducing the data set quality. Therefore, it  
365 would be interesting, in perspective, to evaluate other methods dealing with  
366 non-linearity.

367 In the framework of plant breeding, hyperspectral imaging or field mi-  
368 crospectrometers as tools for high-throughput plant phenotyping could be  
369 considered in real time with this method. In an applicative aspect, parS-  
370 ketch procedure is **parallelisable**, which shows the possibility of fast real-time

371 prediction of a large amount of data. We used parSketch-PLSDA on spectral  
372 data for close-range plant phenotyping. Other applications to plant breed-  
373 ing (disease, biotic/abiotic stress) or other applications related to precision  
374 agriculture could be considered. More generally, this method can be applied  
375 to any other application in analytical chemistry or metabolomics.

### 376 **Acknowledgement**

377 This work was conducted within the OPTIPAG Project supported by the  
378 grant ANR-16-CE04-0010 from the French Agence Nationale de la Recherche.

### 379 **References**

380 Anne-Katrin Mahlein. Plant disease detection by imaging sensors—parallels  
381 and specific demands for precision agriculture and plant phenotyping.  
382 *Plant disease*, 100(2):241–251, 2016. Publisher: Am Phytopath Society.

383 Pasquale Tripodi, Daniele Massa, Accursio Venezia, and Teodoro Cardi.  
384 Sensing technologies for precision phenotyping in vegetable crops: cur-  
385 rent status and future challenges. *Agronomy*, 8(4):57, 2018. Publisher:  
386 Multidisciplinary Digital Publishing Institute.

387 Lana Awada, Peter W. B. Phillips, and Stuart J. Smyth. The adoption of  
388 automated phenotyping by plant breeders. *Euphytica*, 214(8):148, August  
389 2018. ISSN 0014-2336, 1573-5060. doi: 10.1007/s10681-018-2226-z. URL  
390 <http://link.springer.com/10.1007/s10681-018-2226-z>.

391 Aakash Chawade, Joost van Ham, Hanna Blomquist, Oscar Bagge, Erik  
392 Alexandersson, and Rodomiro Ortiz. High-throughput field-phenotyping

393 tools for plant breeding and precision agriculture. *Agronomy*, 9(5):258,  
394 2019. Publisher: Multidisciplinary Digital Publishing Institute.

395 Nadia Shakoor, Scott Lee, and Todd C. Mockler. High throughput pheno-  
396 typing to accelerate crop breeding and monitoring of diseases in the field.  
397 *Current opinion in plant biology*, 38:184–192, 2017. Publisher: Elsevier.

398 Stijn Dhondt, Nathalie Wuyts, and Dirk Inzé. Cell to whole-plant phe-  
399 notyping: the best is yet to come. *Trends in Plant Science*, 18(8):  
400 428–439, August 2013. ISSN 1360-1385. doi: 10.1016/j.tplants.2013.  
401 04.008. URL [http://www.sciencedirect.com/science/article/pii/  
402 S1360138513000812](http://www.sciencedirect.com/science/article/pii/S1360138513000812).

403 Andrew M. Mutka, Sarah J. Fentress, Joel W. Sher, Jeffrey C. Berry,  
404 Chelsea Pretz, Dmitri A. Nusinow, and Rebecca Bart. Quantitative,  
405 image-based phenotyping methods provide insight into spatial and tem-  
406 poral dimensions of plant disease. *Plant Physiology*, page pp.00984.2016,  
407 July 2016. ISSN 0032-0889, 1532-2548. doi: 10.1104/pp.16.00984. URL  
408 <http://www.plantphysiol.org/lookup/doi/10.1104/pp.16.00984>.

409 Nathalie Vigneau, Martin Ecartot, Gilles Rabatel, and Pierre Roumet. Po-  
410 tential of field hyperspectral imaging as a non destructive method to as-  
411 sess leaf nitrogen content in wheat. *Field Crops Research*, 122(1):25–31,  
412 April 2011. ISSN 03784290. doi: 10.1016/j.fcr.2011.02.003. URL [http://  
413 //linkinghub.elsevier.com/retrieve/pii/S0378429011000451](http://linkinghub.elsevier.com/retrieve/pii/S0378429011000451).

414 Sylvain Jay, Nathalie Gorretta, Julien Morel, Fabienne Maupas, Ryad Ben-  
415 doula, Gilles Rabatel, Dan Dutartre, Alexis Comar, and Frédéric Baret.



- 416 Estimating leaf chlorophyll content in sugar beet canopies using millimeter-  
417 to centimeter-scale reflectance imagery. *Remote Sensing of Environment*,  
418 198:173–186, 2017. Publisher: Elsevier.
- 419 Jinzhu Lu, Reza Ehsani, Yeyin Shi, Ana Isabel de Castro, and Shuang  
420 Wang. Detection of multi-tomato leaf diseases ( late blight , target  
421 and bacterial spots ) in different stages by using a spectral-based sen-  
422 sor. *Scientific Reports*, 8(1):2793, February 2018. ISSN 2045-2322. doi:  
423 10.1038/s41598-018-21191-6. URL [https://www.nature.com/articles/  
424 s41598-018-21191-6](https://www.nature.com/articles/s41598-018-21191-6). Number: 1 Publisher: Nature Publishing Group.
- 425 Kang Yu, Jonas Andereg, Alexey Mikaberidze, Petteri Karisto, Fabio  
426 Mascher, Bruce A. McDonald, Achim Walter, and Andreas Hund. Hyper-  
427 spectral Canopy Sensing of Wheat Septoria Tritici Blotch Disease. *Fron-*  
428 *tiers in Plant Science*, 9, 2018. ISSN 1664-462X. doi: 10.3389/fpls.2018.  
429 01195. URL [https://www.frontiersin.org/articles/10.3389/fpls.  
430 2018.01195/full](https://www.frontiersin.org/articles/10.3389/fpls.2018.01195/full). Publisher: Frontiers.
- 431 Jan Behmann, Jörg Steinrücken, and Lutz Plümer. Detection of early  
432 plant stress responses in hyperspectral images. *ISPRS Journal of Pho-*  
433 *togrammetry and Remote Sensing*, 93:98–111, 2014. URL [http://www.  
434 sciencedirect.com/science/article/pii/S092427161400094X](http://www.sciencedirect.com/science/article/pii/S092427161400094X).
- 435 Lene K. Christensen, Shrinivasa K. Upadhyaya, Bernie Jahn, David C.  
436 Slaughter, Eunice Tan, and David Hills. Determining the Influence of  
437 Water Deficiency on NPK Stress Discrimination in Maize using Spec-  
438 tral and Spatial Information. *Precision Agriculture*, 6(6):539–550, De-

439 cember 2005. ISSN 1573-1618. doi: 10.1007/s11119-005-5643-7. URL  
440 <https://doi.org/10.1007/s11119-005-5643-7>.

441 Hui Yan and Heinz W. Siesler. Hand-held near-infrared spectrometers: State-  
442 of-the-art instrumentation and practical applications. *NIR news*, 29(7):  
443 8–12, 2018. Publisher: SAGE Publications Sage UK: London, England.

444 Krzysztof B. Beć, Justyna Grabska, Heinz W. Siesler, and Christian W. Huck.  
445 Handheld near-infrared spectrometers: Where are we heading? *NIR news*,  
446 31(3-4):28–35, 2020. Publisher: SAGE Publications Sage UK: London,  
447 England.

448 Puneet Mishra, Santosh Lohumi, Haris Ahmad Khan, and Alison Nor-  
449 don. Close-range hyperspectral imaging of whole plants for digital phe-  
450 notyping: Recent applications and illumination correction approaches.  
451 *Computers and Electronics in Agriculture*, 178:105780, November 2020.  
452 ISSN 01681699. doi: 10.1016/j.compag.2020.105780. URL [https://](https://linkinghub.elsevier.com/retrieve/pii/S016816992031869X)  
453 [linkinghub.elsevier.com/retrieve/pii/S016816992031869X](https://linkinghub.elsevier.com/retrieve/pii/S016816992031869X).

454 Fabio Fiorani and Ulrich Schurr. Future Scenarios for Plant Phenotyping.  
455 *Annual Review of Plant Biology*, 64(1):267–291, April 2013. ISSN  
456 1543-5008, 1545-2123. doi: 10.1146/annurev-arplant-050312-120137.  
457 URL [http://www.annualreviews.org/doi/abs/10.1146/](http://www.annualreviews.org/doi/abs/10.1146/annurev-arplant-050312-120137)  
458 [annurev-arplant-050312-120137](http://www.annualreviews.org/doi/abs/10.1146/annurev-arplant-050312-120137).

459 Ewa Szymańska. Modern data science for analytical chemical data – A com-  
460 prehensive review. *Analytica Chimica Acta*, 1028:1–10, October 2018.

461 ISSN 0003-2670. doi: 10.1016/j.aca.2018.05.038. URL <http://www.sciencedirect.com/science/article/pii/S0003267018306421>.

463 Svante Wold, Michael Sjöström, and Lennart Eriksson. PLS-regression: a basic tool of chemometrics. *Chemometrics and intelligent laboratory systems*, 58(2):109–130, 2001.

466 Howard Mark and Jerry Workman. *Chemometrics in spectroscopy*. Elsevier/Academic Press, Amsterdam, 2007. ISBN 978-0-12-374024-3. OCLC: 255877127.

469 Pierre Dardenne, George Sinnaeve, and Vincent Baeten. Multivariate Calibration and Chemometrics for near Infrared Spectroscopy: Which Method? *Journal of Near Infrared Spectroscopy*, 8(4):229–237, October 2000. ISSN 0967-0335, 1751-6552. doi: 10.1255/jnirs.283. URL <http://journals.sagepub.com/doi/10.1255/jnirs.283>.

474 E. Bertran, M. Blanco, S. MasPOCH, M. C. Ortiz, M. S. Sánchez, and L. A. Sarabia. Handling intrinsic non-linearity in near-infrared reflectance spectroscopy. *Chemometrics and Intelligent Laboratory Systems*, 49(2): 215–224, October 1999. ISSN 0169-7439. doi: 10.1016/S0169-7439(99)00043-X. URL <http://www.sciencedirect.com/science/article/pii/S016974399900043X>.

480 Wangdong Ni, Lars Nørgaard, and Morten Mørup. Non-linear calibration models for near infrared spectroscopy. *Analytica Chimica Acta*, 813:1–14, February 2014. ISSN 0003-2670. doi: 10.1016/j.aca.2013.

- 483 12.002. URL [http://www.sciencedirect.com/science/article/pii/](http://www.sciencedirect.com/science/article/pii/S0003267013015158)  
484 [S0003267013015158](http://www.sciencedirect.com/science/article/pii/S0003267013015158).
- 485 D. Pérez-Marín, A. Garrido-Varo, and J. E. Guerrero. Non-linear regression  
486 methods in NIRS quantitative analysis. *Talanta*, 72(1):28–42, April 2007.  
487 ISSN 0039-9140. doi: 10.1016/j.talanta.2006.10.036. URL [http://www.](http://www.sciencedirect.com/science/article/pii/S0039914006007119)  
488 [sciencedirect.com/science/article/pii/S0039914006007119](http://www.sciencedirect.com/science/article/pii/S0039914006007119).
- 489 F. Davrieux, D. Dufour, P. Dardenne, J. Belalcazar, M. Pizarro, J. Luna,  
490 L. Londoño, A. Jaramillo, T. Sanchez, N. Morante, F. Calle, L.A. Be-  
491 cerra Lopez-Lavalle, and H. Ceballos. LOCAL Regression Algorithm  
492 Improves near Infrared Spectroscopy Predictions When the Target Con-  
493 stituent Evolves in Breeding Populations. *Journal of Near Infrared Spec-*  
494 *troscopy*, 24(2):109–117, April 2016. ISSN 0967-0335, 1751-6552. doi:  
495 10.1255/jnirs.1213. URL [http://journals.sagepub.com/doi/10.1255/](http://journals.sagepub.com/doi/10.1255/jnirs.1213)  
496 [jnirs.1213](http://journals.sagepub.com/doi/10.1255/jnirs.1213).
- 497 Tormod. Naes, Tomas. Isaksson, and Bruce. Kowalski. Locally weighted  
498 regression and scatter correction for near-infrared reflectance data. *Ana-*  
499 *lytical Chemistry*, 62(7):664–673, April 1990. ISSN 0003-2700, 1520-6882.  
500 doi: 10.1021/ac00206a003. URL [https://pubs.acs.org/doi/abs/10.](https://pubs.acs.org/doi/abs/10.1021/ac00206a003)  
501 [1021/ac00206a003](https://pubs.acs.org/doi/abs/10.1021/ac00206a003).
- 502 Maxime Metz, Matthieu Lesnoff, Florent Abdelghafour, Reza Akbarinia,  
503 Florent Masseglia, and Jean-Michel Roger. A “big-data” algorithm  
504 for KNN-PLS. *Chemometrics and Intelligent Laboratory Systems*, 203:  
505 104076, August 2020. ISSN 0169-7439. doi: 10.1016/j.chemolab.2020.

506 104076. URL [http://www.sciencedirect.com/science/article/pii/](http://www.sciencedirect.com/science/article/pii/S0169743920301908)  
507 [S0169743920301908](http://www.sciencedirect.com/science/article/pii/S0169743920301908).

508 R. Core Team. R: A language and environment for statistical computing.  
509 Vienna, Austria: R Foundation for Statistical Computing. *Available*, 2013.

510 Matthew Barker and William Rayens. Partial least squares for discrimina-  
511 tion. *Journal of Chemometrics*, 17(3):166–173, March 2003. ISSN 0886-  
512 9383, 1099-128X. doi: 10.1002/cem.785. URL [http://doi.wiley.com/](http://doi.wiley.com/10.1002/cem.785)  
513 [10.1002/cem.785](http://doi.wiley.com/10.1002/cem.785).

514 R. A. Fisher. The Use of Multiple Measurements in Taxonomic Problems.  
515 *Annals of Eugenics*, 7(2):179–188, 1936. ISSN 2050-1439. doi: [https://doi.](https://doi.org/10.1111/j.1469-1809.1936.tb02137.x)  
516 [org/10.1111/j.1469-1809.1936.tb02137.x](https://doi.org/10.1111/j.1469-1809.1936.tb02137.x). URL [https://onlinelibrary.](https://onlinelibrary.wiley.com/doi/abs/10.1111/j.1469-1809.1936.tb02137.x)  
517 [wiley.com/doi/abs/10.1111/j.1469-1809.1936.tb02137.x](https://onlinelibrary.wiley.com/doi/abs/10.1111/j.1469-1809.1936.tb02137.x).  
518 \_eprint: [https://onlinelibrary.wiley.com/doi/pdf/10.1111/j.1469-](https://onlinelibrary.wiley.com/doi/pdf/10.1111/j.1469-1809.1936.tb02137.x)  
519 [1809.1936.tb02137.x](https://onlinelibrary.wiley.com/doi/pdf/10.1111/j.1469-1809.1936.tb02137.x).

520 Matthieu Lesnoff, Maxime Metz, and Jean-Michel Roger. Comparison of  
521 locally weighted PLS strategies for regression and discrimination on agro-  
522 nomic NIR data. *Journal of Chemometrics*, page 13, 2020.

523 Svante Wold. Cross-Validatory Estimation of the Number of Components  
524 in Factor and Principal Components Models. *Technometrics*, 20(4):397–  
525 405, 1978. ISSN 0040-1706. doi: 10.2307/1267639. Publisher: [Taylor  
526 & Francis, Ltd., American Statistical Association, American Society for  
527 Quality].

528 José Camacho and Alberto Ferrer. Cross-validation in PCA mod-  
529 els with the element-wise k-fold (ekf) algorithm: theoretical as-  
530 pects. *Journal of Chemometrics*, 26(7):361–373, 2012. ISSN  
531 1099-128X. doi: <https://doi.org/10.1002/cem.2440>. URL <https://onlinelibrary.wiley.com/doi/abs/10.1002/cem.2440>.  
532 [\\_eprint: https://onlinelibrary.wiley.com/doi/abs/10.1002/cem.2440](https://onlinelibrary.wiley.com/doi/abs/10.1002/cem.2440).  
533 <https://onlinelibrary.wiley.com/doi/pdf/10.1002/cem.2440>.

534 Jun-Li Xu, Alexia Gobrecht, Daphné Héran, Nathalie Gorretta, Marie  
535 Coque, Aoife A. Gowen, Ryad Bendoula, and Da-Wen Sun. A polar-  
536 ized hyperspectral imaging system for in vivo detection: Multiple applica-  
537 tions in sunflower leaf analysis. *Computers and Electronics in Agriculture*,  
538 158:258–270, March 2019. ISSN 0168-1699. doi: 10.1016/j.compag.2019.  
539 02.008. URL [http://www.sciencedirect.com/science/article/pii/](http://www.sciencedirect.com/science/article/pii/S0168169918314455)  
540 [S0168169918314455](http://www.sciencedirect.com/science/article/pii/S0168169918314455).