



HAL
open science

RoBoost-PLS2-R: An extension of RoBoost-PLSR method for multi-response

Maxime Metz, Maxime Ryckewaert, Sílvia Mas Garcia, Ryad Bendoula, Pierre Dardenne, Matthieu Lesnoff, Jean Michel Roger

► To cite this version:

Maxime Metz, Maxime Ryckewaert, Sílvia Mas Garcia, Ryad Bendoula, Pierre Dardenne, et al..
RoBoost-PLS2-R: An extension of RoBoost-PLSR method for multi-response. *Chemometrics and
Intelligent Laboratory Systems*, 2022, 222, 10.1016/j.chemolab.2022.104498 . hal-03582272

HAL Id: hal-03582272

<https://hal.inrae.fr/hal-03582272>

Submitted on 22 Jul 2024

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 4.0 International License

1 RoBoost-PLS2-R : An extension of RoBoost-PLSR
2 method for multi-response

3 Maxime Metz^{a,b}, Maxime Ryckewaert^{a,b}, Silvia Mas Garcia^{a,b}, Ryad
4 Bendoula^{a,b}, Pierre Dardenne^{d,b}, Matthieu Lesnoff^{c,b}, Jean-Michel Roger^{a,b}

^a*ITAP Univ Montpellier INRAE Institut Agro Montpellier France*

^b*ChemHouse Research Group Montpellier France*

^c*SELMET Univ Montpellier CIRAD INRAE Montpellier SupAgro Montpellier France*

^d*Wallon Agricultural Research Centre Gembloux*

5 **Abstract**

6 Recently, a novel robust PLSR method was developed to address the
7 problem of outliers in the data. In this paper, an extension of this method,
8 called RoBoost-PLS2-R is proposed to predict multi-response variables.
9 Robustness and efficiency of this new approach have been validated on
10 two simulated data sets and one real data set containing different outlier
11 scenarios. Its performance was also compared with reference methods
12 (PLS2-R and RSIMPLS) for predicting multi-response variables. Results
13 confirm that RoBoost-PLS2-R greatly reduces prediction errors when data
14 contain outliers. Prediction performances of RoBoost-PLS2-R are close
15 to the optimal model (PLS2-R) calibrated without outliers and also to
16 RSIMPLS method. This method seems to be a reliable and a competitive
17 robust regression tool for predicting multi-response variables.

18 *Keywords:* Robust regression methods, outliers, multi-response,
19 multivariate data analysis

20 **1. Introduction**

21 Partial Least Square Regression (PLSR) [1] is a common data analysis
22 method and a well-established tool in chemometrics. PLSR calculates
23 a linear relationship between explanatory variables (\mathbf{X}) and response
24 variables (\mathbf{Y}). PLSR can be used to predict one response (PLS1) or
25 several responses (PLS2). PLSR is particularly useful for processing
26 high-dimensional data, especially when the number of explanatory variables

Preprint submitted to Chemometrics and Intelligent Laboratory Systems 13 janvier 2022

27 exceeds the number of samples. This method is widely used in analytical
28 chemistry for predicting constituent concentrations of a sample based on its
29 spectrum obtained by spectroscopic techniques, such as near-infrared (NIR)
30 spectroscopy, Fluorescence spectroscopy and ultraviolet (UV) spectroscopy.
31 The PLSR model is known to be affected by the presence of atypical
32 observations (outliers) in the data set. Outliers can negatively affect the
33 calibration of PLSR models. To deal with outliers, several robust PLSR
34 methods were proposed in the literature [2–12]. These methods were
35 particularly developed to deal with outliers when the response matrix is
36 uni-dimensional [13] (PLS1-R). However, robust methods that address the
37 case of multi-responses (PLS2) are few. Among them, RSIMPLS is one
38 of the most used method [14]. RSIMPLS proposes to robustly estimate
39 the cross-covariance matrix \mathbf{C}_{xy} and the empirical covariance matrix \mathbf{C}_x
40 used in SIMPLS algorithm. For this, a robust principal component analysis
41 (ROBPCA) is performed on the concatenated data matrix of \mathbf{X} and \mathbf{Y} .
42 RSIMPLS uses additional information from the previous ROBPCA step to
43 perform a reweighted multiple linear regression.
44 Recently, a new robust method called RoBoost-PLSR has been developed
45 [15]. RoBoost-PLSR aims at determining the measure of relevance of
46 the samples for PLSR model calibration. Indeed, in practical cases, the
47 samples of a database are not defined as outliers, i.e. not relevant for the
48 calibration of a PLSR model. RoBoost-PLSR proposes to calculate a weight
49 on each latent variable to define the relevance of the samples. The relevance
50 measurement is defined according to three criteria calculated for each latent
51 variable (X -residuals, Y -residuals, leverage). This method has proven to be
52 effective for outliers in both \mathbf{Y} and \mathbf{X} . However, this algorithm was only
53 developed for a one-dimensional PLSR response variable (PLS1). This paper
54 contributes to the RoBoost-PLSR method which will be able to manage
55 outliers in a multiple response context.
56 The first section introduces the extension of RoBoost-PLSR named
57 RoBoost-PLS2-R and the associated algorithm. The following section
58 presents the data and the methods used to evaluate and compare the
59 predictive ability of RoBoost-PLS2-R. Finally, the prediction performance
60 of RoBoost-PLS2-R and its comparison with standard methods are shown
61 in the last section.

62 **2. Notations**

63 Capital bold characters will be used for matrices, *e.g.* \mathbf{X} ; small bold
64 characters for column vectors, *e.g.* \mathbf{x}_j will denote the j^{th} column of \mathbf{X} ; row
65 vectors will be denoted by the transpose notation, *e.g.* \mathbf{x}_i^T will denote the i^{th}
66 row of \mathbf{X} ; italic characters will be used for scalars, *e.g.* matrix elements x_{ij}
67 or indices i . Constant scalars will be denoted with italicised characters, *e.g.*
68 number of samples n . $\mathbf{1}$ will represent a column vector of ones, of proper
69 dimension. *med* defines the median. \mathbf{X} and \mathbf{Y} are the spectral and the
70 responses matrices. g is the weight function. \mathbf{D} is the matrix of sample weights
71 where the diagonal of the matrix is the sample weight and the other terms
72 are zero.

73 **3. RoBoost-PLSR extension for multi-responses**

74 *3.1. Algorithm*

75 The new algorithm allowing an extension in a multi-response context is
76 the following :

Algorithm RoBoost-PLSR for K LV

For a definite number of K latent variables, the algorithm proceeds as described below :

1: Initialisation step

$$k = 1$$

$$\mathbf{D} = \text{diag}(d_1, d_2, \dots, d_n) \text{ with } d_i = \frac{1}{n}$$

2: Center the data :

$$\mathbf{X}_k = \mathbf{X} - \mathbb{1}\mathbb{1}^T\mathbf{D}\mathbf{X}$$

$$\mathbf{Y}_k = \mathbf{Y} - \mathbb{1}\mathbb{1}^T\mathbf{D}\mathbf{Y}$$

3: Define \mathbf{u}_k as an arbitrary column of \mathbf{Y}

4: Calculate one weighted latent variable NIPALS :

$$\mathbf{w}_k = \frac{\mathbf{X}_k^T \mathbf{D} \mathbf{u}_k}{\|\mathbf{X}_k^T \mathbf{D} \mathbf{u}_k\|}$$

$$\mathbf{t}_k = \mathbf{X}_k \mathbf{w}_k$$

$$\mathbf{p}_k = \frac{\mathbf{X}_k^T \mathbf{D} \mathbf{t}_k}{\mathbf{t}_k^T \mathbf{D} \mathbf{t}_k}$$

$$\mathbf{q}_k = \frac{\mathbf{Y}_k^T \mathbf{D} \mathbf{t}_k}{\mathbf{t}_k^T \mathbf{D} \mathbf{t}_k}$$

$$c_k = \frac{\mathbf{u}_k^T \mathbf{D} \mathbf{t}_k}{\mathbf{t}_k^T \mathbf{D} \mathbf{t}_k}$$

5: Derive (\mathbf{F}) , (\mathbf{E}) , (\mathbf{l}) :

$$\mathbf{E} = \mathbf{X}_k - \mathbf{t}_k \mathbf{p}_k^T$$

$$\mathbf{F} = \mathbf{Y}_k - \mathbf{t}_k \mathbf{q}_k^T$$

$$\mathbf{l} = \mathbf{t}_k$$

6: Update the weights for each $i \in [1, n]$ sample :

$$d_i = \frac{1}{n} \times g(\|\mathbf{e}_i\|, \alpha) \times \prod_{j=1}^m g(f_{ij}, \beta), \times g(l_i, \gamma)$$

7: Return to (step (2) for $k = 1$, otherwise return to step (4)) until convergence of successive c 's.

8: Deflation step

$$\mathbf{X}_{k+1} = \mathbf{X}_k - \mathbf{t}_k \mathbf{p}_k^T$$

$$\mathbf{Y}_{k+1} = \mathbf{Y}_k - \mathbf{t}_k \mathbf{q}_k^T$$

$$\mathbf{u}_{k+1} = \mathbf{Y}_k \mathbf{q}_k$$

set $k = k + 1 \rightarrow$ then go to step (4)

The regression coefficients resulting for K latent variables are estimated as follows :

$$\mathbf{B} = \mathbf{R} \mathbf{c}^T$$

With \mathbf{R} :

$$\mathbf{R} = \mathbf{W}(\mathbf{P}^T \mathbf{W})^{-1}$$

77 *3.2. Theoretical discussions*

78 The algorithm RoBoost-PLS2-R have similar properties to the algorithm
79 proposed in [15], but also new properties :

80

81 — The RoBoost-PLS2-R framework is designed foremost to facilitate
82 the leverage measurement. Leverage is defined as the distance to the
83 centre of the model (see step 6 in the algorithm). In usual strategies,
84 to define distances between the model centre and individuals,
85 different metrics can be used. Euclidean or Mahalanobis distances

86 between scores and the model centre are strategies commonly used
87 in chemometrics. However, in the case of a Euclidean distance, the
88 latest LVs could have a minor contribution to the leverage value.
89 This is due to the decreasing magnitude of scores. Nevertheless,
90 the predictive potential of these latest LVs may not be necessarily
91 negligible. In the case of a Mahalanobis distance, contributions of
92 all LVs become equal in the computation of the leverage value. This
93 can be also detrimental, since the predictive potentials of the LVs are
94 most usually uneven. Considering these limitations, RoBoost-PLSR
95 proposes to estimate the sample leverage for each latent variable. This
96 avoids the need to define specific metrics for the leverage calculation.
97 However, the use of this strategy may make it difficult to assign a low
98 weight to individuals with a leverage effect that is only identifiable
99 with a large number of latent variables.

- 100
101 — The proposed method takes into account X-residuals (see step 6 in
102 the algorithm). Usually only Y-residuals are considered in robust
103 PLS approaches. The inclusion of these residuals provides additional
104 information that cannot be expressed by leverage and Y-residuals
105 alone.
- 106
107 — The algorithm proposed in this article provides regression coefficients.
108 This makes the constructed RoBoost-PLSR models more easily
109 interpretable. Contrary to the first algorithm proposed in [15], the
110 rotation matrix \mathbf{R} used to estimate the regression coefficients can be
111 estimated. This is due to data centring which is only done for the
112 first model with a single latent variable. In the previous algorithm,
113 repeated centring of \mathbf{X} and \mathbf{Y} matrices led to a bias which made it
114 impossible to estimate the rotation matrix.
- 115
116 — Like PLSR, RoBoost-PLSR makes it possible to deduce any of the
117 1 to K LVs models from the calibration of a single K LVs model.
118 This preserves the operability during validation and parameterisation
119 process of the RoBoost-PLSR method. Indeed, from this set of
120 one-variable latent models it is possible to define the rotation matrix
121 \mathbf{R} which enables to compute all previous PLS models.
- 122
123 — The algorithm proposed in [15], determines the convergence with q .

124 However, \mathbf{q} is multidimensional when \mathbf{Y} is multidimensional. In the
 125 new algorithm convergence estimation is facilitated by using c which
 126 is a scalar when responses matrix \mathbf{Y} is multidimensional (see step 7
 127 in the algorithm).

128 — The weights of the sample according to the \mathbf{Y} -residuals are the
 129 product of the estimated weights for each \mathbf{Y} -variable (see step 6 in
 130 the algorithm). A specific sample weight for each residual of each
 131 \mathbf{Y} variable is calculated and then multiply them to give an overall
 132 weight. This strategy enables sample weights to be estimated in a
 133 way that is appropriate to the multivariate nature of \mathbf{Y} . This strategy
 134 takes in consideration the fact that \mathbf{Y} variables may have different
 135 variances. If this aspect is not taking into account, some outliers could
 136 be considered as inliers by the method. For instance, atypical samples
 137 on a specific variable of \mathbf{Y} can mask the outliers of other columns
 138 of \mathbf{Y} which present a lower variability. This strategy also allows a
 139 fast operation by applying the bisquare function on each column of
 140 Y -residuals matrix for each LV according to the β hyperparameter.
 141 Finally, the global weights associated with Y -residuals are defined as
 142 a product of each weight calculated on the Y -residual. This strategy
 143 of combining weights is a commonly used strategy. It is basically
 144 used to combine the weights calculated according to the three criteria
 145 (X -residuals, Y -residuals, leverage) in RoBoost-PLSR. However,
 146 different strategies are possible. Like calculating the Mahalanobis
 147 distances on \mathbf{Y} or making a combination of weights different from the
 148 product. In particular, it is possible to perform a sum of weights, so
 149 that the weighting strategy can eliminate individuals who only have
 150 weights at 0 for each criterion.

151 — In this article, the weight function g is the bisquare function :

$$B(z_i) = (1 - z_i^2)^2 \text{ for } |z_i| < 1 \text{ and } B(z_i) = 0 \text{ for } |z_i| > 1$$

with z_i :

$$\frac{x_i}{c \times \text{med}(|\mathbf{x}|)}$$

152 However, any weight function can be considered and tested in order
 153 to improve the algorithm to obtain better predictive capacity. In
 154 RoBoost-PLS2-R x_i (associated with the bisquare function) is specific

155 according to the chosen statistic. This means that when the weights
156 are calculated according to the residuals of \mathbf{X} , x_i corresponds to
157 the norm of the vector \mathbf{e}_i and \mathbf{x} to the norms of the individuals of
158 \mathbf{E} . When the residuals \mathbf{Y} are taken into account, x_i is the value
159 of the residual y_{ij} and \mathbf{x} is the vector of residuals f_j . Finally, the
160 leverage effect is taken into account, x_i corresponds to the score of
161 a latent variable t_{ik} and \mathbf{x} is the vector of scores \mathbf{t}_k for all samples.
162 Furthermore, the constant c in the bisquares function corresponds to
163 the parameters α , β and γ in step 6 of the algorithm. This constant
164 has to be adjusted according to the type of outlier.

165 4. Materials and methods

166 4.1. Simulated Data

167 To evaluate the performance of RoBoost-PLS2-R in comparison with
168 standard PLS2-R and RSIMPLS, two simulations were performed. The
169 first simulation represents the Y -outlier case and the second simulation the
170 X -outlier case. For each simulation, 1000 samples were generated according
171 to the framework proposed by [16]. Among these samples, 200 outliers
172 were generated. The spectral signatures used for the simulations were the
173 spectral signatures of water, ethanol and glucose estimated in [16]. Using
174 this approach, the matrix of explanatory variables (\mathbf{X}) was generated by :

$$\mathbf{X} = \mathbf{t}_u \mathbf{p}_u^t + \mathbf{T}_d \mathbf{P}_d^t + \mathbf{E} \quad (1)$$

175 And the relationship f between \mathbf{X} and \mathbf{Y} by :

$$\mathbf{Y} = f(\mathbf{t}_u) + \mathbf{F} \quad (2)$$

176 Where \mathbf{p}_u is the spectral signature in the useful space and \mathbf{P}_d are
177 spectral signatures in the detrimental space. \mathbf{t}_u and \mathbf{T}_d are their associated
178 contributions. The \mathbf{E} and \mathbf{F} matrices are defined as gaussian noises of \mathbf{X}
179 and \mathbf{Y} , respectively.

180 The parameters of the simulations are represented in tables (Table 1
181 and Table 2) where differences between simulated inliers and outliers were
182 highlighted in bold in the tables. Scripts of the simulations are available at
183 this link : https://github.com/maxmetz/data_simulation

184 4.1.1. Simulation 1, Y-outliers

185 The Y-outliers were defined by their relationship f between \mathbf{X} and \mathbf{Y} .
 186 All other simulation parameters were common between inliers and outliers.
 187 The construction of the simulated data set 1 is represented in table 1.

TABLE 1 – The different choices in the simulation 1

	Inliers	Outliers
\mathbf{P}_u	Pure spectrum of glucose	
\mathbf{t}_u	Folded-normal distribution	
\mathbf{P}_d	Pure spectrum of water Pure spectrum of ethanol Spectrum of water-ethanol Interaction 10 Artificial spectra	
\mathbf{T}_d	Folded-normal distribution Folded-normal distribution Product between T_{water} and $T_{ethanol}$ Folded-normal distribution	
\mathbf{E}	Gaussian distribution	
f	$Y_1 = 10 * T_{ethanol}$ $Y_2 = 10 * T_{glucose}$ $Y_3 = 10 * T_{water}$	$Y_1 = 10 * T_{ethanol}$ $\mathbf{Y}_2 = -10 * \mathbf{T}_{glucose}$ $Y_3 = 10 * T_{water}$
\mathbf{F}	Gaussian distribution	

188 4.1.2. Simulation 2, X-outliers

189 The X-outliers were defined by others artificial spectral signatures.
 190 These signatures correspond to minority compounds. All other simulation
 191 parameters were common between inliers and outliers. The simulation is
 192 represented in table 2.

TABLE 2 – The different choices in the simulation 2

	Inliers	Outliers
\mathbf{P}_u	Pure spectrum of glucose	
\mathbf{t}_u	Folded-normal distribution	
\mathbf{P}_d	Pure spectrum of water Pure spectrum of ethanol Spectrum of water-ethanol Interaction 10 Artificial spectra	Pure spectrum of water Pure spectrum of ethanol Spectrum of water-ethanol Interaction 10 Artificial spectra 10 Artificial spectra
\mathbf{T}_d	Folded-normal distribution Folded-normal distribution Product between T_{water} and $T_{ethanol}$ Folded-normal distribution	Folded-normal distribution Folded-normal distribution Product between T_{water} and $T_{ethanol}$ Folded-normal distribution Folded-normal distribution
\mathbf{E}	Gaussian distribution	
f	$Y_1 = 10 * T_{ethanol}$ $Y_2 = 10 * T_{glucose}$ $Y_3 = 10 * T_{water}$	$Y_1 = 10 * T_{ethanol}$ $Y_2 = 10 * T_{glucose}$ $Y_3 = 10 * T_{water}$
\mathbf{F}	Gaussian distribution	

193 *4.2. Real data set*

194 The real data set was formed by 261 spectra of raw cow milk collected
195 from farms in Wallonia in 2014 and 2015. Spectra were recorded over
196 a spectral range 397-4000 cm-1 with a resolution of 4 cm-1 by using a
197 FTIR spectrometer (Delta LactoScope, PerkinElmer). For each sample,
198 chemical measurements were performed to obtain two-responses variable :
199 fat content and protein content. Fat and Protein content were determined in
200 accordance with reference methods "ISO 1211 :2010 [IDF 1 :2010]" and "ISO
201 8968-1 :2014 [IDF 20-1 :2014]", respectively. This database is particularly
202 interesting because it contains missing data whose values have been replaced
203 by 0.

204 *4.2.1. Evaluation strategies*

205 RoBoost-PLS2-R was evaluated and compared with two standard
206 regression algorithms : PLS2-R and RSIMPLS.

207 In the case of the simulations, the 1000 samples were divided into two
208 groups : 800 for calibration and 200 for validation. The reference method in

209 terms of prediction performance was PLS2-R calibrated without outliers. For
210 the real data set, calibration set was composed of 209 samples. The validation
211 was conducted on 52 samples. These samples were selected from a study of the
212 data in order to represent the samples as well as possible without containing
213 potential outliers. The reference method in terms of prediction performance
214 was RSIMPLS.

215 The method performance was evaluated according to the validation sets
216 and Root Mean Square Error of Prediction (RMSEP) as a figure of merit.
217 Only the results achieved using the optimal parameters (*i.e.* the parameters
218 that provide the minimum value of the RMSEP) of RoBoost-PLS2-R and
219 RSIMPLS were presented.

220 The evaluation strategy also aimed at assessing the weights attributed to
221 each sample. Indeed, the RoBoost-PLS2-R method allows the visualisation
222 of the weight given to each sample for each LV. In this work, the parameters
223 of the methods RoBoost-PLS2-R and RSIMPLS such as the constants used
224 in the weight functions were adjusted to obtain the minimum RMSEP.

225 *4.3. Software*

226 PLS2-R was performed with “[rnirs](#)” and RoBoost-PLS2-R is available
227 [RoBoost-PLSR](#) functions available in R. RSIMPLS was performed using the
228 function of the LIBRA package available in MALTLAB.

229 **5. Results and discussions**

230 *5.1. Simulation set 1*

231 *5.1.1. Data visualisation*

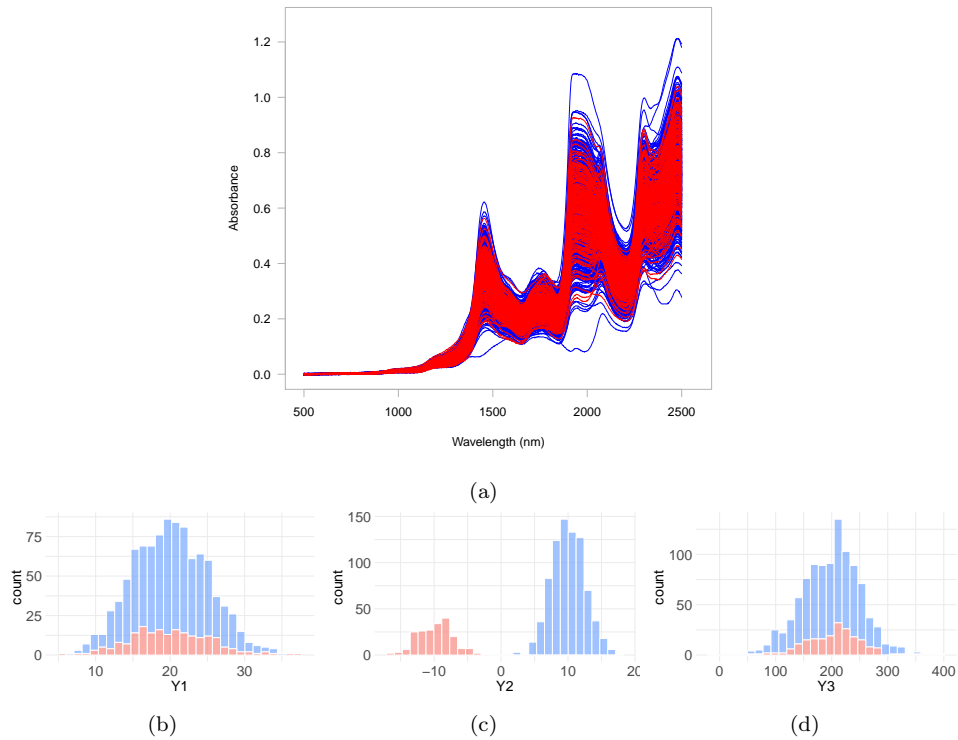


FIGURE 1 – Graphical representation of simulation 1 : (a) spectral data (b) value distribution of Y1 response variable (c) value distribution of Y2 response variable (d) value distribution of Y3 response variable. Outliers are shown in red and inliers in blue.

232 Figure 1 shows the graphical representation of simulation 1. From the
233 spectra plot (Figure 1a), it can be seen that is difficult to identify outliers
234 (in red) from a simple visual inspection. In this case, the outliers were
235 defined by a distinct relation f on one of the response variables (see Table
236 1). Therefore, no spectral difference between the two groups is expected.
237 From the plot of value distributions of the response variables (see Figure
238 1b,c,d) it can be observed that Y1 and Y3 variables present the same
239 distribution for both outliers and inliers. However, different distribution for

240 these two groups is presented in Y2 variable. Moreover, the variances of Y1
 241 are smaller than the variance of Y3.

242 *5.1.2. Method evaluation*

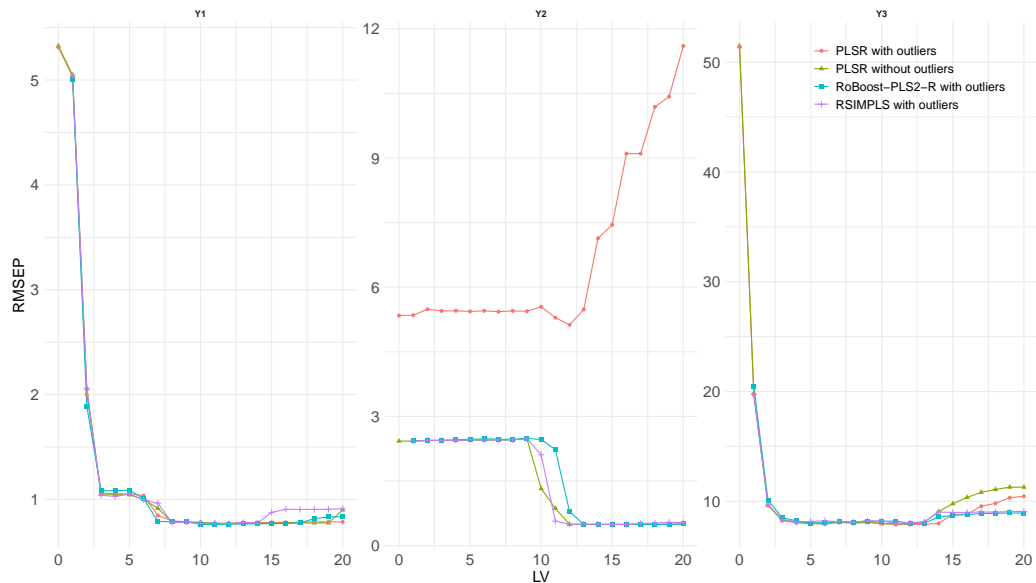


FIGURE 2 – Evolution of the RMSEP as a function of the number of latent variables for the PLS2-R with and without outliers, RSIMPLS and RoBoost-PLS2-R for the simulation 1 set

243 Figure 2 shows the prediction performances for each method and response
 244 variable Y on the basis of simulation 1. For the variables Y1 and Y3, the
 245 error curves obtained by PLS2-R with and without outliers, RSIMPLS and
 246 RoBoost-PLS2-R are similar. This is due to the fact that outliers are only
 247 atypical on Y2 and hence, no impact on the Y1 and Y3 predictions is
 248 expected. For the variables Y2 the error curves obtained by PLS2-R with
 249 and without outliers are different. The PLS2-R model calibrated with outliers
 250 perform poorly in inliers prediction. The prediction performance of RSIMPLS
 251 is close to the PLS2-R without outliers. This means that the RSIMPLS
 252 method can deal with these outliers and provides satisfactory results. These
 253 results show that RoBoost-PLS2-R performs as well as RSIMPLS on this
 254 dataset. Therefore, RoBoost-PLS2-R can handle the presence of outliers in
 255 the response variables regardless of the variance of the responses.

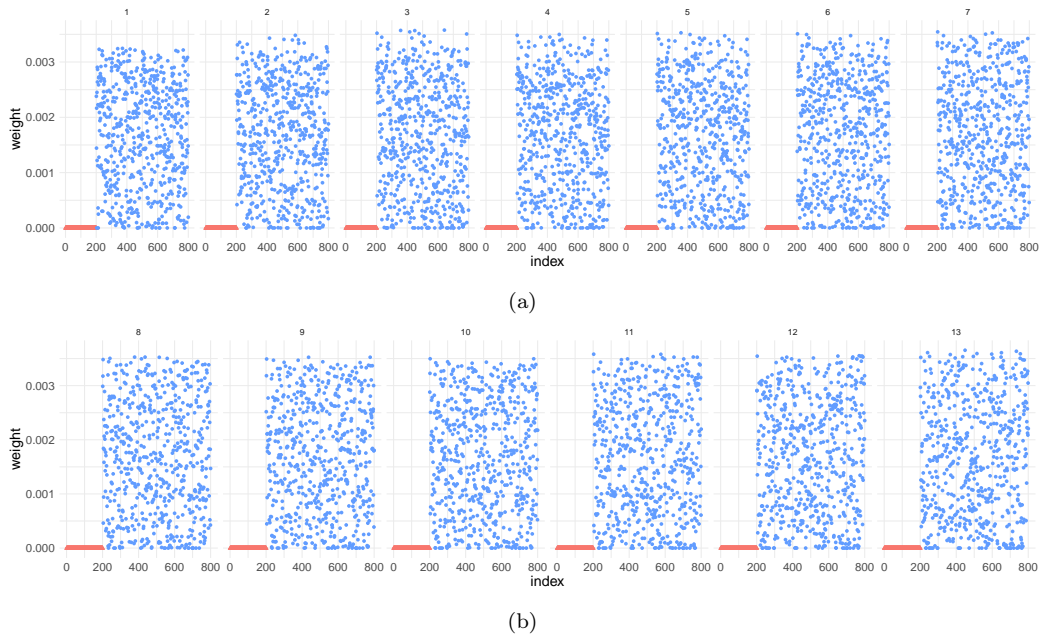


FIGURE 3 – Weights assigned to samples by the RoBoost-PLS2-R method for the simulation set 1 according to the number of LV from 1 to 13. Outliers and inliers are in red and blue, respectively.

256 Figure 3 shows the weights assigned to the samples of simulation 1 by
 257 the RoBoost-PLS2-R method as a function of the number of LV with the
 258 best performing hyperparameters. It can be noted that outliers have a very
 259 low weight while some inliers have a weight close to zero. This may be due
 260 to three reasons. Firstly, the hyperparameters of bisquare function must
 261 be strict enough to assign a weight close to 0 to the outliers for each LV.
 262 Taking into account that some inliers could be very similar to some outliers,
 263 assignation of low weights to these inliers could be expected. Secondly, the
 264 weights associated to Y -residuals are a combination of weights defined for
 265 each Y variables. The hyperparameter β (see Section 3) is assumed to be
 266 constant for each variable in Y . This means that the higher the number of
 267 variables, the more dispersed the weights assigned to the inliers could be.
 268 To achieve a more homogeneous weighting on the outliers, the multivariate
 269 aspect of Y should be taken into account. For example, a potential solution
 270 can be to calculate the robust Mahalanobis distance at the centre of the data
 271 on the residuals of Y for each Latent Variable. Thirdly, some outliers are not
 272 detrimental to the model but are also irrelevant and can therefore have a

273 low weight without impacting on the prediction performance of the model.
 274 In conclusion, RoBoost-PLS2-R has assigned a low weight to a large number
 275 of samples without impacting on the prediction performance of the model.
 276 However, it is potentially possible to improve this approach by modifying the
 277 weighting criteria associated with the Y residuals.

278 *5.2. Simulation 2*

279 *5.2.1. Data visualisation*

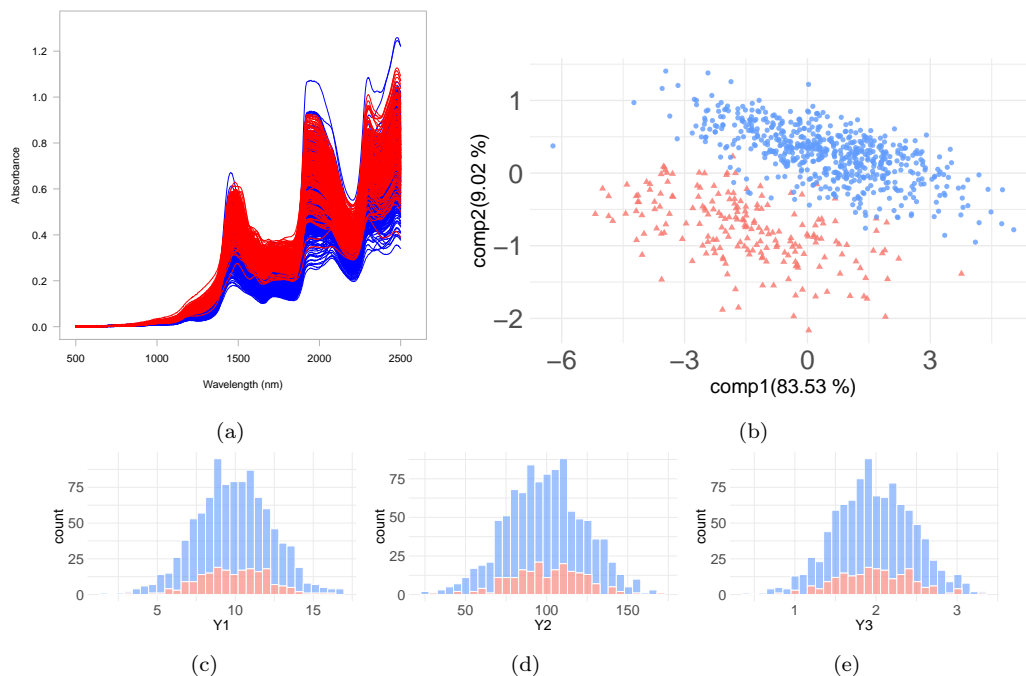


FIGURE 4 – Graphical representation of simulation 2 : (a) spectral data (b) PCA score plot of the two first components (c) value distribution of Y1 response variable (d) value distribution of Y2 response variable (e) value distribution of Y3 response variable. Outliers are shown in red and inliers in blue.

280 Figure 4 shows the graphical representation of simulation 2. From spectra
 281 plot of the sample (Figure 4 a), it can be seen that outliers are not identifiable.
 282 Indeed, in this simulation, outliers are different only for spectral signatures
 283 and hence, they contribute slightly to the construction of the spectra. Figures
 284 4b represents the score plot on the two first principal components. Two

285 centroids can be seen but there is no clear separation between outliers and
 286 inliers. This is due to the outliers having their major compounds in common
 287 (see Table 2). From the value distributions plot of the responses (see : Figures
 288 4c,d,e), it can be seen that outliers and inliers present similar distribution
 289 in all Y response variables. Outliers are different only on the basis of the
 290 spectral signatures that compose them.

291 *5.2.2. Method evaluation*

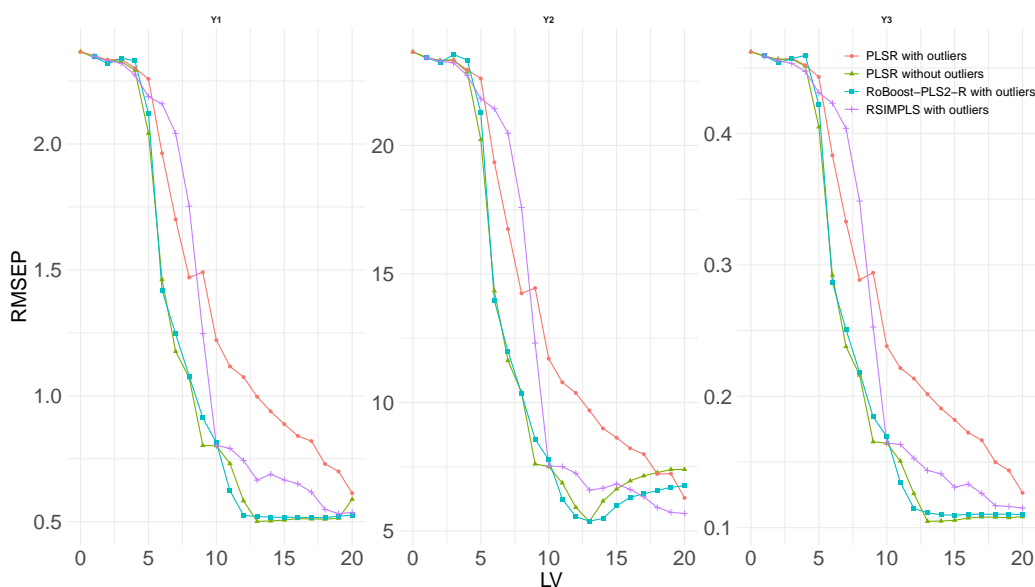


FIGURE 5 – Evolution of the RMSEP as a function of the number of latent variables for the PLS2-R with and without outliers, RSIMPLS and RoBoost-PLS2-R for the simulation 2 set

292 The figure 5 represents the prediction performances of the applied
 293 methods on validation set for each response variable on the basis of the
 294 simulation. As expected, the outliers impact negatively the predictive
 295 capacity of the PLS2-R for all responses. For the RSIMPLS method, all
 296 performance curves are between those of the PLS2-R method with and
 297 without outliers. However, with a large number of latent variables, the
 298 prediction performances of RSIMPLS approach the best performance of
 299 PLS2-R without outliers. This may be due to the fact that RSIMPLS does
 300 not directly take into account the residuals of X but also that the estimation

301 of the leverage effect is not directly taken into account. Indeed, in RSIMPLS
 302 it is the cross-covariance matrices $\mathbf{C}_{\mathbf{xy}}$ and the empirical covariance matrix
 303 $\mathbf{C}_{\mathbf{x}}$ that are robustly estimated.

304 For the RoBoost-PLS2-R method, it can also be seen that for the three
 305 responses, performance curves are close to those of PLS2-R without outliers.
 306 However the optimal number of components is higher for RoBoost-PLS2-R
 307 than the PLS2-R without outliers. To conclude, these results highlight
 308 the fact that RoBoost-PLS2-R can reach the best performance of PLS2-R
 309 without outliers. Thus, RoBoost-PLS2-R can handle these X -outliers for
 310 the prediction of multiple responses.

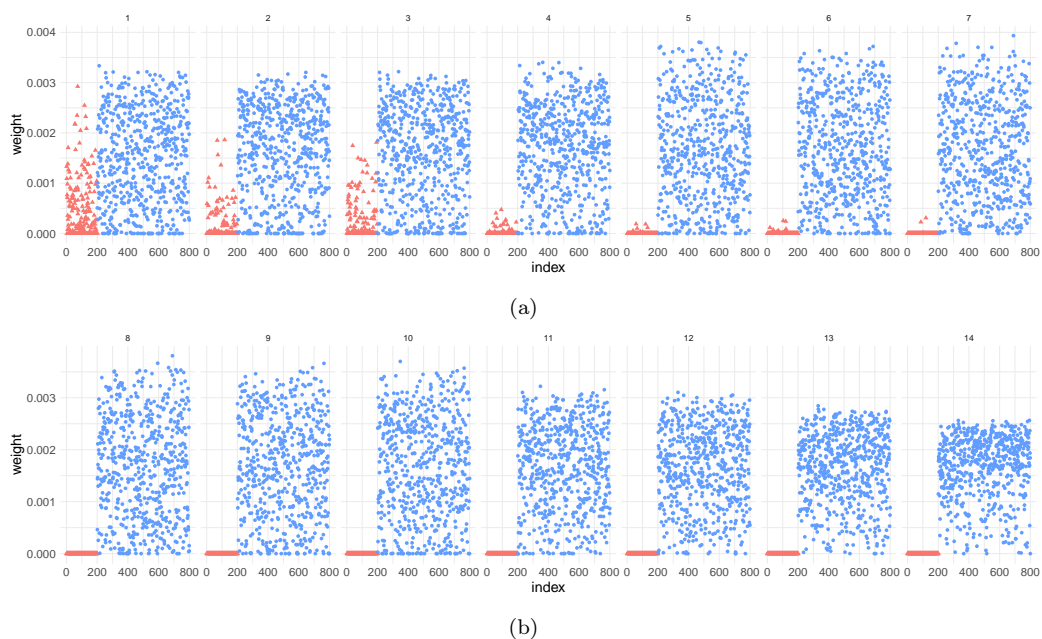


FIGURE 6 – Weights assigned to samples in simulation set 2 according to the chosen number of latent variables from 1 to 14. Outliers and inliers are in red and blue, respectively

311 Figure 6 shows the weight assigned to samples by RoBoost-PLS2-R
 312 according to the number of LV. It can be observed that the weights of
 313 outliers decrease progressively when the number of LV increases. This
 314 gradual decrease is partly explained by the fact that both outliers and inliers
 315 were simulated using common majority spectral signatures. Indeed, only
 316 some minor spectral signatures differentiate the inliers from the outliers (see
 317 Section 4). After 8 latent variables, all outliers have a weight equal to 0,

318 whereas almost all inliers present a high weight. Nevertheless, it is possible
 319 to note that the majority of the inliers have a strong weight and therefore a
 320 large number of them are used to calculate the model.

321 5.3. Real data set

322 5.3.1. Data visualisation

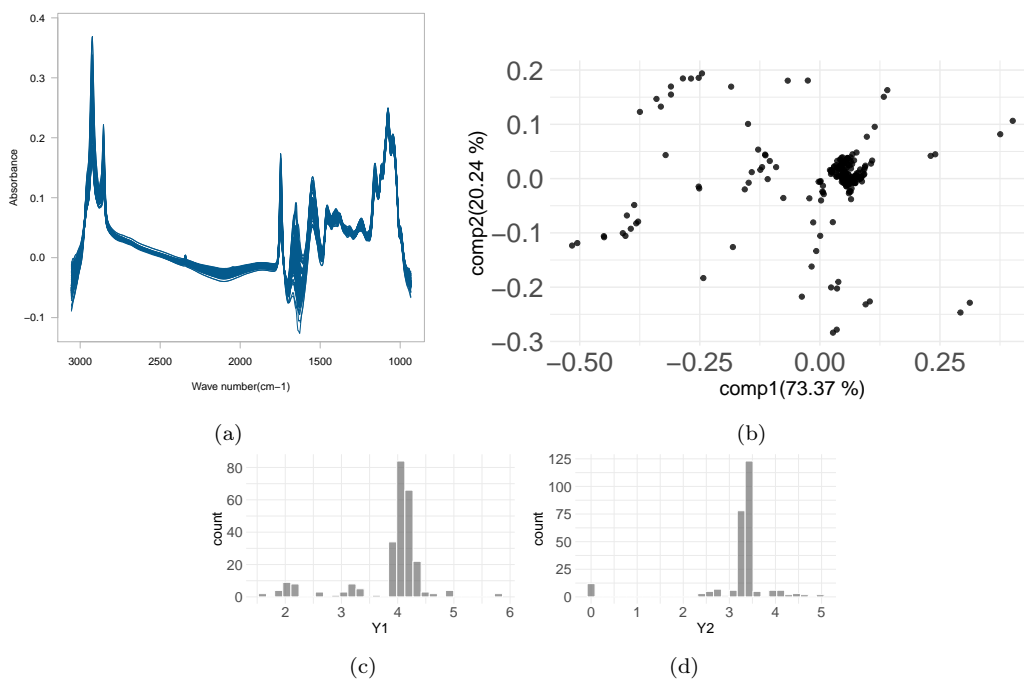


FIGURE 7 – Graphical representation of real data set : (a) spectral data (b) PCA score plot of the two first components (c) value distribution of Y1 (c) value distribution of Y2

323 Figure 7 shows the graphical representation of real data set. From the
 324 spectra plot (Figure. 7a), it can be seen that there is no visible atypical
 325 spectrum. This means that is not possible to identify or detect outliers in this
 326 data set based on spectra visualisation. Figure 7b shows the PCA score plot
 327 of the two first components. It can be observed that some samples scores are
 328 really different from those of other samples. It is possible that some atypical
 329 samples are outliers but some sample can be also relevant to calculate a
 330 model. From the value distributions plot of the responses (see Figures 7c,d),
 331 it can be seen that some samples show extreme response values in Y1 and

332 Y2. In conclusion, this real data set potentially contains samples that are
333 detrimental to the model.

334 *5.3.2. Method evaluation*

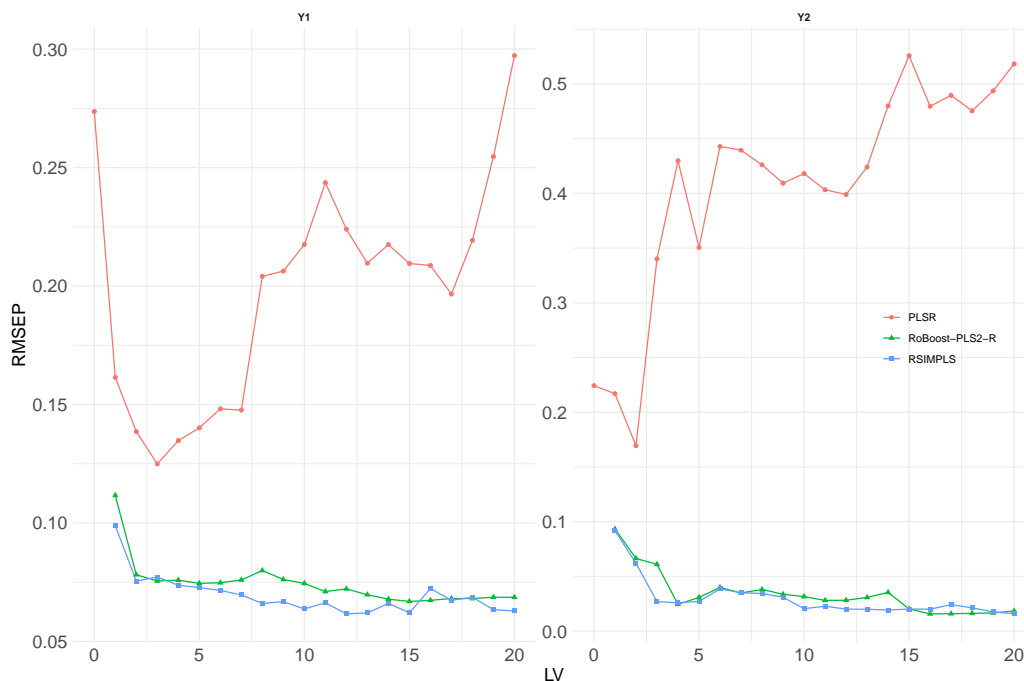


FIGURE 8 – Evolution of the RMSEP as a function of the number of latent variables for the PLS2-R, RSIMPLS and RoBoost-PLS2-R for the real data set

335 Figure 8 represents the prediction performances of the methods on
336 validation set for each reference Y. As there are not all known outliers
337 in the calibration set, it was not possible to define a PLS2-R with and
338 without outliers. Therefore, only the PLS2-R has been calculated on the
339 data with potential outliers. In the figure 8 it can be seen that for both
340 responses the PLSR performance curve is higher than those of the two
341 robust methods. This means that RSIMPLS and RoBoost-PLS2-R method
342 have higher prediction performances than the PLS2-R method applied on
343 this data set. Therefore, some samples are detrimental in the calibration
344 set to the calculation of a PLS2-R model that predicts the samples in the
345 validation set. The two methods RoBoost-PLS2-R and RSIMPLS have close
346 results in terms of RMSEP for a number of latent variables close to 15. This

347 means that both methods were able to deal with potential outliers samples
348 and therefore enable more accurate predictions.

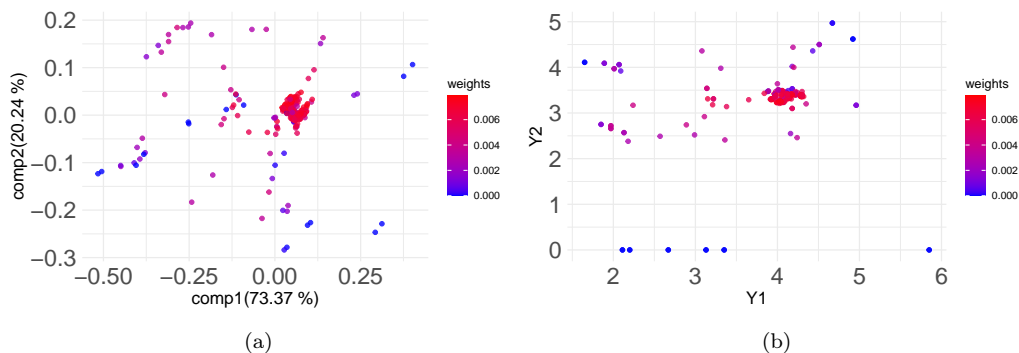


FIGURE 9 – Graphical Representation of the mean weights (for 15 LV) assigned by RoBoost-PLS2-R through PCA score plot of the first two components(a) and Y2 as a function of Y1(b). A colour gradient from blue to red represents the weights assigned to the samples (smallest to largest).

349 Figure 9 shows the weights assigned to the samples by RoBoost-PLS2-R
350 through PCA score plot of the first two components and the Y2 as a function
351 of Y1 plot. It can be seen in figure 9a that not all samples far from the centre
352 were considered as potential outliers (*i.e.* with low weights). Some extreme
353 samples seem to be relevant for the model and were therefore given high
354 weights. The figure 9b shows that some samples have extreme Y-values (0).
355 These samples have a 0 average weight in RoBoost-PLS2-R. This is due to
356 missing value. In this data set, missing data has a value of 0 assigned. It can
357 be concluded through these observations that the RoBoost-PLS2-R method
358 can eliminate outliers on Y but also on X while limiting the assignment of
359 low weights to extreme samples.

360 6. Conclusion

361 In this paper, RoBoost-PLS2-R method is proposed to predict
362 multi-response. This method was evaluated and compared to reference
363 methods on two simulated data sets and one real data set containing
364 different outlier scenarios. For all data sets, prediction performances
365 of RoBoost-PLS2-R are close to those of PLS2-R models calibrated
366 without outliers and to RSIMPLS method. Simulations have shown that
367 RoBoost-PLS2-R extension was very effective when outliers are defined

368 by their spectral properties. In the case of real data, results obtained for
369 both robust methods are better than the PLS2-R method. To conclude,
370 RoBoost-PLS2-R seems to be a reliable and robust regression tool for
371 predicting multi-response variables when data potentially contain outliers.
372 However, some method developments are possible. First of all, the estimation
373 of the criterion evaluated on the Y -residuals can be estimated in another
374 way to take into account the multivariate aspect of Y . In addition, the
375 optimisation of the hyperparameters allowing the weighting of the individuals
376 is complex, it would be relevant to look at automatic parameterisation
377 approaches. Moreover, it could be interesting to use the formalism of the
378 RoBoost-PLS2-R method for cases of categorical variables and thus propose
379 a robust discriminant method. Finally, new RoBoost-PLS2-R algorithm now
380 enables the estimation of regression coefficients contrary to the previous
381 algorithm proposed for RoBoost-PLS1-R. It would be interesting to study
382 these regression coefficients to assess the method's behaviour outside the
383 prediction capacities. In future work, it would be relevant to use the RoBoost
384 formalism for concrete applications involving multi-response variables.

385 It would also be interesting to modify the strategy for visualising the
386 weights of individuals in the calibration. Indeed, here the weights are
387 displayed for each latent variable, so it could be interesting to find a strategy
388 to obtain a weight for each individual allowing to summarise all the weights
389 of each latent variable.

390 **Références**

- 391 [1] S. Wold, M. Sjostrom, L. Eriksson, PLS regression : a basic tool of
392 chemometrics, *Chemometrics and Intelligent Laboratory Systems* 58 (2)
393 (2001) 109–130. [doi:10.1016/S0169-7439\(01\)00155-1](https://doi.org/10.1016/S0169-7439(01)00155-1).
- 394 [2] M. Griep, I. Wakeling, P. Vankeerberghen, D. Massart, Comparison of
395 semirobust and robust partial least squares procedures, *Chemometrics*
396 and *Intelligent Laboratory Systems* 29 (1) (1995) 37–50. [doi:10.1016/0169-7439\(95\)80078-N](https://doi.org/10.1016/0169-7439(95)80078-N).
- 398 [3] I. Stanimirova, S. Serneels, P. J. Van Espen, B. Walczak, How to
399 construct a multiple regression model for data with missing elements
400 and outlying objects, *Analytica Chimica Acta* 581 (2) (2007) 324–332.
401 [doi:10.1016/j.aca.2006.08.014](https://doi.org/10.1016/j.aca.2006.08.014).

- 402 [4] R. J. Pell, Multiple outlier detection for multivariate calibration using
403 robust statistical techniques, *Chemometrics and Intelligent Laboratory*
404 *Systems* 52 (1) (2000) 87–104. doi:10.1016/S0169-7439(00)00082-4.
- 405 [5] J. A. Gil, R. Romera, On robust partial least squares (PLS) methods,
406 *Journal of Chemometrics* 12 (6) (1998) 365–378. doi:10.1002/(SICI)
407 1099-128X(199811/12)12:6<365::AID-CEM519>3.0.CO;2-G.
- 408 [6] J. González, D. Peña, R. Romera, A robust partial least squares
409 regression method with applications, *Journal of Chemometrics* 23 (2)
410 (2009) 78–90. doi:10.1002/cem.1195.
- 411 [7] I. N. Wakelinc, H. J. H. Macfie, A robust PLS procedure, *Journal of*
412 *Chemometrics* 6 (4) (1992) 189–198. doi:10.1002/cem.1180060404.
- 413 [8] J. Peng, S. Peng, Y. Hu, Partial least squares and random sample
414 consensus in outlier detection, *Analytica Chimica Acta* 719 (2012) 24–29.
415 doi:10.1016/j.aca.2011.12.058.
- 416 [9] P. Filzmoser, R. Maronna, M. Werner, Outlier identification in high
417 dimensions, *Computational Statistics & Data Analysis* 52 (3) (2008)
418 1694–1711. doi:10.1016/j.csda.2007.05.018.
- 419 [10] M. Hubert, K. V. Branden, Robust methods for partial least squares
420 regression, *Journal of Chemometrics* 17 (10) (2003) 537–549. doi:10.
421 1002/cem.822.
- 422 [11] U. Kruger, Y. Zhou, X. Wang, D. Rooney, J. Thompson, Robust
423 partial least squares regression : Part II, new algorithm and benchmark
424 studies, *Journal of Chemometrics* 22 (1) (2008) 14–22, _eprint :
425 https://onlinelibrary.wiley.com/doi/pdf/10.1002/cem.1095. doi:10.
426 1002/cem.1095.
- 427 [12] I. Hoffmann, S. Serneels, P. Filzmoser, C. Croux, Sparse partial robust
428 M regression, *Chemometrics and Intelligent Laboratory Systems* 149
429 (2015) 50–59. doi:10.1016/j.chemolab.2015.09.019.
- 430 [13] P. Filzmoser, S. Serneels, R. Maronna, C. Croux, Robust multivariate
431 methods in *Chemometrics*, arXiv :2006.01617 [stat] (2020)
432 393–430ArXiv : 2006.01617. doi:10.1016/B978-0-12-409547-2.
433 14642-6.

- 434 [14] M. Hubert, K. V. Branden, [Robust methods for partial least](#)
435 [squares regression](#), *Journal of Chemometrics* 17 (10) (2003) 537–549.
436 [doi:10.1002/cem.822](#).
437 URL [https://onlinelibrary.wiley.com/doi/abs/10.1002/cem.](https://onlinelibrary.wiley.com/doi/abs/10.1002/cem.822)
438 [822](#)
- 439 [15] M. Metz, F. Abdelghafour, J.-M. Roger, M. Lesnoff, [A novel](#)
440 [robust PLS regression method inspired from boosting principles:](#)
441 [RoBoost-PLSR](#), *Analytica Chimica Acta* (2021) 338823 [doi:](#)
442 [10.1016/j.aca.2021.338823](#).
443 URL [https://linkinghub.elsevier.com/retrieve/pii/](https://linkinghub.elsevier.com/retrieve/pii/S0003267021006498)
444 [S0003267021006498](#)
- 445 [16] M. Metz, A. Biancolillo, M. Lesnoff, J.-M. Roger, [A note on spectral](#)
446 [data simulation](#), *Chemometrics and Intelligent Laboratory Systems* 200
447 (2020) 103979. [doi:10.1016/j.chemolab.2020.103979](#).