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

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

^aITAP Univ Montpellier INRAE Institut Agro Montpellier France ^bChemHouse Research Group Montpellier France ^cSELMET Univ Montpellier CIRAD INRAE Montpellier SupAgro Montpellier France ^dWallon Agricultural Research Centre Gembloux

5 Abstract

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

18 Keywords: Robust regression methods, outliers, multi-response,

19 multivariate data analysis

20 1. Introduction

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

Preprint submitted to Chemometrics and Intelligent Laboratory Systems 13 janvier 2022

exceeds the number of samples. This method is widely used in analytical 27 chemistry for predicting constituent concentrations of a sample based on its 28 spectrum obtained by spectroscopic techniques, such as near-infrared (NIR) 29 spectroscopy, Fluorescence spectroscopy and ultraviolet (UV) spectroscopy. 30 The PLSR model is known to be affected by the presence of atypical 31 observations (outliers) in the data set. Outliers can negatively affect the 32 calibration of PLSR models. To deal with outliers, several robust PLSR 33 methods were proposed in the literature [2-12]. These methods were 34 particularly developed to deal with outliers when the response matrix is 35 uni-dimensional [13] (PLS1-R). However, robust methods that address the 36 case of multi-responses (PLS2) are few. Among them, RSIMPLS is one 37 of the most used method [14]. RSIMPLS proposes to robustly estimate 38 the cross-covariance matrix C_{xy} and the empirical covariance matrix C_x 39 used in SIMPLS algorithm. For this, a robust principal component analysis 40 (ROBPCA) is performed on the concatenated data matrix of \mathbf{X} and \mathbf{Y} . 41 RSIMPLS uses additional information from the previous ROBPCA step to 42 perform a reweighted multiple linear regression. 43

Recently, a new robust method called RoBoost-PLSR has been developed 44 [15]. RoBoost-PLSR aims at determining the measure of relevance of 45 the samples for PLSR model calibration. Indeed, in practical cases, the 46 samples of a database are not defined as outliers, i.e. not relevant for the 47 calibration of a PLSR model. RoBoost-PLSR proposes to calculate a weight 48 on each latent variable to define the relevance of the samples. The relevance 49 measurement is defined according to three criteria calculated for each latent 50 variable (X-residuals, Y-residuals, leverage). This method has proven to be 51 effective for outliers in both Y and X. However, this algorithm was only 52 developed for a one-dimensional PLSR response variable (PLS1). This paper 53 contributes to the RoBoost-PLSR method which will be able to manage 54 outliers in a multiple response context. 55

The first section introduces the extension of RoBoost-PLSR named RoBoost-PLS2-R and the associated algorithm. The following section presents the data and the methods used to evaluate and compare the predictive ability of RoBoost-PLS2-R. Finally, the prediction performance of RoBoost-PLS2-R and its comparison with standard methods are shown in the last section.

62 2. Notations

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

73 3. RoBoost-PLSR extension for multi-responses

74 3.1. Algorithm

The new algorithm allowing an extension in a multi-response context is 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

$$k = 1$$

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

2: Center the data :

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

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

- 3: Define \mathbf{u}_k as an arbitrary column of \mathbf{Y}
- 4: Calculate one weighted latent variable NIPALS :

$$\mathbf{w}_{k} = \frac{\mathbf{X}_{k}^{\mathrm{T}} \mathbf{D} \mathbf{u}_{k}}{||\mathbf{X}_{k}^{\mathrm{T}} \mathbf{D} \mathbf{u}_{k}||}$$
$$\mathbf{t}_{k} = \mathbf{X}_{k} \mathbf{w}_{k}$$
$$\mathbf{p}_{k} = \frac{\mathbf{X}_{k}^{\mathrm{T}} \mathbf{D} \mathbf{t}_{k}}{\mathbf{t}_{k}^{\mathrm{T}} \mathbf{D} \mathbf{t}_{k}}$$
$$\mathbf{q}_{k} = \frac{\mathbf{Y}_{k}^{\mathrm{T}} \mathbf{D} \mathbf{t}_{k}}{\mathbf{t}_{k}^{\mathrm{T}} \mathbf{D} \mathbf{t}_{k}}$$
$$c_{k} = \frac{\mathbf{u}_{k}^{\mathrm{T}} \mathbf{D} \mathbf{t}_{k}}{\mathbf{t}_{k}^{\mathrm{T}} \mathbf{D} \mathbf{t}_{k}}$$

5: Derive (F), (E), (l):

$$\mathbf{E} = \mathbf{X}_k - \mathbf{t}_k \mathbf{p}_k^{ ext{T}}$$
 $\mathbf{F} = \mathbf{Y}_k - \mathbf{t}_k \mathbf{q}_k^{ ext{T}}$

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

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

$$egin{aligned} \mathbf{X}_{k+1} &= \mathbf{X}_k - \mathbf{t}_k \mathbf{p}_k^{ ext{T}} \ \mathbf{Y}_{k+1} &= \mathbf{Y}_k - \mathbf{t}_k \mathbf{q}_k^{ ext{T}} \ \mathbf{u}_{k+1} &= \mathbf{Y}_k \mathbf{q}_k \end{aligned}$$

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

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

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

With \mathbf{R} :

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

77 3.2. Theoretical discussions

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

80 81

82

83

84

85

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

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

- 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.
- The algorithm proposed in this article provides regression coefficients. 107 This makes the constructed RoBoost-PLSR models more easily 108 interpretable. Contrary to the first algorithm proposed in [15], the 109 rotation matrix \mathbf{R} used to estimate the regression coefficients can be 110 estimated. This is due to data centring which is only done for the 111 first model with a single latent variable. In the previous algorithm, 112 repeated centring of X and Y matrices led to a bias which made it 113 impossible to estimate the rotation matrix. 114
- 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

115

100

106

— The algorithm proposed in [15], determines the convergence with q.

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

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

151

124

125

126

127

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

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

with z_i :

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

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

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

¹⁶⁵ 4. Materials and methods

166 4.1. Simulated Data

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

$$\mathbf{X} = \mathbf{t}_{\mathbf{u}} \mathbf{p}_{\mathbf{u}}^{\ t} + \mathbf{T}_{\mathbf{d}} \mathbf{P}_{\mathbf{d}}^{\ t} + \mathbf{E}$$
(1)

And the relationship f between **X** and **Y** by :

$$\mathbf{Y} = f(\mathbf{t}_{\mathbf{u}}) + \mathbf{F} \tag{2}$$

Where $\mathbf{p}_{\mathbf{u}}$ is the spectral signature in the useful space and $\mathbf{P}_{\mathbf{d}}$ are spectral signatures in the detrimental space. $\mathbf{t}_{\mathbf{u}}$ and $\mathbf{T}_{\mathbf{d}}$ are their associated contributions. The **E** and **F** matrices are defined as gaussian noises of **X** and **Y**, respectively.

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

184 4.1.1. Simulation 1, Y-outliers

The Y-outliers were defined by their relationship f between **X** and **Y**. All other simulation parameters were common between inliers and outliers. The construction of the simulated data set 1 is represented in table 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	$Y_1 = 10 * T_{ethanol}$	$Y_1 = 10 * T_{ethanol}$	
	$Y_2 = 10 * T_{glucose}$	$\mathbf{Y_2} = -10*\mathbf{T_{glucose}}$	
	$Y_3 = 10 * T_{water}$	$Y_3 = 10 * T_{water}$	
\mathbf{F}	Gaussian distribution		

TABLE 1 – The different choices in the simulation 1

188 4.1.2. Simulation 2, X-outliers

The X-outliers were defined by others artificial spectral signatures. These signatures correspond to minority compounds. All other simulation parameters were common between inliers and outliers. The simulation is represented in table 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 water
	Pure spectrum of ethanol	Pure spectrum of ethanol
	Spectrum of water-ethanol Interaction	Spectrum of water-ethanol Interaction
	10 Artificial spectra	10 Artificial spectra
		10 Artificial spectra
\mathbf{T}_d	Folded-normal distribution	Folded-normal distribution
	Folded-normal distribution	Folded-normal distribution
	Product between T_{water} and $T_{ethanol}$	Product between T_{water} and $T_{ethanol}$
	Folded-normal distribution	Folded-normal distribution
		Folded-normal distribution
Е	Gaussian distribution	
f	$Y_1 = 10 * T_{ethanol}$	$Y_1 = 10 * T_{ethanol}$
	$Y_2 = 10 * T_{glucose}$	$Y_2 = 10 * T_{glucose}$
	$Y_3 = 10 * T_{water}$	$Y_3 = 10 * T_{water}$
F	Gaussian distribution	

TABLE 2 – The different choices in the simulation 2

193 4.2. Real data set

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

204 4.2.1. Evaluation strategies

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

In the case of the simulations, the 1000 samples were divided into two groups : 800 for calibration and 200 for validation. The reference method in terms of prediction performance was PLS2-R calibrated without outliers. For
the real data set, calibration set was composed of 209 samples. The validation
was conducted on 52 samples. These samples were selected from a study of the
data in order to represent the samples as well as possible without containing
potential outliers. The reference method in terms of prediction performance
was RSIMPLS.

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

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

225 4.3. Software

PLS2-R was performed with "rnirs" and RoBoost-PLS2-R is available RoBoost-PLSR functions available in R. RSIMPLS was performed using the 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.

Figure 1 shows the graphical representation of simulation 1. From the spectra plot (Figure 1a), it can be seen that is difficult to identify outliers (in red) from a simple visual inspection. In this case, the outliers were defined by a distinct relation f on one of the response variables (see Table 1). Therefore, no spectral difference between the two groups is expected.

From the plot of value distributions of the response variables (see Figure 1b,c,d) it can be observed that Y1 and Y3 variables present the same distribution for both outliers and inliers. However, different distribution for

these two groups is presented in Y2 variable. Moreover, the variances of Y1are 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

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



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.

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

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

278 5.2. Simulation 2





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.

Figure 4 shows the graphical representation of simulation 2. From spectra plot of the sample (Figure 4 a), it can be seen that outliers are not identifiable. Indeed, in this simulation, outliers are different only for spectral signatures and hence, they contribute slightly to the construction of the spectra. Figures the score plot on the two first principal components. Two centroids can be seen but there is no clear separation between outliers and inliers. This is due to the outliers having their major compounds in common (see Table 2). From the value distributions plot of the responses (see : Figures **4**c,d,e), it can be seen that outliers and inliers present similar distribution in all Y response variables. Outliers are different only on the basis of the 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

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

of the leverage effect is not directly taken into account. Indeed, in RSIMPLS it is the cross-covariance matrices C_{xy} and the empirical covariance matrix C_x that are robustly estimated.

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

Figure 6 shows the weight assigned to samples by RoBoost-PLS2-R according to the number of LV. It can be observed that the weights of outliers decrease progressively when the number of LV increases. This gradual decrease is partly explained by the fact that both outliers and inliers were simulated using common majority spectral signatures. Indeed, only some minor spectral signatures differentiate the inliers from the outliers (see Section 4). After 8 latent variables, all outliers have a weight equal to 0, whereas almost all inliers present a high weight. Nevertheless, it is possible to note that the majority of the inliers have a strong weight and therefore a 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

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

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





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

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

means that both methods were able to deal with potential outliers samplesand 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).

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

360 6. Conclusion

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

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

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

390 Références

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

- [4] R. J. Pell, Multiple outlier detection for multivariate calibration using
 robust statistical techniques, Chemometrics and Intelligent Laboratory
 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.
- [6] J. González, D. Peña, R. Romera, A robust partial least squares
 regression method with applications, Journal of Chemometrics 23 (2)
 (2009) 78–90. doi:10.1002/cem.1195.
- [7] I. N. Wakelinc, H. J. H. Macfie, A robust PLS procedure, Journal of
 Chemometrics 6 (4) (1992) 189–198. doi:10.1002/cem.1180060404.
- [8] J. Peng, S. Peng, Y. Hu, Partial least squares and random sample
 consensus in outlier detection, Analytica Chimica Acta 719 (2012) 24–29.
 doi:10.1016/j.aca.2011.12.058.
- [9] P. Filzmoser, R. Maronna, M. Werner, Outlier identification in high
 dimensions, Computational Statistics & Data Analysis 52 (3) (2008)
 1694–1711. doi:10.1016/j.csda.2007.05.018.
- [10] M. Hubert, K. V. Branden, Robust methods for partial least squares
 regression, Journal of Chemometrics 17 (10) (2003) 537-549. doi:10.
 1002/cem.822.
- [11] U. Kruger, Y. Zhou, X. Wang, D. Rooney, J. Thompson, Robust partial least squares regression : Part II, new algorithm and benchmark studies, Journal of Chemometrics 22 (1) (2008) 14–22, __eprint : https ://onlinelibrary.wiley.com/doi/pdf/10.1002/cem.1095. doi:10.
 1002/cem.1095.
- I. Hoffmann, S. Serneels, P. Filzmoser, C. Croux, Sparse partial robust
 M regression, Chemometrics and Intelligent Laboratory Systems 149
 (2015) 50–59. doi:10.1016/j.chemolab.2015.09.019.
- [13] P. Filzmoser, S. Serneels, R. Maronna, C. Croux, Robust multivariate
 methods in Chemometrics, arXiv :2006.01617 [stat] (2020)
 393-430ArXiv : 2006.01617. doi:10.1016/B978-0-12-409547-2.
 14642-6.

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