



HAL
open science

Further investigations on the relationship between the OPLS preprocessing and the NAS

Jean-Claude Boulet, Robert Sabatier

► **To cite this version:**

Jean-Claude Boulet, Robert Sabatier. Further investigations on the relationship between the OPLS preprocessing and the NAS. *Chemometrics and Intelligent Laboratory Systems*, 2020, 206, 10.1016/j.chemolab.2020.104159 . hal-03108878

HAL Id: hal-03108878

<https://hal.inrae.fr/hal-03108878>

Submitted on 26 Sep 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 4.0 International License

Further investigations on the relationship between the O-PLS preprocessing and the NAS.

Jean-Claude Boulet ^{*†}
Robert Sabatier [‡]

*INRAE, UMR1083 Sciences For Enology, SupAgro, University of Montpellier, F-34060
Montpellier, France

†Chemhouse Research Group, Montpellier, France

‡University of Montpellier, Institut of Functionnal Genomic, CNRS, UMR5203, Inserm
U1191, F-34060 Montpellier, France

Abstract

O-PLS is a preprocessing that was presented as an improvement of the PLS algorithm it was issued from. Nevertheless, the bibliography did not confirm, neither for prediction, nor for interpretation. To contribute to a better understanding, we investigated the relationship between O-PLS and the NAS net analyte signal. For four numerical applications, the matrix obtained after the O-PLS deflation tended towards a matrix of rank 1 when the number of removed dimensions increases. Therefore, the line-vectors of this matrix are colinear to the NAS, and so the usual one-latent-variable PLS1 regression after the O-PLS preprocessing can be replaced by almost any regression taking into account from one to all of the variables.

Keywords: PLS, O-PLS, OSC, NAS, prediction, interpretation

1 Introduction

2 Chemometrics is largely based on the Partial Least Squares or Projection
3 to Latent Structures (PLS) method. Hereafter, all the discussion will ap-
4 ply to the prediction of a single variable. Let \mathbf{X} and \mathbf{y} be a matrix and a
5 column-vector containing the values of variables measured on the same ob-
6 servations. PLS regression was developed to predict the values of \mathbf{y} using \mathbf{X} ,
7 specially when the variables in \mathbf{X} are highly correlated. For example, PLS
8 regression allowed the prediction of gluten in a wheat sample from a simple
9 near infrared spectrum. Many PLS algorithms have been proposed. Nine
10 of them were compared in Andersson [1]. We proposed a tenth algorithm
11 more focused on the geometry of PLS [2]. Countless applications of PLS can
12 be found in industrial and scientific applications. Nevertheless, in certain
13 fields such as metabolomics, another method is being replacing PLS. This
14 method is based on a preprocessing called Orthogonal Projection to Latent
15 Structures (O-PLS). The O-PLS algorithm was obtained by a slightly mod-
16 ification of the PLS algorithm. To understand it, let us come back to the
17 origin.

18 The principle of PLS is to extract the subspace in \mathbf{X} that contains the
19 information related to the explanation of \mathbf{y} . The idea of the Orthogonal
20 Signal Correction (OSC) preprocessings is to help PLS by removing from \mathbf{X}
21 all information not related to \mathbf{y} . OSC are followed by a regression, usually
22 a PLS. A first OSC algorithm, derived from the NIPALS-PLS algorithm,
23 was proposed by Wold et al [3]. This algorithm was latter improved by
24 Trygg and Wold, yielding to a new method called O-PLS [4, 5, 6]. But
25 while other OSC algorithms were proposed [7, 8], the ability of OSC to
26 improve the predictive ability of a regular PLS was being discussed [9, 10].
27 More attention was paid to O-PLS. It was observed that one O-PLS factor

28 removed spared one PLS latent variable, without gain in prediction errors
29 [7, 11, 9]. Then Verron et al [12] mathematically confirmed this observation,
30 and Kemsley & Tapp [13] showed a direct relationship between O-PLS and
31 PLS scores. As a conclusion, O-PLS followed by a PLS regression yields
32 exactly the same prediction than a single PLS regression. Then O-PLS
33 benefits were focused on interpretation rather than prediction, e.g. [14],
34 but even the added value of interpretation was challenged [15], conclusions
35 we agree with. Nevertheless, there is a gap between O-PLS users on the
36 one hand, and mathematical considerations on O-PLS on the other hand.
37 Maybe a geometrical approach could contribute to a better understanding
38 of O-PLS? In particular, Goicoechea et al [10] and Ni et al [16] studied
39 the relationship between OSC methods and the Net Analyte Signal (NAS)
40 introduced by Lorber [17]. We propose to further investigate this property.
41 Notations are presented in the figure 1.

42 **2 Theory.**

43 The name O-PLS let think to a PLS-like method, i.e. a regression method.
44 Certain so-called "O-PLS" functions associate an O-PLS preprocessing fol-
45 lowed by a PLS regression with one latent variable. These names are mis-
46 leading, because O-PLS as introduced in the original papers of Trygg and
47 Wold is clearly a preprocessing, not a calibration method. In order to obtain
48 predictions, O-PLS needs to be followed by a calibration step, e.g. a classical
49 least squares or a PLS regression [10] usually processed with a single latent
50 variable.

51 O-PLS calculation is iterative. Let A_{max} be a predetermined large num-
52 ber of dimensions to be removed. Scores \mathbf{t}_i and loadings \mathbf{p}_i are computed
53 for each value of i between 1 and A_{max} , then the correction is performed

54 by deflation[9]: $\mathbf{X}_{\text{opls}} = \mathbf{X} - \sum_{i=1}^{A_{\text{max}}} \mathbf{t}_i \mathbf{p}_i'$. Each dimension i is supposed
55 to contain information unrelated to \mathbf{y} . But while i is increasing, there will
56 be a breakpoint for which no more information unrelated to \mathbf{y} will be left
57 in \mathbf{X}_{opls} . Let A be this value of i ; $A \leq A_{\text{max}}$. A definition of the NAS
58 is *the part of the signal orthogonal to the other constituents* [17]. This
59 definition was made for spectra, but it can be extended to other signals.
60 Let \mathbf{s}_{nas} be the NAS associated to the compound of interest whose val-
61 ues form \mathbf{y} . An interesting property of the NAS is that the product of
62 \mathbf{X}_{opls} by \mathbf{s}_{nas} should be proportional to \mathbf{y} . It derives from the following
63 equation [10] : $\mathbf{X}_{\text{opls}} = k \mathbf{y} \mathbf{s}_{\text{nas}}'$ with k a constant. This situation should oc-
64 cur when an appropriate number of dimensions A has been removed by the
65 OSC. Therefore, the corresponding matrix \mathbf{X}_{opls} should be of rank 1, all the
66 line-vectors of \mathbf{X}_{opls} being collinear. Under these conditions, the classical
67 one-latent-variable PLS1 regression which follows the O-PLS preprocessing
68 could be replaced by *any* other regression without significant loss of pre-
69 dictive ability, even with a simple linear regression. On the other hand, if
70 the theory is not verified, if \mathbf{X}_{opls} is not close enough from a rank-1 matrix,
71 then the predictions should strongly depend on the choice of the regression
72 method. To assess this property, three methods based on O-PLS were com-
73 pared. The first method noted *OPLS-classic* is the classical combination of
74 an O-PLS preprocessing followed by a PLS regression with 1 latent variable.
75 The second method noted *OPLS-ones* is a modified version of the previous
76 *OPLS-classic*. The modification consisted in replacing the weight vector \mathbf{w}
77 which begins the PLSR algorithm by an arbitrary vector of ones of same
78 dimension. *OPLS-classic* and *OPLS-ones* algorithms are summarized in Ta-
79 ble 2. The third method noted *OPLS-univ* is a classical OPLS followed by
80 an univariate regression. The selected variable was chosen as the one with

81 the largest variability in \mathbf{X}_{opls} after the O-PLS preprocessing removed A_{max}
82 components.

83 **3 Material and methods**

84 Four datasets were selected. Two datasets were challenges proposed by
85 Pierre Dardenne and Vincent Baeten, CRA-Wallonie, at the French annual
86 chemometrics conference in 2007 and 2018. The goal of the 2007 challenge
87 [18] was to predict gluten concentrations. The test dataset was picked up,
88 since it contained a large number of observations (2000). The 2018 chal-
89 lenge consisted in 3908 calibration plus 429 test samples of near infrared
90 spectra. The compound to quantify was not identified. A third dataset
91 was provided by Sylvie Bureau, INRAE, UMR408. Mid infrared spectra
92 ($4000 - 650\text{cm}^{-1}$) were acquired on 750 apricots, for which several reference
93 analysis were performed including refractive index. The last dataset was
94 provided by INRAE, UMR ITAP and SPO. Visible-near infrared spectra
95 were acquired in transmittance and Brix degrees were also measured on 250
96 grape berries. Processing needed a calibration and a test dataset. Thus,
97 the 2007 challenge, the apricots and grape berries datasets were splitted,
98 the 1500, 600 and 200 observations were assigned to the calibration dataset,
99 the last 500, 150 and 50 observations were assigned to the test dataset, re-
100 spectively. After centering, the three regression methods described above,
101 *OPLS-classic*, *OPLS-ones* and *OPLS-univ* were processed on the calibra-
102 tion datasets with 1 to A_{max} dimensions removed by OPLS, A_{max} being set
103 to 50 for the apricot datasets and to 40 for the three other datasets. The
104 variable selected for the *OPLS-univ* model was the one with the largest vari-
105 ability after removal of A_{max} dimensions by OPLS. The predictive abilities
106 of the three models were assessed by the calibration errors (RMSEC), the

107 cross-validation errors (RMSECV) and the prediction errors (RMSEP). The
108 RMSECVs were the average of a random 2-blocks cross-validation repeated
109 50 times.

110 4 Results

111 Results are presented as figures 3, 4, 5 and 6, one for each of the four
112 datasets. After removal of the A_{max} OPLS components, the residual spec-
113 tra in subfigures (a) present a similar shape, sometimes symmetrical along
114 the axis $y = 0$ which is consistent with a low rank of the corresponding
115 matrix. These subfigures were used to choose the univariate variable of the
116 *OPLS-univ* regression. RMSECs in subfigures (b) are similar for the four
117 datasets. When i is low, a few differences between the three models can be
118 observed. But when i increases, the curves overlay. RMSECV in subfigures
119 (c) also present the same trend for all datasets, but with a different behav-
120 ior according to the models. The RMSECV of *OPLS-univ* decreases steadily
121 when i increases. On the other hand, *OPLS-classic* and *OPLS-ones* present
122 the classical shape of a decrease, then an increase. When i is high, those
123 two curves overlay. Then, RMSEP in subfigures (d) can present differences
124 between models. But when i is high, the curves of the three models overlay.

125 5 Discussion and conclusion

126 The hypothesis that, by increasing the number i of removed dimensions, the
127 rank of the resulting matrix \mathbf{X}_{opls} gets close to 1, still stands up after the
128 benchmark on 4 datasets. Similar results, not presented, confirmed this ob-
129 servation. When i increases, all the line vectors of \mathbf{X}_{opls} tend to be colinear.
130 This situation occurred when i was higher than $A = 7, 15, 40$ and 5 for the

131 2007 and 2018 challenges, the apricots and grape berries datasets respec-
 132 tively. These values of A need to be compared to A_{opt} the optimal number
 133 of dimensions to be removed, determined by the RMSECV curves, around
 134 13–20, 15–20, 7–12 and 15–22 respectively. For the apricot dataset, A_{opt}
 135 corresponds to a situation where the three models *OPLS-classic*, *OPLS-ones*
 136 and *OPLS-univ* do not yield the same prediction, and where prediction er-
 137 rors remain high. Figure 5(d) suggests that a higher number of dimensions,
 138 A_{opt} around 20–22, would have yielded better predictions, but such values
 139 remain before the convergence region obtained for $A = 40$ ($A > A_{opt}$). For
 140 the three other datasets, convergence of the three models had already been
 141 achieved: $A < A_{opt}$, and the models were found more robust when applied
 142 to the test datasets. Maybe a clue to identify robust ($A < A_{opt}$) from not
 143 robust ($A > A_{opt}$) models?

144 Finally, all the information about the NAS is supported by the \mathbf{X}_{opls} ma-
 145 trix, issued from the O-PLS preprocessing. The weights \mathbf{w} and the regression
 146 vector \mathbf{b} of the PLSR in the *OPLS-classic* method are linear combination
 147 of the lines of \mathbf{X}_{opls} , weighed by the \mathbf{y} values. Setting \mathbf{w} to a vector of
 148 ones in *OPLS-ones* is another way to weight the lines of \mathbf{X}_{opls} , each line
 149 has the same weight. When \mathbf{X}_{opls} is close to a rank-1 matrix, these calcula-
 150 tions yield the same result. The \mathbf{X}_{osc} matrix or the regression coefficients \mathbf{b}
 151 can all be interpreted as an estimation of the NAS. According to its defini-
 152 tion, the NAS is obtained after an orthogonal (or oblique) projection, and
 153 therefore inherits of the properties of such calculations. In particular, the
 154 NAS is expected to present features from the compound of interest (whose
 155 concentrations are \mathbf{y}), but more annoying it can present features from the
 156 other coexisting constituents [19], and eventually it can have dropped fea-
 157 tures that were present in the compound of interest. So the interpretation

158 of the NAS remains tricky and should be confirmed by other informations.
159 These considerations also apply to the residual matrix, $\mathbf{X} - \mathbf{X}_{opls}$, which is
160 the result of an orthogonal projection, too. Note that the same interpre-
161 tation of the NAS calculated from O-PLS can also be performed with the
162 regression coefficients of a PLS regression, already reported as an estimation
163 of the NAS [20].

164 **6 Acknowledgements**

165 Thanks to Vincent Baeten, Pierre Dardenne and Cécile Barron, for providing
166 datasets; and Christelle Reynès and JM Roger for useful comments on the
167 manuscript.

168 **References**

- 169 [1] Martin Andersson. A comparison of nine pls1 algorithms. *Journal of*
170 *Chemometrics*, 23:518–529, 2009.
- 171 [2] Jean-Claude Boulet, Dominique Bertrand, Gérard Mazerolles, Robert
172 Sabatier, and Jean-Michel Roger. A family of regression methods de-
173 rived from standard plsr. *Chemometrics and Intelligent Laboratory Sys-*
174 *tems*, 120:13–25, 2013.
- 175 [3] Svante Wold, Henrik Antti, Frederik Lindgren, and Jerker Ohman. Or-
176 thogonal signal correction of near infrared spectra. *Chemometrics and*
177 *Intelligent Laboratory Systems*, 44:175–185, 1998.
- 178 [4] J.Trygg. *Parsimonious multivariate models*. PhD thesis, Umea Univer-
179 sity, Sweden, 2001.

- 180 [5] Yohan Trygg and Svante Wold. Orthogonal projections to latent struc-
181 tures (o-pls). *Journal of Chemometrics*, 16:119–128, 2002.
- 182 [6] Johan Trygg and Svante Wold. Orthogonal signal projection. Patent
183 US20030200040A1, 2003.
- 184 [7] Tom Fearn. On orthogonal signal correction. *Chemometrics and Intel-*
185 *ligent Laboratory Systems*, 50:47–52, 2000.
- 186 [8] Johan A. Westerhuis, Sijmen De Jong, and Age K. Smilde. Direct
187 orthogonal signal correction. *Chemometrics and Intelligent Laboratory*
188 *Systems*, 23:13–25, 2001.
- 189 [9] O. Svensson, T. Kourti, and J.F. MacGregor. An investigation of or-
190 thogonal signal correction algorithms and their characteristics. *Journal*
191 *of Chemometrics*, 16:176–188, 2002.
- 192 [10] Hector C. Goicoechea and Alejandro C. Olivieri. A comparison of or-
193 thogonal signal correction and net analyte preprocessing methods. the-
194 oretical and experimental study. *Chemometrics and Intelligent Labora-*
195 *tory Systems*, 56:73–81, 2001.
- 196 [11] Agnar Hoskuldsson. Variable and subset selection in pls regression.
197 *Chemometrics and Intelligent Laboratory Systems*, 55:23–38, 2001.
- 198 [12] Thomas Verron, Robert Sabatier, and Richard Joffre. Some theoretical
199 properties of the o-pls method. *Journal of Chemometrics*, 18:62–68,
200 2004.
- 201 [13] E.K. Kemsley and H.S. Tapp. Opls filtered data can be obtained directly
202 from non-orthogonalized pls1. *Journal of Chemometrics*, 23:263–264,
203 2009.

- 204 [14] Matthias Hedenstrom, Susanne Wiklund, Bjorn Sundberg, and Ulf Ed-
205 lund. Visualization and interpretation of o-pls models based on 2d nmr
206 data. *Chemometrics and Intelligent Laboratory Systems*, 92:110–117,
207 2008.
- 208 [15] Ulf G. Indahl. The o-pls methodology for orthogonal signal correction
209 - is it correcting or confusing. *Journal of Chemometrics*, pages 1–14,
210 2017.
- 211 [16] Wangdong Ni, Steven.D. Brown, and Ruilin Man. The relationship
212 between net analyte signal/preprocessing and orthogonal signal cor-
213 rection algorithms. *Chemometrics and Intelligent Laboratory Systems*,
214 98:97–107, 2009.
- 215 [17] Avraham Lorber, Klass Faber, and Bruce R. Kowalski. Net analyte
216 signal calculation in multivariate calibration. *Analytical Chemistry*,
217 69:1620–1626, 1997.
- 218 [18] J.A. Fernandez Pierna, F. Chauchard, S. Preys, J.M. Roger, O. Galtier,
219 V. Baeten, and P. Dardenne. How to build a robust model against per-
220 turbation factors with only a few reference values: a chemometric chal-
221 lenge at chimie 2007. *Chemometrics and Intelligent Laboratory*
222 *Systems*, 106:152–159, 2011.
- 223 [19] J.C. Boulet and J.M. Roger. Pretreatments by means of orthogonal
224 projections. *Chemometrics and Intelligent Laboratory Systems*, 117:61–
225 69, 2012.
- 226 [20] Joan Ferré and Nicolaas M. Faber. Net analyte signal calculation for
227 multivariate calibration. *Chemometrics and Intelligent Laboratory Sys-*
228 *tems*, 69:123–136, 2003.

229 7 Captions

230 Figure 1: Notations

231 Figure 2: *OPLS-classic* a regular O-PLS followed by a PLS1 with one latent
232 variable (left); *OPLS-ones* a regular O-PLS followed by a w-modified PLS1,
233 \mathbf{w} being set to a vector of ones (right)

234 Figure 3: Chimie 2007 challenge, NIR spectra: (a) spectra after 40
235 OPLS components removed, the vertical line representing the variable se-
236 lected for *OPLS-univ*; (b-d) RMSEC, RMSECV and RMSEP for the 3 mod-
237 els: *OPLS-classic* (red), *OPLS-ones* (green) and *OPLS-univ* (blue)

238 Figure 4: Chimie 2018 challenge, NIR spectra: (a) spectra after 40
239 OPLS components removed, the vertical line representing the variable se-
240 lected for *OPLS-univ*; (b-d) RMSEC, RMSECV and RMSEP for the 3 mod-
241 els: *OPLS-classic* (red), *OPLS-ones* (green) and *OPLS-univ* (blue)

242 Figure 5: Apricots MIR spectra: (a) spectra after 50 OPLS components
243 removed, the vertical line representing the variable selected for *OPLS-univ*;
244 (b-d) RMSEC, RMSECV and RMSEP for the 3 models: *OPLS-classic* (red),
245 *OPLS-ones* (green) and *OPLS-univ* (blue)

246 Figure 6: Grape fruit NIR reflectances: (a) spectra after 40 OPLS compo-
247 nents removed, the vertical line representing the variable selected for *OPLS-*
248 *univ*; (b-d) RMSEC, RMSECV and RMSEP for the 3 models: *OPLS-classic*

249 (red), *OPLS-ones* (green) and *OPLS-univ* (blue)

250

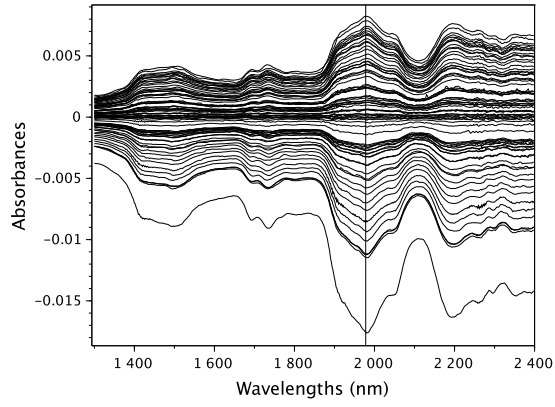
251 **8 Figures**

N	number of observations
Q	number of variables
i	number of O-PLS components removed
A	minimum value of i for which \mathbf{X}_{osc} can be approximated to a rank-1matrix
A_{opt}	value or range of values of i for which the RMSECV is minimum
A_{max}	a value of i larger than A and A_{opt}
\mathbf{X}	matrix of N lines and Q columns
\mathbf{y}	column vector of N elements
\mathbf{X}_{osc}	the matrix \mathbf{X} after OSC correction
\mathbf{X}_{opls}	the matrix \mathbf{X} after O-PLS correction
\mathbf{t}_i	i^{th} score vector for OSC/OPLS
\mathbf{p}_i	i^{th} loading vector for OSC/OPLS
\mathbf{w}	weight vector of PLS
\mathbf{t}	score vector of PLS
\mathbf{b}	regression vector of PLS
$\mathbf{1}_Q$	a vector of ones of length Q
\mathbf{s}_{nas}	the NAS net analyte signal

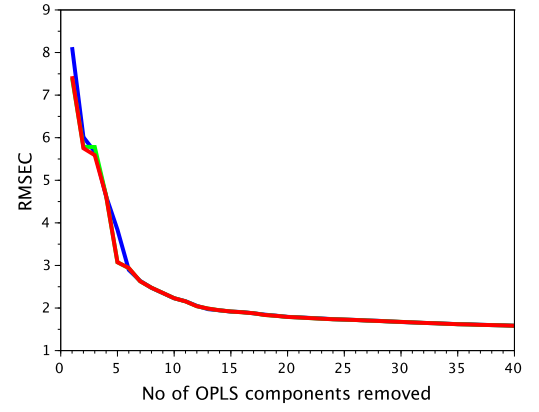
Figure 1:

\mathbf{X}_{opls} calculated with O-PLS	
$\mathbf{w} = \mathbf{X}'_{opls}\mathbf{y}$	$\mathbf{w} = \mathbf{1}_Q$
$\mathbf{t} = \mathbf{X}_{opls}\mathbf{w}$	$\mathbf{t} = \mathbf{X}_{opls}\mathbf{1}_Q$
$c = \mathbf{y}'\mathbf{t}(\mathbf{t}'\mathbf{t})^{-1}$	$c = \mathbf{y}'\mathbf{t}(\mathbf{t}'\mathbf{t})^{-1}$
$\mathbf{b} = c\mathbf{w}$	$\mathbf{b} = c\mathbf{1}_Q$

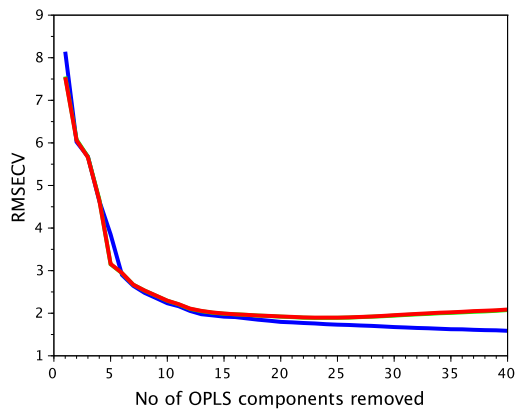
Figure 2:



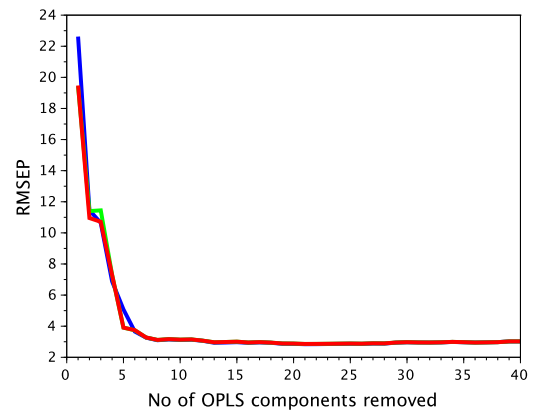
(a)



(b)

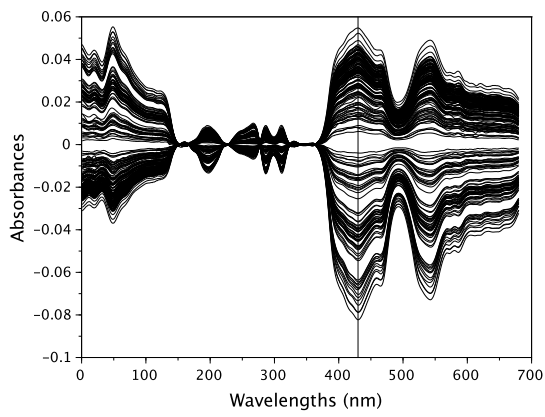


(c)

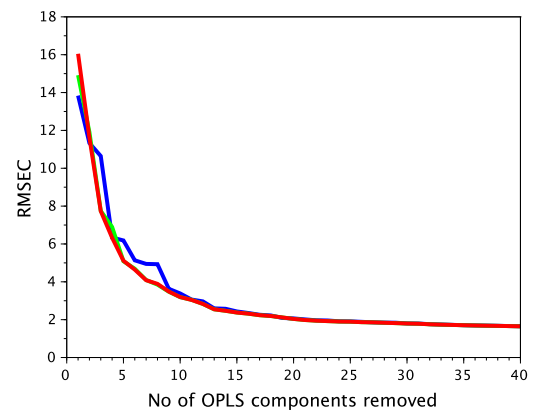


(d)

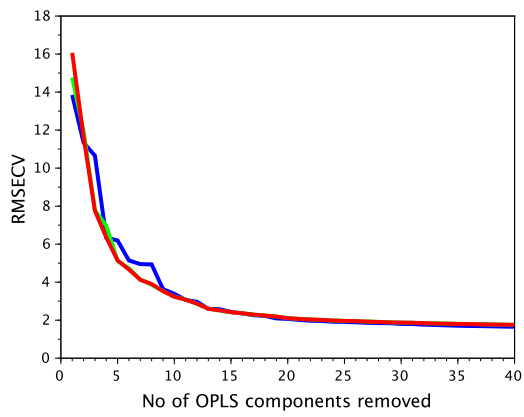
Figure 3:



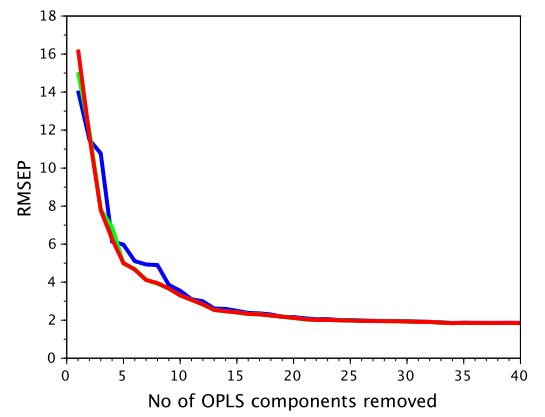
(a)



(b)

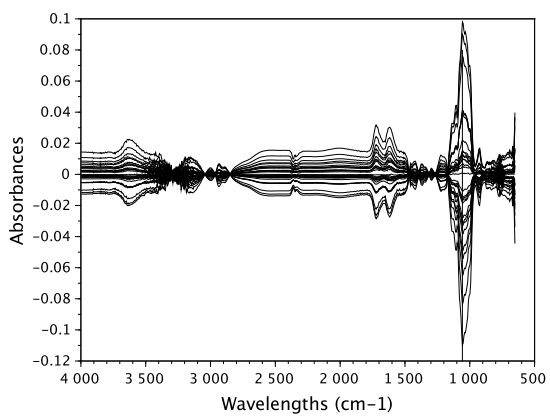


(c)

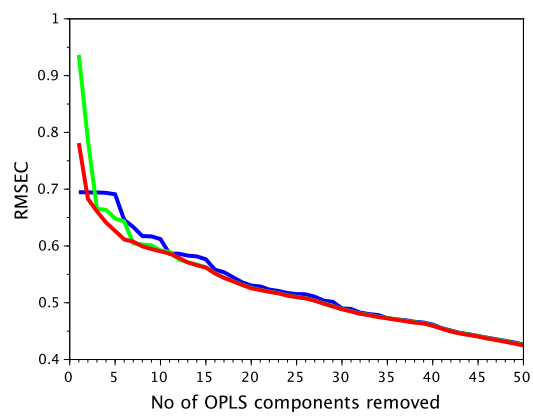


(d)

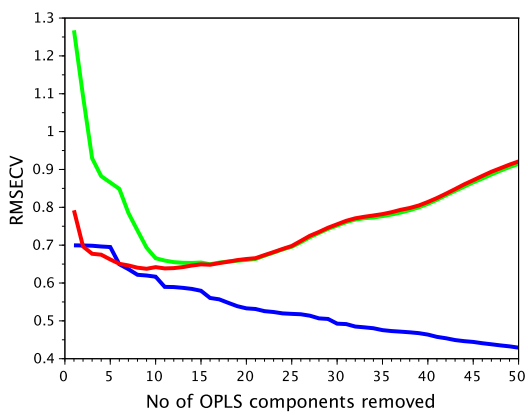
Figure 4:



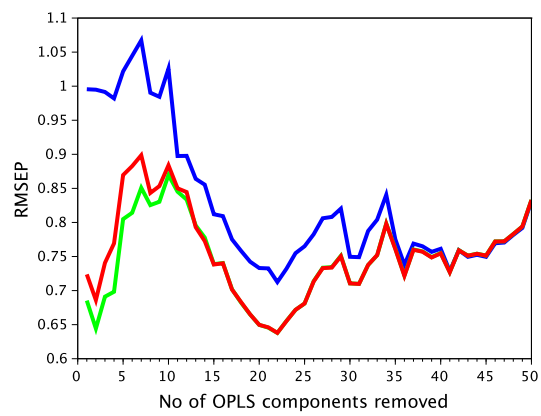
(a)



(b)

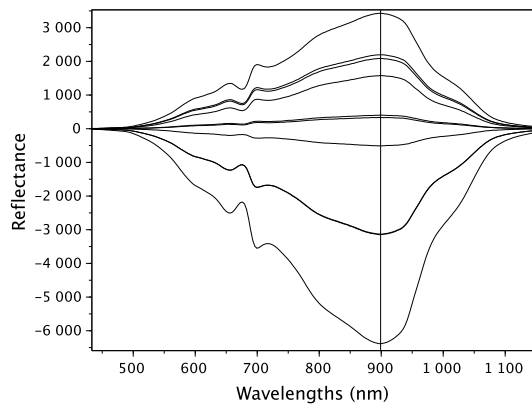


(c)

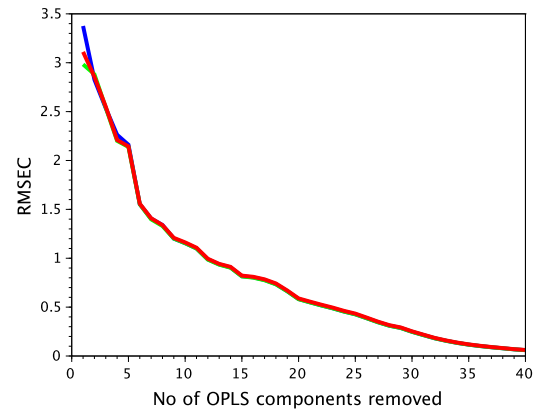


(d)

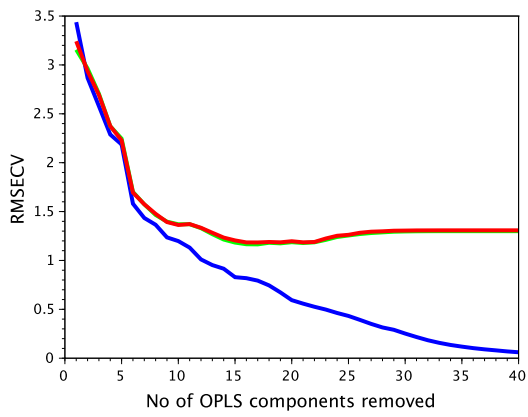
Figure 5:



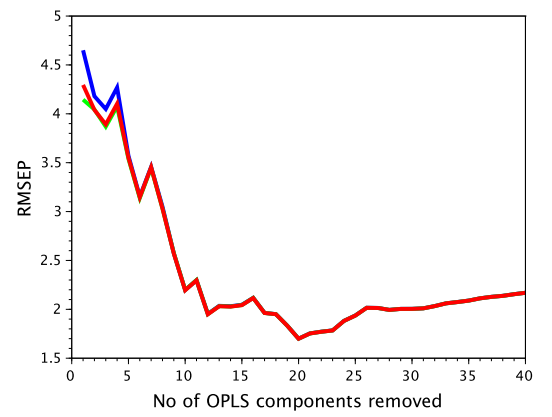
(a)



(b)



(c)



(d)

Figure 6: