



HAL
open science

A family of regression methods derived from standard PLSR

J.C. Boulet, D. Bertrand, G. Mazerolles, Rodolphe Sabatier, J.M. Roger

► **To cite this version:**

J.C. Boulet, D. Bertrand, G. Mazerolles, Rodolphe Sabatier, J.M. Roger. A family of regression methods derived from standard PLSR. *Chemometrics and Intelligent Laboratory Systems*, 2013, 120, p. 116 - p. 125. 10.1016/j.chemolab.2012.11.002 . hal-00780076

HAL Id: hal-00780076

<https://hal.science/hal-00780076>

Submitted on 23 Jan 2013

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.

A family of regression methods derived from standard PLSR

Jean-Claude Boulet^{a,1}, Dominique Bertrand^{b,2}, Gérard Mazerolles^a, Robert Sabatier^c, Jean-Michel Roger^d

^aINRA, UMR1083 Sciences pour l'oenologie, F-34060 Montpellier, France

^bINRA, UR1268 Biopolymères interactions assemblages, F-44316 Nantes, France

^cUM1, EA2415 Biostatistique épidémiologie recherche clinique, F-34093 Montpellier, France

^dIRSTEA, UMR ITAP Information et technologie pour les agroprocédés, F-34191 Montpellier, France

Abstract

We present a new regression method derived from standard PLSR which has a geometric point of view and consists of two projections. In the first, scores are obtained after an oblique projection of the spectra onto the loadings. In the second, the vector of response values is projected orthogonally onto the scores. A metric is introduced for the oblique projection and a new algorithm for calculating the loadings into the variable space is proposed. This work also puts forward a new parameter, a vector, whose different values lead to different regression models with their own prediction abilities, and one of them is the exact form of standard PLSR. This method (called vector orientation decided through knowledge assessment, or VODKA regression) is another way to build least squares regressions using only a few latent variables. We propose two

Email addresses: bouletjc@supagro.inra.fr (Jean-Claude Boulet),
domibertrand@free.fr (Dominique Bertrand), mazeroll@supagro.inra.fr (Gérard Mazerolles),
sabatier@univ-montp1.fr (Robert Sabatier), jean-michel.roger@irstea.fr (Jean-Michel Roger)

¹corresponding author

²present address: Data_ Frame, 25 rue Stendhal, F-44300 Nantes, France

applications to illustrate its performance capabilities.

Keywords: PLSR; metric; orthogonal; oblique; projection; Vodka

1. Introduction

Many current analytical methods are based on spectroscopic techniques such as near-infrared (NIR), mid-infrared(MIR) or raman spectroscopy. The data consist of a set of observations, such as spectra acquired from several samples, and a set of the corresponding analytical results obtained using generally time-consuming analytical techniques. The observations form a first matrix, and the analytical results form a second matrix, which contains the quantitative amounts of one or more compounds of interest, the response variables. The prediction of the response variables can be done using the observations associated with the calibration method. This is a main goal for the development of online, fast and non-destructive analytical methods. Among the proposed methods, partial least squares regression or projection to latent structures regression (PLSR) is the most popular. PLSR is a linear indirect calibration method. PLSR-1 and PLSR-2 are associated with the prediction of one or several response variables respectively. This work concerns only PLSR-1 which is noted PLSR for simplification.

The nonlinear iterative partial least squares (NIPALS) algorithm was proposed by H.Wold for principal component analysis (PCA) calculations [1]. Modifications of this algorithm led H.Wold, S.Wold and H.Martens to develop the first PLSR algorithm [2], which is referred to as "standard PLSR" [3, 4] to avoid confusion with NIPALS for PCA. Later, other algorithms were proposed,

22 such as non-orthogonalized scores PLSR by Martens [5] [6], and SIMPLS by
23 De Jong [3]. The goal of most PLSR algorithms including PLSR-1 has been to
24 produce results close to standard PLSR, at least for PLSR-1. As a consequence,
25 Andersson [7] compared the respective performances of nine PLSR algorithms
26 for speed and the numerical stability. Standard PLSR belonged to the four most
27 stable algorithms, and thus confirmed its status as a reference method.

28 Standard PLSR has been presented from different points of view, such as an
29 application of the Heisenberg uncertainty principle [8], statistical modeling [9],
30 and its geometry [4]. However, most presentations have concerned the algorithm
31 itself [2, 7, 10, 11] and the calculation of different parameters, such as loadings \mathbf{P}
32 and \mathbf{c} , weights \mathbf{W} and scores \mathbf{T} . The parameters $\{ \mathbf{w}_i, \mathbf{t}_i, \mathbf{p}_i, c_i \}$ are calculated
33 simultaneously for a latent variable i , then $\{ \mathbf{w}_{i+1}, \mathbf{t}_{i+1}, \mathbf{p}_{i+1}, c_{i+1} \}$ are calcu-
34 lated for the latent variable $i + 1$, and so on. We show that the same algorithm
35 can be written differently such that each item is calculated separately. A metric
36 $\mathbf{\Sigma}$ is defined as the Moore-Penrose pseudo-inverse of $\mathbf{X}'\mathbf{X}$. The loadings \mathbf{P} are
37 obtained using two elements: a square matrix $\mathbf{X}'\mathbf{X}$ and a vector $\mathbf{r} = \mathbf{X}'\mathbf{y}$, but
38 in a manner different from the Krylov sequence. Then, the matrix \mathbf{T} and the
39 regression vector \mathbf{b} are obtained with \mathbf{P} and $\mathbf{\Sigma}$. Neither \mathbf{W} nor \mathbf{c} are necessary,
40 so they are not calculated.

41 VODKA regressions, which comprise a new family of regression methods, are
42 derived from this new presentation of standard PLSR. The vector $\mathbf{r} = \mathbf{X}'\mathbf{y}$ is
43 considered to be a parameter, which can be replaced by any other vector of the
44 same dimension to calculate the loadings. Each value of \mathbf{r} is associated with a

45 different regression model whose accuracy depends strongly on a relevant choice
46 for \mathbf{r} . Several approaches are proposed for choosing \mathbf{r} and two applications
47 illustrate the proposed method.

48 **2. Theory**

49 The theory is divided into three parts: the standard PLSR algorithm; a
50 rewriting of the standard PLSR including a new algorithm for the calculation
51 of the loadings into the variable space; and the proposal of a new regression
52 method. Vectors are noted in bold lowercase, matrices in bold uppercase, scalars
53 in normal uppercase, variables in normal characters. A spectrum is represented
54 as a column vector, but several spectra form the rows of a matrix, e.g. in \mathbf{X} or
55 \mathbf{X}_G . On the other hand, vectors issued from calculations form the columns of
56 the matrices which gather them, e.g. \mathbf{P} or \mathbf{W} . The transposed forms of vector
57 \mathbf{m} and matrix \mathbf{M} are respectively noted \mathbf{m}' and \mathbf{M}' . Table B.1 summarizes the
58 main notations. In the general case, if optional pretreatments are necessary, e.g.
59 centering, smoothing, orthogonal projection, they should be applied to the raw
60 data previously to yield the calibration dataset (\mathbf{X}, \mathbf{y}) .

61 A data \mathbf{X} of dimension $(N \times P)$ can be explained into the \mathbb{R}^N space spanned
62 by its P column vectors of dimensions $(N \times 1)$, or into the \mathbb{R}^P space spanned by
63 its N line vectors of dimensions $(P \times 1)$. We will focus on this second issue. A
64 metric represented by a square and symmetrical matrix $\mathbf{\Sigma}$ of dimensions $(P \times P)$
65 is associated with \mathbb{R}^P to form a vectorial space, and is used to calculate inner
66 (dot) products and distances as well as to perform projections. To simplify the

67 notations, the term metric is also used for pseudo-metric, that is, when Σ is not
 68 of full rank. The usual Euclidian space is associated with the identity: $\Sigma = \mathbf{I}_P$
 69 and with *orthogonal* projections. Oblique Euclidian spaces are associated with
 70 $\Sigma \neq \mathbf{I}_P$ and with *oblique* projections. Two types of projectors can handle
 71 both orthogonal and oblique projections: *projectors onto* a subspace, and *anti-*
 72 *projectors to* a subspace. For example, the oblique Σ *projector onto* the subspace
 73 spanned by the column vectors of \mathbf{M} is: $\mathcal{P}_M = \Sigma \mathbf{M}(\mathbf{M}'\Sigma\mathbf{M})^{-1}\mathbf{M}'$, and the
 74 oblique Σ *anti-projector to* the subspace spanned by the column vectors of \mathbf{M}
 75 is: $\mathcal{P}_M^\perp = \mathbf{I}_P - \mathcal{P}_M$. The terms \mathcal{P} and \mathcal{T} designate projectors; \mathcal{P}^\perp and \mathcal{T}^\perp
 76 designate anti-projectors. Due to the context, \mathcal{P} and \mathcal{P}^\perp are defined into \mathbb{R}^P
 77 with the metric Σ , whereas \mathcal{T} and \mathcal{T}^\perp are defined into \mathbb{R}^N with the usual
 78 Euclidian metric.

79 *2.1. The standard PLSR algorithm*

80 Standard-PLSR has been described several times, for instance by Geladi
 81 [11]. It aims at building a model for A latent variables. To start, $\mathbf{X}_{1:0} = \mathbf{X}$ and
 82 $\mathbf{y}_{1:0} = \mathbf{y}$. Then a loop calculates the PLSR parameters at each iteration. For
 83 $i = 1, 2, 3, \dots, A$:

$$\mathbf{w}_i = \mathbf{X}'_{1:i-1} \mathbf{y}_{1:i-1} \quad (1)$$

$$\|\mathbf{w}_i\| = 1 \quad (2)$$

$$\mathbf{t}_i = \mathbf{X}_{1:i-1} \mathbf{w}_i \quad (3)$$

$$c_i = \mathbf{y}'_{1:i-1} \mathbf{t}_i (\mathbf{t}'_i \mathbf{t}_i)^{-1} \quad (4)$$

$$\mathbf{p}_i = \mathbf{X}'_{1:i-1} \mathbf{t}_i (\mathbf{t}'_i \mathbf{t}_i)^{-1} \quad (5)$$

$$\mathbf{X}_{1:i} = \mathbf{X}_{1:i-1} - \mathbf{t}_i \mathbf{p}'_i = (\mathbf{I}_N - \mathbf{t}_i (\mathbf{t}'_i \mathbf{t}_i)^{-1} \mathbf{t}'_i) \mathbf{X}_{1:i-1} \quad (6)$$

$$\mathbf{y}_{1:i} = \mathbf{y}_{1:i-1} - \mathbf{t}_i \mathbf{c}'_i = (\mathbf{I}_N - \mathbf{t}_i (\mathbf{t}'_i \mathbf{t}_i)^{-1} \mathbf{t}'_i) \mathbf{y}_{1:i-1} \quad (7)$$

84 then the algorithm returns to equation (1) incrementing i by 1. After A
 85 iterations, the A weight vectors \mathbf{w}_i , the A loadings-for- \mathbf{X} vectors \mathbf{p}_i , the A score
 86 vectors \mathbf{t}_i and the A loadings-for- \mathbf{y} scalars c_i are gathered into the respective
 87 matrices and vector \mathbf{W} , \mathbf{P} , \mathbf{T} and \mathbf{c} . The calibration model for A latent variables
 88 is represented by a regression vector of b-coefficients \mathbf{b} which verifies: $\mathbf{y} =$
 89 $\mathbf{X}\mathbf{b} + \mathbf{e}$, with a vector of errors \mathbf{e} . Let $\hat{\mathbf{b}}$ be an estimation of \mathbf{b} , and let $\hat{\mathbf{y}}$ be
 90 the estimation of \mathbf{y} using $\hat{\mathbf{b}}$. Thus, $\hat{\mathbf{y}} = \mathbf{X}\hat{\mathbf{b}}$, with:

$$\hat{\mathbf{b}} = \mathbf{W}(\mathbf{P}'\mathbf{W})^{-1}\mathbf{c} \quad (8)$$

91 2.2. Rewriting standard PLSR

92 PLSR decomposes a matrix \mathbf{X} into matrices \mathbf{T} , \mathbf{P} and a residual matrix \mathbf{E}
 93 such that: $\mathbf{X} = \mathbf{TP}' + \mathbf{E} = \mathbf{X}^U + \mathbf{E}$, where \mathbf{X}^U represents the information
 94 from \mathbf{X} which is useful for the prediction of \mathbf{y} . Several properties of standard
 95 PLSR, recalled or demonstrated, are reported in Appendix A. Let $\mathbf{\Sigma}$ be the
 96 Moore-Penrose pseudo-inverse of $(\mathbf{X}'\mathbf{X})$; $\mathbf{\Sigma} = (\mathbf{X}'\mathbf{X})^+$. In Property 5, Eq. A.6
 97 gives a new expression of \mathbf{T} , that is: $\mathbf{T} = \mathbf{X}\mathbf{\Sigma}\mathbf{P}(\mathbf{P}'\mathbf{\Sigma}\mathbf{P})^{-1}$. So:

$$\mathbf{X}^U = \mathbf{TP}' = \mathbf{X}\mathbf{\Sigma}\mathbf{P}(\mathbf{P}'\mathbf{\Sigma}\mathbf{P})^{-1}\mathbf{P}' \quad (9)$$

98 This means that the useful information from \mathbf{X} is obtained by an oblique
99 projection of \mathbf{X} onto the loadings \mathbf{P} into \mathbb{R}^P . The matrix \mathbf{T} contains the scores
100 of the observations in the basis of the useful space spanned by the column-
101 vectors of \mathbf{P} . PLSR is also a regression (or orthogonal projection into \mathbb{R}^N) of
102 the reference values \mathbf{y} onto the scores \mathbf{T} , i.e. $\hat{\mathbf{y}} = \mathbf{T}(\mathbf{T}'\mathbf{T})^{-1}\mathbf{T}\mathbf{y}$ [4]. The new
103 expression of \mathbf{T} leads to a new expression of $\hat{\mathbf{b}}$ (see property 6, appendix A):

$$\hat{\mathbf{b}} = \Sigma\mathbf{P}(\mathbf{P}'\Sigma\mathbf{P})^{-1}\mathbf{P}'\Sigma\mathbf{X}'\mathbf{y} \quad (10)$$

104 From a geometric point of view, standard PLSR consists of two projections
105 (see Fig.B.1):

- 106 • First, an oblique projection into the variable space \mathbb{R}^P
107 The vector \mathbf{p}_i is the i^{th} element of a basis of the subspace of \mathbb{R}^P which
108 contains the relevant information for the prediction of \mathbf{y} . The A first vec-
109 tors \mathbf{p}_i form the loading matrix \mathbf{P} of dimensions $(P \times A)$. The information
110 from \mathbf{X} which is useful for the prediction of \mathbf{y} is extracted by an oblique
111 Σ projection of \mathbf{X} onto \mathbf{P} , yielding \mathbf{X}^U ; the scores of its observations into
112 the basis $\{\mathbf{p}_1, \mathbf{p}_2, \dots, \mathbf{p}_A\}$ are given by \mathbf{T} .
- 113 • Then, an orthogonal projection into the observation space \mathbb{R}^N
114 The predicted vector $\hat{\mathbf{y}}$ is the orthogonal projection of \mathbf{y} onto \mathbf{T} . The
115 regression vector $\hat{\mathbf{b}}$ is deduced from this last equation.

116 The calculation of Σ is straightforward. Thus, \mathbf{P} remains the only parameter
117 to calculate to obtain a PLSR model. According to Eq. 6, the deflation of \mathbf{X} at

118 step i is performed into \mathbb{R}^N by multiplying the anti-projector orthogonal to $\mathbf{T}_{1:i}$
119 by \mathbf{X} . Thus, the calculation of the loadings \mathbf{p}_i implies successive steps into \mathbb{R}^N
120 and into \mathbb{R}^P . However, it is shown in Property 7 and Eq. A.7 that the deflation
121 of \mathbf{X} at step i can also be performed into \mathbb{R}^P :

$$\mathbf{X}_{1:i} = \mathcal{T}_{1:i}^\perp \mathbf{X} = \mathbf{X} \mathcal{P}_{1:i}^\perp$$

122 Thus, the calculation steps into \mathbb{R}^N are no longer mandatory, and it becomes
123 possible to rewrite standard PLSR as a calculation of the loadings only into \mathbb{R}^P
124 and independently from the parameters \mathbf{T} , \mathbf{W} or \mathbf{c} . An algorithm is obtained
125 and described in Appendix B. This new approach for the calculation of the
126 loadings raises the following issues.

127 • Geometric building of the \mathbf{p}_i

128 Let z be a point of \mathbb{R}^P and \mathbf{q} the vector from the origin of \mathbb{R}^P to z . Let
129 γ be a positive scalar. The set of points z which verify: $\mathbf{q}'\mathbf{X}'\mathbf{X}\mathbf{q} = \gamma$
130 form an ellipsoid. The value of γ tunes the size of the ellipsoid. We are
131 more interested in the shape of the ellipsoid, whose main directions are the
132 eigenvectors of $\mathbf{X}'\mathbf{X}$. According to Fig. B.2 adapted from [4], the direction
133 of \mathbf{p}_1 is the result of a tangent rotation of $\mathbf{X}'\mathbf{y}$ towards the main direction
134 of the ellipsoid, i.e. the eigenvector associated with the largest eigenvalue
135 of $\mathbf{X}'\mathbf{X}$. The other loadings \mathbf{p}_i are calculated similarly using the deflated
136 forms of $\mathbf{X}'\mathbf{X}$ and $\mathbf{X}'\mathbf{y}$ (see Appendix B). They are compromises between
137 the main spectral information and the spectral information that explains

138 \mathbf{y} .

139 • The Krylov sequences

140 The Krylov sequences can yield the matrix \mathbf{P} directly [12]. A sequence
141 is obtained by the multiplication of two terms, $\mathbf{X}'\mathbf{X}$ and $\mathbf{X}'\mathbf{y}$ for the
142 loadings, the first being risen to power i for loading i . For example,
143 $\mathbf{p}_{K,1} = \mathbf{X}'\mathbf{X}\mathbf{X}'\mathbf{y}$, $\mathbf{p}_{K,2} = \mathbf{X}'\mathbf{X}\mathbf{X}'\mathbf{X}\mathbf{X}'\mathbf{y}$ and so on; $\mathbf{p}_{K,i} = (\mathbf{X}'\mathbf{X})^i\mathbf{X}'\mathbf{y}$.
144 The loadings obtained by the Krylov sequence have no particular proper-
145 ties, moreover they can be highly (but not completely) colinear, so it is
146 convenient to replace them by vectors \mathbf{v}_i obtained after a Gram-Schmidt
147 orthogonalization of the $\mathbf{p}_{K,i}$ [3]. One difference between standard PLSR
148 and the Krylov sequence is that the \mathbf{v}_i are orthogonal with a Euclidian
149 metric: $\mathbf{v}'_i\mathbf{v}_j = 0$ for $i \neq j$, whereas the \mathbf{p}_i are orthogonal with the Σ
150 metric: $\mathbf{p}'_i\Sigma\mathbf{p}_j = 0$ for $i \neq j$. The \mathbf{p}_i and the \mathbf{v}_i span the same subspaces,
151 but only the \mathbf{p}_i constitute a Σ orthonormal basis, according to Eq. 10.
152 Thus the proposed algorithm for the loadings calculation is not based on
153 the Krylov sequences.

154 • The weights \mathbf{W}

155 The relationships between the vectors \mathbf{p}_i and \mathbf{w}_i can be determined di-
156 rectly from the standard PLSR algorithm or from [13]. However, they
157 remain complex and it has not been possible to identify a relationship
158 between the matrices \mathbf{P} and \mathbf{W} of same dimension ($P \times A$). The ap-
159 parent proximity between Eqs. 8 and 10 is misleading because $(\mathbf{P}'\mathbf{W})$ is
160 a bidiagonal matrix [14] and $(\mathbf{P}'\Sigma\mathbf{P})$ is strictly diagonal (see Appendix

161 A, Property 4). Therefore, clearly $\mathbf{W} \neq \mathbf{\Sigma}\mathbf{P}$ and \mathbf{W} is neither used nor
162 explained here.

163 • Properties of $\mathbf{\Sigma}$

164 If $\mathbf{\Sigma}$ is the identity or if $\mathbf{\Sigma}$ is built such that $\mathbf{\Sigma} = \mathbf{M}'\mathbf{M}$ for any matrix \mathbf{M} ,
165 the inner or dot product defined by $f(\mathbf{u}, \mathbf{v}) = \mathbf{u}'\mathbf{\Sigma}\mathbf{v}$ verifies the conditions
166 of a metric or pseudo-metric. Using the singular value decomposition, it
167 applies also to $\mathbf{\Sigma} = (\mathbf{M}'\mathbf{M})^+$. So, under our conditions $\mathbf{\Sigma} = (\mathbf{X}'\mathbf{X})^+$ has
168 the position and the properties of a metric. Moreover, it is a metric asso-
169 ciated with a Mahalanobis distance [15]. Each observation i is represented
170 by a point i into \mathbb{R}^P , or a vector \mathbf{x}_i between the origin and the point i . In
171 the particular case where \mathbf{X} has been previously centered, $(\mathbf{X}'\mathbf{X})/N$ is the
172 variance-covariance matrix of the P spectral variables and $\mathbf{\Sigma} = (\mathbf{X}'\mathbf{X})^+$
173 represents its pseudo-inverse, the coefficient $1/N$ can be dropped. The
174 origin is located at the center of the cloud of points representing the ob-
175 servations of \mathbf{X} . It is important to note that this metric depends on the
176 data, which is not the case for several other regression methods such as
177 Multiple Linear Regression (MLR) or Principal Component Regression
178 (PCR). If spectra are removed or added to \mathbf{X} , then the metric and the
179 inner product are changed. As a consequence, the scores and loadings are
180 different. For instance, suppose that two PLSR models are built: the first
181 one with (\mathbf{X}, \mathbf{y}) for A latent variables, yielding the loadings \mathbf{P}_1 ; the sec-
182 ond one with \mathbf{X}^U defined above and \mathbf{y} , yielding \mathbf{P}_2 . The two metrics are
183 different, $(\mathbf{X}'\mathbf{X})^+$ and $(\mathbf{X}^{U'}\mathbf{X}^U)^+$ respectively, and so are the two models.

184 It is not possible to recover the same models, as reported previously [16],
185 but this is logical and not an artifact. With orthogonal and oblique pro-
186 jections, standard PLSR is based on very simple and basic rules of matrix
187 algebra, and its calculation is correct [16, 17].

188 The added value of a generalized Euclidian metric such as Σ is to remove or
189 decrease the weights associated with certain directions of the space. Thus,
190 these directions become more weak. The Mahalanobis distance lightens
191 the weights of the variables which are the most variable and also takes into
192 account their relationships with the other variables. This property is rele-
193 vant for prediction problems, in which spectral variables of low variability
194 may be good predictors. The Mahalanobis metric can accentuate them
195 whereas the usual Euclidian metric, used in PCR for instance, cannot.

196 Other multivariate methods based on the Mahalanobis distance, such as
197 the Hotelling test [15], are perfectly in accord with standard PLSR.

198 • The regression vector \mathbf{b}

199 In studying the geometry of PLSR [4], Phatak concluded that “ *the PLS*
200 *estimator of the vector of b-coefficients...is an oblique projection of the*
201 *ordinary least square (OLS) estimator* ” . If Σ is of full rank, the OLS
202 estimator is $\hat{\mathbf{b}}_{OLS} = \Sigma \mathbf{X}' \mathbf{y}$ [6, 4]. From Eq. 10, it is straightforward
203 that the oblique projector is $\Sigma \mathbf{P} (\mathbf{P}' \Sigma \mathbf{P})^{-1} \mathbf{P}'$. Thus, the PLSR regression
204 vector is the oblique Σ projection of $\hat{\mathbf{b}}_{OLS}$ onto the space spanned by the
205 loadings \mathbf{p}_i .

- 206 • The scores \mathbf{t}_i and the prediction of new observations

207 In the standard PLSR algorithm, a score \mathbf{t}_i is calculated before the corre-
 208 sponding loading \mathbf{p}_i , which suggests that the loadings calculation depends
 209 on the scores [16]. In fact, the reciprocal is true: from Eq. A.4, scores also
 210 depend on the loadings for their calculation. To conclude, scores and load-
 211 ings are calculated together and they depend on each other. An example
 212 has been shown previously for SIMPLS [3]: standardization of the scores
 213 in SIMPLS simultaneously standardizes the loadings. It is also concluded
 214 from Eq. 9 that the so-called scores are true scores: for an observation \mathbf{x}
 215 and for $i = 1$ to $i = A$, the score value t_i represents the expansion of the
 216 loading \mathbf{p}_i .

217 For a new observation \mathbf{x}_v which does not belong to \mathbf{X} , the prediction \hat{y}_v
 218 is deduced from Eq. 10 :

$$\hat{y}_v = \mathbf{x}'_v \Sigma \mathbf{P} (\mathbf{P}' \Sigma \mathbf{P})^{-1} \mathbf{P}' \Sigma \mathbf{X}' \mathbf{y}$$

$$\hat{y}_v = \mathbf{x}'_v^U \Sigma \mathbf{X}' \mathbf{y} \quad (11)$$

219 The prediction can be viewed as two steps:

- 220 – First step: an oblique Σ projection of \mathbf{x}_v onto \mathbf{P} , yielding the useful
 221 part of \mathbf{x}_v which is \mathbf{x}_v^U ; the scores of \mathbf{x}_v^U in the basis $\{\mathbf{p}_1, \mathbf{p}_2, \dots, \mathbf{p}_A\}$
 222 are: $\mathbf{t}_v = \mathbf{x}_v \Sigma \mathbf{P} (\mathbf{P}' \Sigma \mathbf{P})^{-1}$;

223 – Second step: \hat{y}_v is the Σ inner product between \mathbf{x}_v^U and $\mathbf{X}'\mathbf{y}$.

224 The term $\mathbf{X}'\mathbf{y}$ appears in the regression coefficients of PLSR, PCR and
225 OLSR [4, 6]. However, it also has a specific place in PLSR when build-
226 ing the loadings, as seen in Appendix B, which leads us to propose the
227 following method.

228 2.3. *VODKA regressions: an outcome of the new presentation of standard PLSR*

229 PLSR aims at determining scores that maximize $(\mathbf{t}'_i\mathbf{y})^2$ under the condition:
230 $\|\mathbf{w}_i\| = 1$ [18]. Using Eq. A.4 and the normalization of the \mathbf{t}_i in the proposed
231 algorithm, this constraint can be switched from \mathbb{R}^N to \mathbb{R}^P and expressed as:
232 maximizing $\mathbf{p}'_i\Sigma\mathbf{X}'\mathbf{y}$ under the condition $\mathbf{p}'_i\Sigma\mathbf{p}_i = 1$. The question is whether
233 $\mathbf{X}'\mathbf{y}$ is truly the best vector? Would another vector \mathbf{r} of the same dimension be
234 more representative of the relevant information from \mathbf{X} which explains \mathbf{y} ?

235 This issue has been discussed previously in a context independent from
236 PLSR. The net analyte signal (NAS) [19] is the most condensed spectral in-
237 formation about the compound to be predicted, and it also constitutes the
238 principle of direct calibration [20]. Two definitions of the NAS have been pro-
239 posed [21]: (1) *the NAS for a component is the part of its pure spectrum which is*
240 *orthogonal to the pure spectra of the other constituents;* (2) *The NAS is the part*
241 *of the gross spectrum that is useful for prediction.* According to the first defini-
242 tion, if the pure spectrum \mathbf{k} of the compound to be quantified is known, and if
243 all other influences have been characterized as spectra or loadings and merged
244 into the matrix \mathbf{D} , the NAS can be estimated: $\mathbf{s}_{nas} = (\mathbf{I}_P - \mathbf{D}(\mathbf{D}'\mathbf{D})^{-1}\mathbf{D}')\mathbf{k}$.

245 The regression coefficients obtained from PLSR or other regression models con-
246 stitute another estimation of the NAS [21]. Moreover, in certain conditions, the
247 regression vector of PLSR can be exactly the NAS [9]. Therefore, if a good
248 approximation of the NAS can be obtained with additional information, this
249 justifies using the NAS for value of \mathbf{r} .

250 When no other information than (\mathbf{X}, \mathbf{y}) is available, this vector \mathbf{r} can be built
251 using only \mathbf{X} and \mathbf{y} . Moreover, \mathbf{r} is set in the form: $\mathbf{r} = \mathbf{X}'g(\mathbf{y})$, where $g(\mathbf{y})$ is
252 a vector function of \mathbf{y} . The vector $\mathbf{X}'\mathbf{y}$ from standard PLSR with $g(\mathbf{y}) = \mathbf{y}$ is
253 collinear to the mean of the observations, weighted by the concentrations of the
254 response variable. The higher the concentrations are, the more the correspond-
255 ing spectra will contribute to \mathbf{r} . When \mathbf{X} or \mathbf{y} are centered, the intermediate
256 values in \mathbf{y} are not well represented in \mathbf{r} , which is obtained mainly by subtracting
257 the spectra associated with the lowest values of the response variable from the
258 spectra associated with the highest values of the response variable. A different
259 weighting is possible which is useful for taking into account the nonlinearities
260 between the response variable and the observations, or for simply overexpressing
261 the contribution of the response variable. For example, four different functions
262 $g(\mathbf{y})$ can be developed along with several other possibilities. Let $g_1(\mathbf{y}) = \mathbf{y}^2$,
263 $g_2(\mathbf{y}) = \exp(\mathbf{y})$, $g_3(\mathbf{y}) = \sqrt{\mathbf{y}}$ and $g_4(\mathbf{y}) = \log(\mathbf{y})$ be four functions in which
264 each result is a vector obtained after: raising each of the elements of \mathbf{y} to the
265 power of 2 (g_1); taking the exponential of the elements of \mathbf{y} (g_2); the square
266 root of the elements of \mathbf{y} (g_3); or the logarithm of the elements of \mathbf{y} (g_4). If we
267 suppose that all the elements of \mathbf{y} are larger than 1, then using $g_1(\mathbf{y})$ or $g_2(\mathbf{y})$

268 will accentuate the weight of the highest concentrations whereas using $g_3(\mathbf{y})$ or
269 $g_4(\mathbf{y})$ will reduce the weight of the highest concentrations.

270 New regression models can be obtained by switching the vector $\mathbf{X}'\mathbf{y}$ of stan-
271 dard PLSR with another vector \mathbf{r} of the same dimension. The loadings are first
272 calculated with the following algorithm derived from Appendix B:

273 • Step 1:

$$\mathbf{p}_1 = \mathbf{X}'\mathbf{X}\mathbf{r}$$

$$\mathbf{p}_1 = \mathbf{p}_1(\mathbf{p}'_1\Sigma\mathbf{p}_1)^{-0.5}$$

274 • Step $i + 1$:

$$\mathbf{p}_{i+1} = \mathcal{P}'_{1:i}\mathbf{X}'\mathbf{X}\mathcal{P}'_{1:i}\mathbf{r}$$

$$\mathbf{p}_{i+1} = \mathbf{p}_{i+1}(\mathbf{p}'_{i+1}\Sigma\mathbf{p}_{i+1})^{-0.5}$$

275 Other parameters, such as scores and regression coefficients are obtained
276 with Eqs. A.6 and 10 respectively. From Figure B.2, it is straightforward that a
277 choice for \mathbf{r} different from $\mathbf{X}'\mathbf{y}$ leads to a new orientation of the \mathbf{p}_i vectors and of
278 the useful space they span. For the models to be relevant, however, the choice of
279 \mathbf{r} cannot be left to chance. It is advisable to first review the available knowledge
280 regarding the data to be processed, and then to choose the most useful elements
281 (e.g. an estimation of the NAS) to build \mathbf{r} . This step is summarized in the
282 proposed name for this method: vector orientation decided through knowledge
283 assessment (VODKA) regression.

284 **3. Applications**

285 We used VODKA in two applications: the quantification of ethanol in fer-
286 menting wines and of mannoproteins in polysaccharide extracts. These appli-
287 cations presented several respective differences: the samples were liquid and
288 powder, the spectra were NIR and MIR, and they employed different numbers
289 of samples for calibration. However, the main difference was in the nature of
290 the compounds of interest: ethanol is a single compound, whereas the manno-
291 proteins are a family of similar molecules with slight differences.

292 *3.1. Material and methods*

293 For each of the two applications, spectra \mathbf{X} were centered, but not the ref-
294 erence values \mathbf{y} which contained the raw values. Four models were calculated
295 depending on the choice of the vector \mathbf{r} , see Table B.2. For the classical PLSR
296 model m_{plsR} , $g(\mathbf{y}) = \mathbf{y}$, so \mathbf{r} was set to $\mathbf{X}'\mathbf{y}$, which is exactly the same whether
297 \mathbf{y} is centered or not, provided \mathbf{X} is centered. Two other functions were chosen,
298 leading to the model m_{y^2} with $g(\mathbf{y}) = \mathbf{y}^2$ and $\mathbf{r} = \mathbf{X}'\mathbf{y}^2$, and model $m_{exp(\mathbf{y})}$
299 with $g(\mathbf{y}) = exp(\mathbf{y})$ and $\mathbf{r} = \mathbf{X}'exp(\mathbf{y})$. Finally, for model m_{nas} , $\mathbf{r} = \mathbf{s}_{nas}$ where
300 \mathbf{s}_{nas} is an estimation of the NAS. The root-mean square error of cross-validation
301 ($RMSECV$) was calculated for each of the four models and used to determine
302 the optimal number of latent variables, which were chosen among the lowest
303 values and according to the shape of the $RMSECV$. All the calibration mod-
304 els were applied to the validation dataset $(\mathbf{X}_V, \mathbf{y}_V)$, then characterized by the
305 root-mean square error of prediction ($RMSEP$) and the norm of the regression

306 vectors \mathbf{b} . Low *RMSEP* values and low norms for \mathbf{b} are expected for the best
307 models.

308 *3.1.1. Application 1: quantification of the ethanol concentration in fermenting*
309 *musts*

310 This application aimed at quantifying ethanol in wines using near-infrared
311 spectroscopy. Spectra were acquired with a Jasco spectrophotometer (optical
312 length 1 mm, wavelength range 500 – 2500 nm, acquisition step 2 nm, water ref-
313 erence) at the Skalli-Fortant de France winery (Sète, France). The wavelength
314 range was reduced to 500 – 1898 nm such that $P = 700$ spectral variables, for
315 reasons of compatibility with fiber optics. A vertical shift of baselines was per-
316 formed such that each corrected absorbance at 1170 nm was null. The original
317 data was split into two datasets: (1) \mathbf{X} : the N first 480 samples of ferment-
318 ing musts, for calibration; and (2) \mathbf{X}_V : the last 1000 samples of fermenting
319 musts, for validation. A subset of \mathbf{X} was called \mathbf{X}_m and contained 165 musts
320 before fermentation, without ethanol. The reference values were expressed as
321 a percent volume, e.g. 12% vol., and were measured with a WineScan MIR
322 spectrophotometer (Foss) whose standard error of prediction was estimated to
323 be 0.2%. The two datasets formed vectors \mathbf{y} and \mathbf{y}_V of dimensions (480×1)
324 and (1000×1) respectively, and were associated with \mathbf{X} and \mathbf{X}_V . These data
325 were completed by the pure spectra of ethanol, water, glycerol and lactic acid
326 (\mathbf{k}_{etch} , \mathbf{k}_{water} , \mathbf{k}_{glyc} , \mathbf{k}_{lact}) acquired against air using the same Jasco spectropho-
327 tometer.

328 Matrices \mathbf{X} and \mathbf{X}_V were centered around the row mean of \mathbf{X} . The pure

329 spectrum of ethanol yielded \mathbf{k} . A matrix \mathbf{D} of dimensions (700×7) gathered all
330 the influences to remove: the first four eigenvectors of a PCA performed onto
331 \mathbf{X}_m ; and the spectra \mathbf{k}_{water} , \mathbf{k}_{glyc} and \mathbf{k}_{lact} . The NAS represented by the vector
332 \mathbf{s}_{nas} was calculated using \mathbf{k} and \mathbf{D} , as described above. The *RMSECV* and
333 the *RMSEP* were calculated from 1 to 20 latent variables.

334 3.1.2. Application 2: quantification of mannoproteins in wines

335 Wines produced from healthy grapes contain three main families of polysac-
336 charides: arabinogalactan-proteins (AGPs), rhamno-galacturonan II (RG-II)
337 and mannoproteins (MPs). We were interested in quantifying the MPs which
338 impact on the physical and sensory properties of wines.

339 The calibration dataset of 40 samples was built with powder mixtures of
340 four pure fractions: RG-II, neutral MP (MP0), neutral (AGP0) and acidic
341 (AGP4) AGPs, as previously described [22]. The validation dataset consisted
342 of powdered polysaccharide extracts from 65 wines. Spectra acquisition were
343 performed with an Avatar 360 MIR spectrophotometer (Nicolet) equipped with
344 an attenuated total reflectance cell and germanium crystal. The spectra in the
345 range of $950 - 1850 \text{ cm}^{-1}$ were standardized to adjust the absorbance to 0 and
346 1 for the respective wavenumbers of 1850 cm^{-1} and 1035 cm^{-1} , with the latter
347 corresponding to the highest glucoside bond absorbance peak. MP reference
348 values were deduced for the calibration dataset from the experimental design,
349 and were obtained for the validation dataset by a chemical method involving
350 hydrolysis and quantification by gas chromatography of the alditol acetates of
351 the neutral sugars. The raw data were \mathbf{X} , (40×951) , and \mathbf{X}_V , (65×951) , with

352 the associated response variables \mathbf{y} and \mathbf{y}_V . The NAS was estimated with \mathbf{k} ,
353 the pure spectrum of MP0, and \mathbf{D} , which contained the spectra of RG-II and
354 seven different fractions of AGPs.

355 3.2. Results

356 3.2.1. Application one

357 *Standard errors of cross validation and of prediction.* According to the *RMSECV*
358 values, the optional numbers of latent variables chosen for models m_{plsr} , m_{nas} ,
359 m_{y^2} and $m_{exp(y)}$ were 10, 10, 5 and 10 respectively (Table B.3). The corre-
360 sponding *RMSEP* were similar: 0.92, 0.92, 0.92 and 0.90. The chosen number
361 of latent variables was optimal for PLSR because it corresponded to the lower
362 *RMSEP*. However this was not the case for the three other models, which
363 presented their best *RMSEPs* for several latent variables within a range equiv-
364 alent to or slightly better than the best one of m_{plsr} . Thus, a small amount of
365 error in the number of latent variables may be less serious for m_{nas} , m_{y^2} and
366 $m_{exp(y)}$ than for m_{plsr} .

367 *Comparison of the regression vectors \mathbf{b} .* For each of the four models, the regres-
368 sion vectors leading to the lower *RMSEPs* are presented in Figure B.3. It is
369 noteworthy that the vectors of models m_{plsr} and m_{y^2} are very similar, but not
370 equal; they seem to converge towards the same optimal solution. On the other
371 hand, the b-coefficients vector for $m_{exp(y)}$ presents certain common peaks with
372 the m_{plsr} and the m_{y^2} vectors, but its global shape is very different. The shape
373 of the b-coefficients vector for m_{nas} has nothing in common with the three oth-

374 ers; but it is very close to the NAS (Figure B.3(b)). Thus three different (but
375 equivalent in terms of prediction) solutions were identified by the four models.

376 The shape of \mathbf{r} provides little information about the shape of the b-coefficients
377 and the quality of prediction. For example, although \mathbf{r} and \mathbf{b} are very different
378 for m_{plsr} and very similar for m_{nas} , these two models are nevertheless equivalent
379 in terms of prediction errors.

380 The norms of the regression vectors for the four models and different latent
381 variables are compared in Table B.4. In general, the norms increase with the
382 number of latent variables. However, the standard PLSR model is the only
383 model for which the norm increases steadily. For the three other models, the
384 norm can decrease locally, for example m_{nas} and $m_{exp(y)}$ decrease between 10
385 and 11 latent variables. As a consequence, models with high numbers of latent
386 variables can yield regression vectors of low norm, in contrast to standard PLSR.

387 3.2.2. Application two

388 *Standard errors of cross validation and of prediction.* The results are summa-
389 rized in Table B.5. The *RMSECV* values of the four models, calculated for
390 1 to 30 latent variables, decrease steadily. The choice for A should be close to
391 30, which is not reasonable for only 40 samples. Thus, the *RMSECV* cannot
392 be used here for the determination of A . However, *RMSEP* values shows that
393 good predictions were obtained with a low number of latent variables. In this
394 example, the lowest *RMSEP* for standard PLSR is 0.058 for seven latent vari-
395 ables. The model m_{nas} is difficult to compare to m_{plsr} , since it presents lower
396 *RMSEP* values but was obtained with a higher number of latent variables. On

397 the other hand, it is obvious that m_{y^2} and $m_{exp(y)}$ perform better than m_{plsr} :
398 they present equivalent or lower *RMSEP* values in a range of 12 latent vari-
399 ables, compared to the unique latent variable for m_{plsr} in the same range of
400 *RMSEP* values.

401 *Comparison of the regression vectors.* We compared the vector $\mathbf{X}'\mathbf{y}$ and the
402 NAS to the regression vectors resulting from the four models, Figure B.4. Two
403 spectral regions were identified: (1) $1200 - 950 \text{ cm}^{-1}$, in which the same peaks
404 are found in all six examples, that is $\mathbf{X}'\mathbf{y}$, the NAS and the four regression
405 vectors; and (2) $1850 - 1200 \text{ cm}^{-1}$, in which $\mathbf{X}'\mathbf{y}$ and the NAS have similar
406 shapes, but the b-coefficients resemble noise near the baseline, with m_{plsr} being
407 the most noisy.

408 In this application, the regression vectors tended towards a unique optimal
409 solution, in contrast to Application one in which several solutions were identi-
410 fied. We suspect that the spectral information of each chemical compound in
411 the Application two is more precise, specific peaks more localized within the
412 spectrum and thus more easily identified in the MIR spectra than in the NIR
413 spectra of Application one, leading to a unique solution. The relevant informa-
414 tion in the $1200 - 950 \text{ cm}^{-1}$ region appears to be treated in the same way by
415 the four models. Thus, the differences in their predictive abilities lie in the way
416 noise is minimized in the $1850 - 1200 \text{ cm}^{-1}$ region.

417 We compared the norms of the regression vectors for the four models and the
418 different latent variables (Table B.6) and found the comparisons to be similar
419 to those made for Application one. Moreover, for the selected latent variables,

420 the regression vector of standard PLSR had the highest norm, which is not a
421 good characteristic.

422 4. Discussion and conclusion

423 We have written a new algorithm for standard PLSR. We introduced a new
424 parameter, a matrix Σ which has the position of a metric or pseudo-metric. We
425 also dropped two parameters: the weights \mathbf{W} and the loadings \mathbf{c} for \mathbf{y} . We
426 showed that the deflation of \mathbf{X} into \mathbb{R}^N can be replaced by a deflation into \mathbb{R}^P
427 by means of the loadings and the metric Σ . These results have several con-
428 sequences. Firstly, the extraction of the useful information from \mathbf{X} is the Σ
429 oblique *projection* of \mathbf{X} onto the *latent structures* represented by the loadings,
430 according to PLSR. The work of Phatak is completed by the identification of
431 the metric and the space. Secondly, from a mathematical point of view, a metric
432 and a basis of a subspace are independent elements of a vectorial space. Either
433 Σ or \mathbf{P} can be replaced into Eqs. A.6 and 10 by another metric or another
434 matrix containing an other basis respectively, and the mathematics remain cor-
435 rect. This point specifically addresses the metric that attempts to weight the
436 variables and to take into account their collinearity. The more comprehensive
437 the observational data are within the space spanned by the spectra, the bet-
438 ter the metric will be defined. Suppose that the observations associated with
439 a response variable and represented by \mathbf{X} have been extracted from a much
440 larger database \mathbf{X}_t . For example, the whole database of a spectrometer that
441 has been used for years for calibration is \mathbf{X}_t . Then, if a new parameter is cali-

442 brated, the response values are only available for a small subset of \mathbf{X}_t which is
443 \mathbf{X} . Since a much better representation of the space spanned by the spectra can
444 be obtained from \mathbf{X}_t than from \mathbf{X} , $\mathbf{\Sigma}_t = (\mathbf{X}'_t \mathbf{X}_t)^+$ is a better estimation of the
445 metric to be applied than $\mathbf{\Sigma} = (\mathbf{X}' \mathbf{X})^+$. Thus the PLSR should be built with
446 the metric $\mathbf{\Sigma}_t$ rather than $\mathbf{\Sigma}$. This applies also to the cross-validation, which
447 consists of extracting observations for prediction and calculating a model with
448 the remaining observations. Instead of calculating a new metric at each loop, it
449 seems mathematically more logical to determine a unique metric which is used
450 independently from the set of observations being processed at any given time.
451 This also would increase the calculation speed.

452 This new presentation of standard PLSR allows the deflation of \mathbf{X} into \mathbb{R}^P ,
453 yielding a new algorithm for the building of the loadings. This puts forward an
454 inner parameter, a vector \mathbf{r} whose value is $\mathbf{X}'\mathbf{y}$ for standard PLSR. We show
455 that it is also a tool for orienting the calculation of the loadings, that is, the
456 information contained within the spectra used by the model for prediction. The
457 choice of $\mathbf{r} = \mathbf{X}'\mathbf{y}$ is very likely to yield good predictions, but in certain cases
458 other choices for \mathbf{r} may be possible and yield VODKA regressions. Like stan-
459 dard PLSR, VODKA models use orthogonal scores and $\mathbf{\Sigma}$ -orthogonal loadings.
460 Among the first motivation of the authors of PLSR was to find relevant latent
461 variables (orthogonal scores into \mathbb{R}^N , loadings into \mathbb{R}^P) that could explain an
462 observation [23]. Our proposed modification aims at doing that also, but dif-
463 ferent choices for \mathbf{r} can lead to dramatically different solutions, whereas the
464 different PLSR algorithms tend to converge towards the same solution [7]. To

465 avoid confusion, the PLSR acronym has not been linked to VODKA, even if
466 standard PLSR is one of its particular solutions.

467 Standard PLSR and VODKA both use orthogonal scores into \mathbb{R}^N associated
468 with orthogonal loadings into \mathbb{R}^P with the metric Σ to identify latent variables
469 relevant for prediction. However, these two methods are different in terms of
470 means. When calculating the loadings by the new algorithm, standard PLSR
471 seeks into (\mathbb{R}^P, Σ) vectors of norm 1 whose inner product with $\mathbf{X}'_{1:i}\mathbf{y}$ is maxi-
472 mal; and Vodka seeks into (\mathbb{R}^P, Σ) vectors of norm 1 whose inner product with a
473 vector \mathbf{r}_i is maximal, with $\mathbf{r}_i = \mathcal{P}'_{1:i}\mathbf{r}$. The results of the two methods can differ.
474 Standard PLSR takes into account all the information within the observations
475 which is correlated to the response variable; that is, direct information from the
476 response variable itself, plus indirect information provided by other compounds
477 correlated to the response variables. This property is a strength because good
478 models are often obtained with the indirect information, and sometimes without
479 any contribution by the direct information. It is also a weakness because un-
480 expected modifications in the correlations among the different compounds will
481 modify the predictions and lead to lower robustness. VODKA can be a solution
482 to these situations. One strength is that it allows the introduction of external
483 and selected knowledge to be introduced in the calculation process of the regres-
484 sion, and this enhances performance. In the two applications examined here,
485 VODKA was able to produce models with prediction errors equivalent to or
486 slightly lower than the errors associated with standard PLSR models. Perhaps
487 the best result of VODKA has been to yield large ranges of the best prediction

488 errors, which implies that these models are more stable near the optimal number
489 of latent variables.

490 In our first application of VODKA, we obtained three different models with
491 comparable performances. This raises the question of the unicity of a best re-
492 gression model. According to definition 1 of the NAS, one unique and best
493 solution should exist. However definition 2 allows a wider interpretation and is
494 more in accordance with the functioning of PLSR. If several compounds con-
495 tribute to the prediction of the compound of interest, increasing the contribution
496 of one will decrease the contribution of the others, so several equivalent solutions
497 may exist. An other possible explanation lies in the nature of the spectra, for
498 example NIR vs MIR, as suggested by the second application. However, we do
499 not have enough elements or applications to identify the most likely hypothesis
500 for these differences.

501 To conclude, VODKA regression provides an opportunity to take into ac-
502 count all the available information, not just that from the calibration dataset,
503 allowing regression models to be built which can present some advantages over
504 those produced by standard PLSR.

505 **5. Acknowledgements**

506 This work was made possible within the IRVIN program, was carried out
507 with the Skalli winery and was supported by the Languedoc-Roussillon region.
508 The authors thank all the people who contributed to this manuscript, and es-
509 pecially Steven Brown for his relevant advice.

510 **References**

- 511 [1] H.Wold, in: Multivariate analysis, Krishnaiah (ed), Academic Press, Lon-
512 don, 1966.
- 513 [2] S.Wold, H.Martens, H.Wold, The multivariate calibration method in chem-
514 istry solved by the pls model, in: Proc. conf. matrix pencils, lecture notes
515 in mathematics, Springer-Verlag, Heidelberg, A.Ruhe, B.Kagstrom, 1983,
516 pp. 286–293.
- 517 [3] S.DeJong, Simpls: an alternative approach to partial least squares regres-
518 sion, Chemom. Intell. Lab. Syst. 18 (1993) 251–263.
- 519 [4] A.Phatak, S.DeJong, The geometry of partial least squares, J. Chemom.
520 11 (1997) 311–338.
- 521 [5] H.Martens, T.Naes, in: Near infrared technology in the agricultural and
522 food industries, Williams-Norris (eds), Amer. Assn. of Cereal Chemists, St
523 Paul, 1987.
- 524 [6] H.Martens, T.Naes, Multivariate Calibration, Wiley, Chichester, 1989.
- 525 [7] M.Andersson, A comparison of nine pls1 algorithms, J. Chemom. 23 (2009)
526 518–529.
- 527 [8] A.Hoskuldsson, The h principle: new ideas, algorithms and methods in
528 applied statistics and physics, Chemom. Intell. Lab. Syst. 23 (1994) 1–28.
- 529 [9] B.Nadler, R.R.Coifman, Partial least squares, beer’s law and the net ana-
530 lyte signal: statistical modeling and analysis, J.Chemom. 19 (2005) 45–54.

- 531 [10] A.Hoskuldsson, Pls regression methods, *J.Chemom.* 2 (1988) 211–228.
- 532 [11] P.Geladi, B.Kowalski, Partial least squares regression: a tutorial,
533 *Anal.Chim.Acta* 185 (1986) 1–17.
- 534 [12] I.S.Helland, Partial least squares regression and statistical models, *Scand.*
535 *J. of Statist.* 17 (1990) 97–114.
- 536 [13] R.J.Pell, L.S.Ramos, R.Manne, The model space in partial least squares
537 regression, *J. Chemom.* 21 (2007) 165–172.
- 538 [14] R.Ergon, Re-interpretation of nipals results solves pls regression inconsistency prob-
539 lem, *J. Chemom.* 23(2) (2009) 72–75.
- 540 [15] R.DeMaesschack, D.JouanRimbaud, D.L.Massart, The mahalanobis dis-
541 tance, *Chemom. Intell. Lab. Syst.* 50 (2000) 1–18.
- 542 [16] R.Manne, R.J.Pell, L.S.Ramos, The pls model space: the inconsistency
543 persists, *J.Chemom.* 23 (2009) 76–77.
- 544 [17] S.Wold, M.Hoy, H.Martens, J.Trygg, F.Westad, J.MacGregor, B.M.Wise,
545 The pls model space revisited, *J.Chemom.* 23 (2009) 67–68.
- 546 [18] J.Trygg, Parsimonious multivariate models, Ph.D. thesis, Umea University,
547 Sweden (2001).
- 548 [19] A.Lorber, N.K.M.Faber, B.R.Kowalski, Net analyte signal calculation in
549 multivariate calibration, *Anal. Chem.* 69(8) (1997) 1620–1626.
- 550 [20] R.Marbach, A new method for multivariate calibration., *Journal of Near*
551 *Infrared Spectroscopy* 13 (2005) 241–254.

- 552 [21] J.Ferre, N.K.M.Faber, Net analyte signal calculation for multivariate cali-
553 bration, *Chemom. Intell. Lab. Syst.* 69 (2003) 123–136.
- 554 [22] J.C.Boulet, T.Dococ, J.M.Roger, Improvement of calibration models using
555 two successive orthogonal projection methods, application to quantification
556 of wine polysaccharides, *Chemom. Intell. Lab. Syst.* 87 (2007) 295–302.
- 557 [23] S.Wold, Personal memories of the early pls development, *Chemom. Intell.*
558 *Lab. Syst.* 58 (2001) 83–84.
- 559 [24] M.Tenenhaus, *La régression PLS*, Technip, Paris, 1998.
- 560 [25] B.S.Dayal, J.F.MacGregor, Improved pls algorithms, *J. Chemom.* 11 (1997)
561 73–85.

562 **Appendix A. Properties of standard PLSR**

563 Several properties of standard PLSR are recalled or developed. They are
564 logically ordered to demonstrate the seventh and last property, each of them
565 relying on the previous ones.

- 566 • Property 1: the projection of the \mathbf{t}_i onto the space spanned by the columns
567 of \mathbf{X}

568 The \mathbf{t}_i belong to the subspace of \mathbb{R}^N spanned by the columns of \mathbf{X} [10].

569 Thus their orthogonal projection onto \mathbf{X} is invariant:

$$\mathbf{X}(\mathbf{X}'\mathbf{X})^+\mathbf{X}'\mathbf{t}_i = \mathbf{t}_i \quad (\text{A.1})$$

- 570 • Property 2: the relationship between $(\mathbf{t}'_i \mathbf{t}_i)$ and $(\mathbf{p}'_i \Sigma \mathbf{p}_i)$

571

572 Equation 5 can be simplified [24] leading to:

$$\mathbf{p}_i = \mathbf{X}' \mathbf{t}_i (\mathbf{t}'_i \mathbf{t}_i)^{-1} \quad (\text{A.2})$$

573 Thus:

$$\begin{aligned} \mathbf{p}'_i \Sigma \mathbf{p}_i &= (\mathbf{t}'_i \mathbf{t}_i)^{-1} \mathbf{t}'_i \mathbf{X} (\mathbf{X}' \mathbf{X})^+ \mathbf{X}' \mathbf{t}_i (\mathbf{t}'_i \mathbf{t}_i)^{-1} \\ \mathbf{p}'_i \Sigma \mathbf{p}_i &= (\mathbf{t}'_i \mathbf{t}_i)^{-1} \end{aligned} \quad (\text{A.3})$$

575 In his presentation of SIMPLS [3], De Jong wrote a similar equation:

576 $\mathbf{P}' \Sigma \mathbf{P} = \mathbf{T}' \mathbf{T} = \mathbf{I}_A$, which was developed under the hypothesis that the \mathbf{t}_i

577 had been normed and thus, for this particular case, $(\mathbf{t}'_i \mathbf{t}_i)^{-1} = (\mathbf{t}'_i \mathbf{t}_i) = 1$.

578 However, for standard PLSR and current versions of SIMPLS for which

579 the scores are not normed, only Equation A.3 is valid.

- 580 • Property 3: expression of \mathbf{t}_i in terms of \mathbf{X} , \mathbf{p}_i and Σ

581

582 Each term of Equation A.2 is multiplied on the left by $\mathbf{X} \Sigma$, that is

583 $\mathbf{X} (\mathbf{X}' \mathbf{X})^+$; the terms are permuted and simplified according to Equation

584 A.1 and then Equation A.3:

$$\mathbf{t}_i = \mathbf{X} \Sigma \mathbf{p}_i (\mathbf{p}'_i \Sigma \mathbf{p}_i)^{-1} \quad (\text{A.4})$$

- 585 • Property 4: Σ -orthogonality of the \mathbf{p}_i

586

587 For $i \neq j$ [10]:

$$\mathbf{p}'_i \Sigma \mathbf{p}_j = 0 \quad (\text{A.5})$$

588 Moreover the matrix $\mathbf{P}'\Sigma\mathbf{P}$ is diagonal and the nonzero term at the i^{eme}

589 line and i^{eme} column is $\mathbf{p}'_i \Sigma \mathbf{p}_i$.

- 590 • Property 5: expression of \mathbf{T} in terms of \mathbf{X} , \mathbf{P} and Σ

591

592 Because of Property 4, Equation A.4 leads to the Property 5:

$$\mathbf{T} = \mathbf{X}\Sigma\mathbf{P}(\mathbf{P}'\Sigma\mathbf{P})^{-1} \quad (\text{A.6})$$

- 593 • Property 6: expression of the b-coefficients using \mathbf{X} , \mathbf{y} , Σ and \mathbf{P}

594

595 $\hat{\mathbf{y}}$ is obtained after an orthogonal projection of \mathbf{y} onto \mathbf{T} [4], then the

596 value of \mathbf{T} is replaced by its value from Equation A.6 and the expression

597 simplified:

$$\hat{\mathbf{y}} = \mathbf{X}\Sigma\mathbf{P}(\mathbf{P}'\Sigma\mathbf{P})^{-1}\mathbf{P}'\Sigma\mathbf{X}'\mathbf{y}$$

598 Equation 10 is straightforward.

599 • Property 7: The relationship between $\mathbf{X}\mathcal{P}_{1:i}^\perp$ and $\mathcal{T}_{1:i}^\perp\mathbf{X}$

600

601 Let $\mathcal{P}_{1:i}^\perp$ be the oblique Σ anti-projector to $\mathbf{P}_{1:i}$, and $\mathcal{T}_{1:i}^\perp$ the orthogonal
 602 anti-projector to $\mathbf{T}_{1:i}$. Due to Property 4:

$$\mathcal{P}_{1:i}^\perp = \mathbf{I}_P - \sum_{k=1}^{k=i} \Sigma \mathbf{p}_k (\mathbf{p}'_k \Sigma \mathbf{p}_k)^{-1} \mathbf{p}'_k$$

603 The matrix $\mathbf{X}_{1:i}$ can be written by means of two expressions. From Equa-
 604 tion 6, it is obvious that $\mathbf{X}_{1:i} = \mathcal{T}_{1:i}^\perp \mathbf{X}$. Using the values of \mathbf{t}_i from Equa-
 605 tion A.4, it is possible to substitute the \mathbf{t}_i into Equation 6. Thus, a new
 606 expression of $\mathbf{X}_{1:i}$ is deduced:

$$\begin{aligned} \mathbf{X}_{1:i} &= \mathbf{X}_{1:i-1} - \mathbf{X} \Sigma \mathbf{p}_i (\mathbf{p}'_i \Sigma \mathbf{p}_i)^{-1} \mathbf{p}'_i \\ &= \mathbf{X} - \sum_{k=1}^{k=i} \mathbf{X} \Sigma \mathbf{p}_k (\mathbf{p}'_k \Sigma \mathbf{p}_k)^{-1} \mathbf{p}'_k \\ &= \mathbf{X} \mathcal{P}_{1:i}^\perp \end{aligned}$$

607 Finally:

$$\mathbf{X}_{1:i} = \mathcal{T}_{1:i}^\perp \mathbf{X} = \mathbf{X} \mathcal{P}_{1:i}^\perp \quad (\text{A.7})$$

608 The anti-projection of \mathbf{X} orthogonally to \mathbf{T} into \mathbb{R}^N gives the same result
 609 as its oblique Σ anti-projection to \mathbf{P} into \mathbb{R}^P .

610 **Appendix B. New calculation of the \mathbf{p}_i into \mathbb{R}^P from standard PLSR**

611 The deflation of \mathbf{y} is not necessary when \mathbf{X} is deflated [25], so equation 1
 612 can first be simplified and then written using $\mathcal{T}_{1:i}^\perp$:

$$\begin{aligned}\mathbf{w}_{i+1} &= \mathbf{X}'_{1:i}\mathbf{y} \\ &= \mathbf{X}'\mathcal{T}_{1:i}^\perp\mathbf{y}\end{aligned}$$

613 The combination of Equations 1, 3 and 5 from standard PLSR plus Equation
 614 A.7 leads to:

$$\mathbf{p}_{i+1} = \alpha_{i+1}\mathbf{X}'\mathcal{T}_{1:i}^\perp\mathcal{T}_{1:i}^\perp\mathbf{X}\mathbf{X}'\mathcal{T}_{1:i}^\perp\mathbf{y} \quad (\text{B.1})$$

$$= \alpha_{i+1}\mathbf{X}'\mathcal{T}_{1:i}^\perp\mathbf{X}\mathbf{X}'\mathcal{T}_{1:i}^\perp\mathbf{y} \quad (\text{B.2})$$

$$= \alpha_{i+1}\mathcal{P}'_{1:i}^\perp\mathbf{X}'\mathbf{X}\mathcal{P}'_{1:i}^\perp\mathbf{X}'\mathbf{y} \quad (\text{B.3})$$

615 with α_{i+1} a nonzero scalar associated to \mathbf{p}_{i+1} . The value of α_{i+1} is not
 616 important because it is simplified into Equation 10, but vectors \mathbf{p}_i should have
 617 small norms to improve the stability of the calculation. For this reason, the
 618 \mathbf{p}_i are Σ -normalized to 1 such that they form a Σ -orthonormal basis of the
 619 subspace of \mathbb{R}^P containing the useful information. The new algorithm for the
 620 calculation of the \mathbf{p}_i is thus written:

- 621 • Step 1:

$$\mathbf{p}_1 = \mathbf{X}'\mathbf{X}\mathbf{X}'\mathbf{y}$$

$$\mathbf{p}_1 = \mathbf{p}_1(\mathbf{p}'_1\Sigma\mathbf{p}_1)^{-0.5}$$

622 • Step $i + 1$:

$$\mathbf{p}_{i+1} = \mathcal{P}'_{1:i} \mathbf{X}' \mathbf{X} \mathcal{P}'_{1:i} \mathbf{X}' \mathbf{y}$$

$$\mathbf{p}_{i+1} = \mathbf{p}_{i+1} (\mathbf{p}'_{i+1} \mathbf{\Sigma} \mathbf{p}_{i+1})^{-0.5}$$

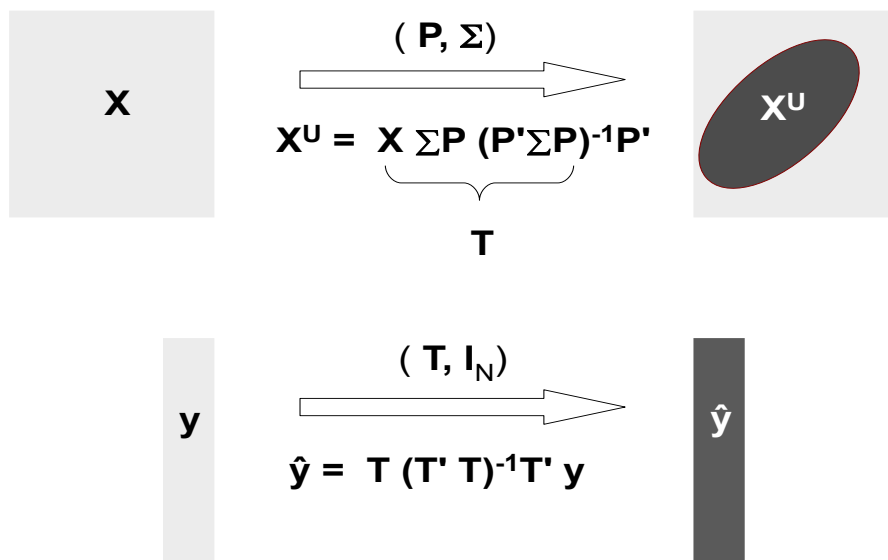


Figure B.1: Standard PLSR as a double projection. The upper panel shows an oblique Σ -projection of \mathbf{X} onto \mathbf{P} yielding \mathbf{X}^U and its scores \mathbf{T} . The lower panel shows an orthogonal projection of \mathbf{y} onto \mathbf{T} yielding $\hat{\mathbf{y}}$.

X	matrix ($N \times P$), N observations and P spectral variables
y	vector ($N \times 1$), the response variable
X_{1:i}	anti-projection of X orthogonally to $\{\mathbf{t}_1, \mathbf{t}_2, \dots, \mathbf{t}_i\}$
y_{1:i}	anti-projection of y orthogonally to $\{\mathbf{t}_1, \mathbf{t}_2, \dots, \mathbf{t}_i\}$
T	matrix ($N \times A$), scores for X
P	matrix ($P \times A$), loadings for X
T_{1:i}, P_{1:i}	matrices containing the i first columns of T and P
W	matrix ($P \times A$), weights for X
c	vector ($A \times 1$), loadings c_1, c_2, \dots, c_A for y
Σ	Moore-Penrose pseudo-inverse of (X' X)
I_N, I_P	identity matrices for \mathbb{R}^N and \mathbb{R}^P spaces
r	vector ($P \times 1$)
T_{1:i}[⊥]	$(N \times N)$ orthogonal anti-projector to $\{\mathbf{t}_1, \mathbf{t}_2, \dots, \mathbf{t}_i\}$
P_{1:i}[⊥]	$(P \times P)$ oblique anti-projector to $\{\mathbf{p}_1, \mathbf{p}_2, \dots, \mathbf{p}_i\}$ with the metric Σ
t_i	score vector at step i of standard PLSR; also i^{th} column vector of T
p_i	loading vector at step i of standard PLSR; also i^{th} column vector of P
w_i	weight vector at step i of standard PLSR; also i^{th} column vector of W
b	regression vector, or vector of b-coefficients, for A latent variables
A	number of latent variables; also dimension of the PLSR model and rank of T , W and P

Table B.1: Main notations

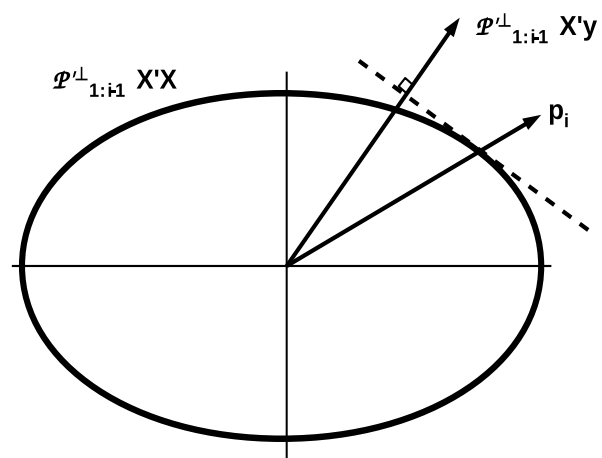


Figure B.2: Geometric building of the loadings \mathbf{p}_i using standard PLSR (from Figure 13 of Phatak [4])

Model	Data	Reference values	\mathbf{r}
m_{pls}	\mathbf{X}	\mathbf{y}	$\mathbf{X}'\mathbf{y}$
m_{nas}	\mathbf{X}	\mathbf{y}	\mathbf{s}_{nas}
m_{y^2}	\mathbf{X}	\mathbf{y}	$\mathbf{X}'\mathbf{y}^2$
$m_{exp(y)}$	\mathbf{X}	\mathbf{y}	$\mathbf{X}'exp(\mathbf{y})$

Table B.2: The four models obtained using VODKA regression with different values for \mathbf{r} .

Latent variables	<i>RMSECV</i>				<i>RMSEP</i>			
	m_{pls}	m_{nas}	m_{y^2}	$m_{exp(y)}$	m_{pls}	m_{nas}	m_{y^2}	$m_{exp(y)}$
4	0.96	1.07	0.68	1.36	1.76	1.99	1.04	2.59
5	0.89	0.99	0.61	1.25	1.50	1.90	0.92	2.31
6	0.73	0.90	0.62	1.22	1.20	1.76	0.91	2.17
7	0.68	0.72	0.61	1.12	1.06	0.87	0.90	2.22
8	0.67	0.68	0.58	1.02	0.99	0.89	0.96	2.19
9	0.64	0.64	0.57	0.72	0.93	0.92	0.96	1.11
10	0.56	0.56	0.56	0.63	0.92	0.92	0.94	0.90
11	0.56	0.56	0.56	0.61	0.95	0.94	0.94	0.89
12	0.56	0.56	0.55	0.57	1.02	0.99	0.96	0.93

Table B.3: Application 1: The standard errors of cross-validation (*RMSECV*) and prediction (*RMSEP*) for the four models and a range including the lower latent variables. *RMSEP* values less than or equal to 0.92 are represented in bold.

Latent variables	m_{plsr}	m_{nas}	m_{y^2}	$m_{exp(y)}$
4	52.8	137.7	51.6	50.7
5	56.3	114.3	54.2	50.1
6	63.7	105.7	57.9	50.9
7	68.8	117.8	81.9	51.3
8	70.8	190.2	94.3	64.7
9	76.7	189.4	100.1	144.5
10	98.4	204.9	98.8	99.3
11	116.0	147.9	110.8	95.2
12	152.2	206.8	142.4	139.9

Table B.4: Application 1: The norms of the regression vectors \mathbf{b} for models m_{plsr} , m_{nas} , m_{y^2} and $m_{exp(y)}$ of Application 1. Values corresponding to the selected number of latent variables are shown in bold

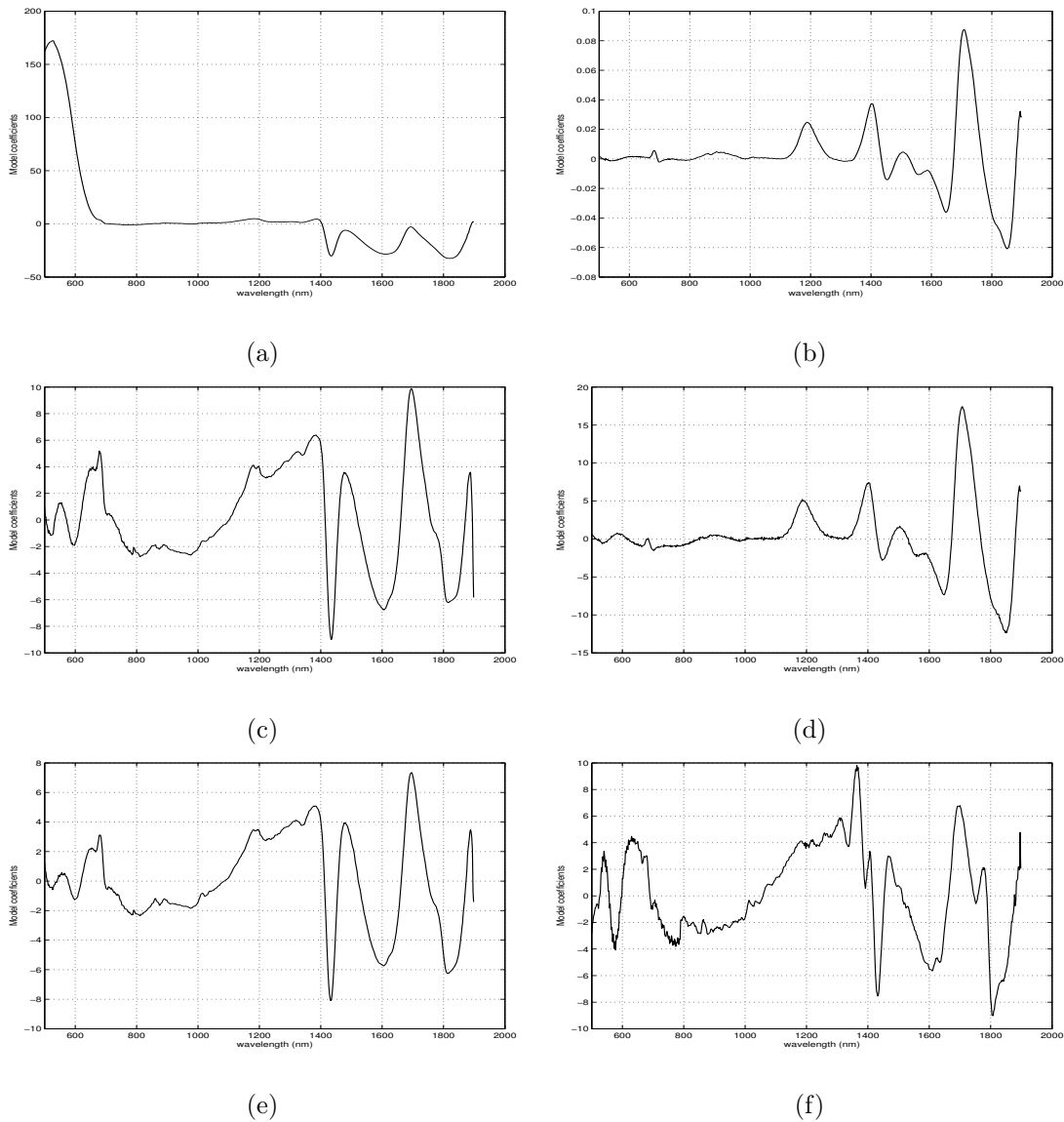


Figure B.3: Spectra and b-coefficient vectors for Application 1. Spectra of $\mathbf{X}'\mathbf{y}$ (a) and \mathbf{s}_{nas} (b). The coefficient vectors for models m_{plsr} (c), m_{nas} (d), m_{y2} (e) and $m_{exp(y)}$ (f) were calculated with 10, 7, 7 and 11 latent variables, respectively.

Latent variables	<i>RMSECV</i>				<i>RMSEP</i>			
	m_{pls}	m_{nas}	m_{y^2}	$m_{exp(y)}$	m_{pls}	m_{nas}	m_{y^2}	$m_{exp(y)}$
6	0.076	0.059	0.056	0.069	0.073	0.15	0.17	0.069
7	0.083	0.060	0.048	0.051	0.058	0.19	0.058	0.050
8	0.053	0.059	0.051	0.062	0.076	0.15	0.058	0.050
9	0.049	0.049	0.048	0.041	0.079	0.12	0.051	0.049
10	0.042	0.047	0.047	0.046	0.090	0.089	0.059	0.047
11	0.039	0.040	0.041	0.038	0.088	0.084	0.059	0.051
12	0.039	0.042	0.042	0.043	0.096	0.071	0.054	0.050
13	0.037	0.039	0.042	0.043	0.095	0.061	0.053	0.050
14	0.035	0.039	0.035	0.036	0.096	0.060	0.053	0.045
15	0.029	0.041	0.035	0.033	0.095	0.059	0.053	0.046
16	0.029	0.033	0.031	0.030	0.094	0.059	0.050	0.049
17	0.029	0.036	0.029	0.029	0.094	0.056	0.052	0.056
18	0.027	0.028	0.029	0.028	0.094	0.054	0.054	0.059
19	0.027	0.030	0.028	0.027	0.094	0.057	0.065	0.058

Table B.5: The standard errors of cross-validation (*RMSECV*) and prediction (*RMSEP*) for the four models and a range including the lower latent variables of Application 2. *RMSEP* values less than or equal to 0.058 are represented in bold.

Latent variables	m_{pls_r}	m_{nas}	m_{y^2}	$m_{exp(y)}$
6	1.44	1.31	1.16	1.15
7	1.52	1.14	1.17	1.15
8	1.61	1.21	1.22	1.19
9	1.71	1.23	1.25	1.20
10	1.77	1.21	1.24	1.25
11	1.8	1.21	1.28	1.24
12	1.85	1.25	1.30	1.24
13	1.86	1.28	1.31	1.32
14	1.86	1.34	1.30	1.36
15	1.86	1.46	1.31	1.34
16	1.86	1.42	1.39	1.50
17	1.86	1.42	1.41	1.52
18	1.86	1.42	1.42	1.52
19	1.86	1.44	1.41	1.53

Table B.6: The norms of the regression vectors \mathbf{b} for models m_{pls_r} , m_{nas} , m_{y^2} and $m_{exp(y)}$ of Application 2. Values corresponding to the selected number of latent variables are shown in bold.

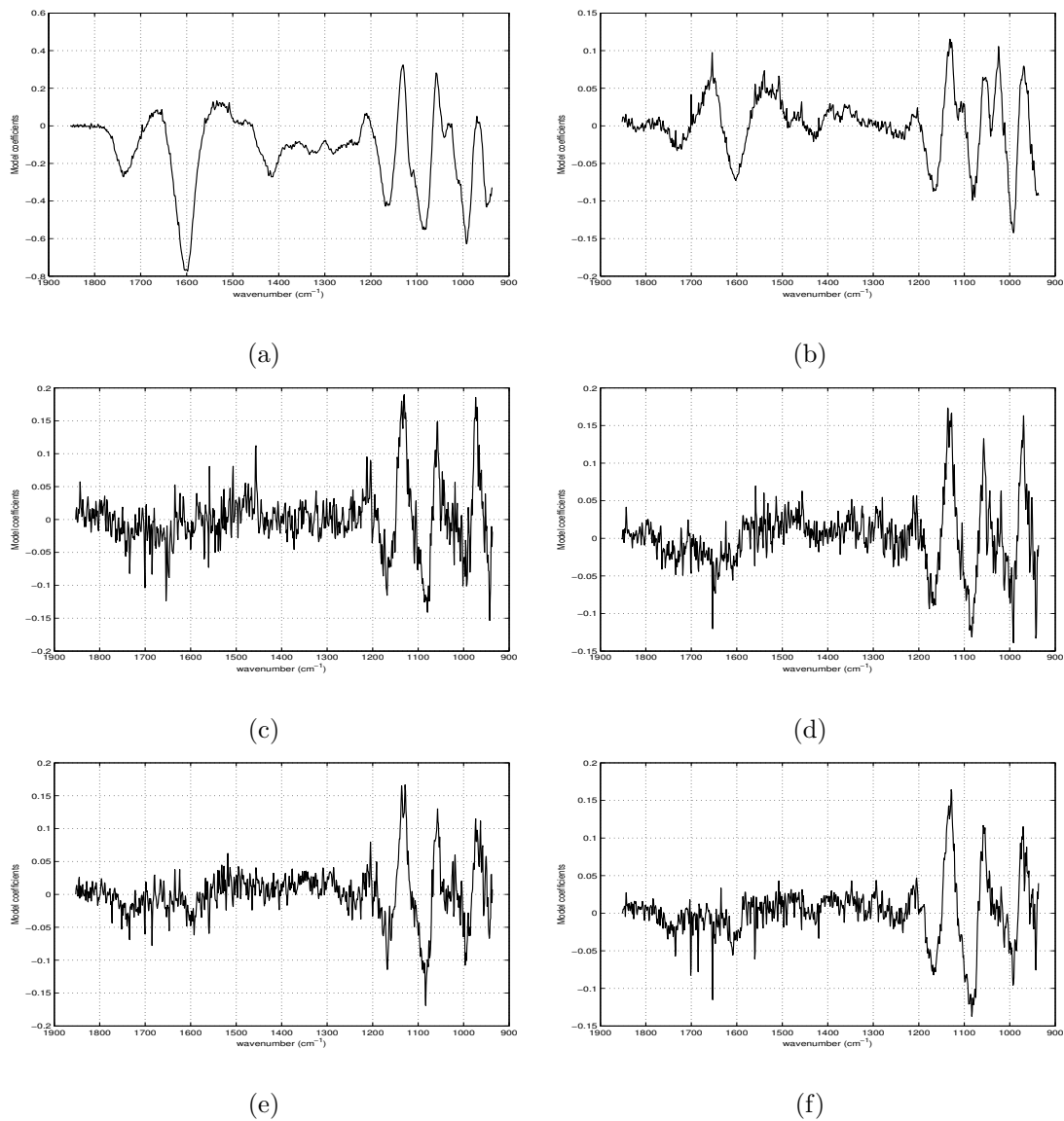


Figure B.4: Spectra and b-coefficient vectors for Application 2. Spectra of $\mathbf{X}'\mathbf{y}$ (a) and \mathbf{s}_{nas} (b). The b-coefficient vectors for models m_{plsr} (c), m_{nas} (d), m_{y2} (e) and $m_{exp(y)}$ (f) were calculated with 7, 18, 9 and 10 latent variables, respectively.