

# PLS regression : a directional signal-noise approach

Pierre Druilhet <sup>a,\*</sup> Alain Mom <sup>b</sup>

<sup>a</sup>*CREST-ENSAI, Campus de Ker Lann, 35 170 BRUZ, France*

<sup>b</sup>*Laboratoire de statistique, Université Rennes II, RENNES, France*

---

## Abstract

We present a new approach of univariate partial least squares regression (PLSR) based on directional signal-noise ratios. We show how PLSR, unlike PCR, takes into account the actual value and not only the variance of the OLS estimator. We find an orthogonal sequence of directions associated with decreasing signal-noise ratios. Then, we state PLS estimators as least squares estimators constrained to be null onto the last directions. We also give another procedure that shows how PLSR rebuilds the OLS estimator iteratively by seeking at each step the direction with largest difference of signals over the noise. The latter approach does not involve any arbitrary scale or orthogonality constraints.

*Key words:* Biased regression, Constrained least squares, Latent variables, Partial least squares, Principal component, Shrinkage

*1991 MSC:* 62J05, 62J07

---

## 1 Introduction

Many different regression methods are now available to improve ordinary least squares (OLS) estimators in multiple linear regression when explanatory variables are strongly correlated. Like principal component regression (PCR), partial least squares regression (PLSR) is based on the construction of latent variables, i.e. new variables that are linear combinations of the original explanatory variables. In PCR, these latent variables have a statistical meaning: they explain the internal covariance structure of the explanatory variables.

---

\* Corresponding author

*Email addresses:* [druilhet@ensai.fr](mailto:druilhet@ensai.fr) (Pierre Druilhet), [alain.mom@uhb.fr](mailto:alain.mom@uhb.fr) (Alain Mom).

In PLSR, the statistical interpretation of the latent variables is not always so clear, but intuitively, they make a compromise between the internal structure of the explanatory variables and their relationship with the response. Continuum regression (Stone and Brooks, 1990) gives a more general latent variables method that includes PLSR and PCR. Latent variables approach is often related to an underlying common joint covariance distribution structure between the explanatory variables and the response for each individual: Helland (1990, 1992), Helland and Almøy (1994) give a population model for PLSR and establish asymptotic properties. However, in some situations, the sampling scheme used to collect data can be complex or unknown; on the other hand, the aim of the study can be to predict new individuals with the same conditional distribution but different joint distributions. In these cases, it is of interest to formulate a conditional model and then to consider the explanatory variables as fixed. See Breiman and Spector (1992) for further discussions and simulation studies upon the difference between fixed and random cases. In this paper, we assume a conditional model and focus on the improvement of the least squares estimator.

Like PCR and Ridge regression, PLSR provides a shrunk version of the OLS estimator (De Jong, 1995 and Goutis, 1996). For a comparison of these three techniques by simulations and case studies, see Frank and Friedman (1993); for a general discussion, see Brown (1993). When the explanatory variables are strongly correlated, the improvement of OLS estimator obtained by shrinking can be quite large w.r.t. the mean-squared error (MSE). However the MSE is a global measure of how the estimate is far from the true parameter and does not show how the estimate is improved on specific directions. Butler and Denham (2000) and Lingjærde and Christophersen (2000) study the shrinkage properties of PLS estimators onto directions given by the singular values decomposition of the design matrix.

In this paper, we aim to show that PLSR singles out some directions to shrink OLS estimators. More precisely, while PCR shrinks the OLS estimator onto directions with large variance (under a scale constraint), we show how PLSR do the same onto some directions corresponding to small signal-noise ratios. This signal-noise ratio is directly related with optimal directional shrinkage factors. We also show how the sequence of PLS estimators rebuilds successively the OLS estimator by seeking the direction of maximal difference of signals over the noise.

In section 2, we give a brief review of the classical constrained estimators. In section 3, we recall some known results on the link between shrinkage coefficients and signal-noise ratios for a one-dimensional parameter. Then, we apply these results to linear models. In section 4, we use a maximization procedure to obtain a sequence of orthogonal directions with decreasingly ordered signal-noise ratios. Then, we obtain PLS estimators as OLS estimators

under the constraint of nullity onto the last directions. In section 5, we propose a constraint free algorithm based on the difference of signals over the noise. This algorithm gives iteratively the PLS estimators.

## 2 Multiple linear model and biased estimation

We assume the following model for the  $n$  individuals:

$$\dot{y}_i = \mu + \dot{x}'_i \beta + \varepsilon_i, \quad i = 1, \dots, n, \quad (1)$$

where  $\dot{y}_i$  is the real response observed on the  $i^{\text{th}}$  individual,  $\dot{x}_i$  are the  $p$ -vector of its explanatory variables considered as fixed,  $\beta$  is the unknown  $p$ -vector of parameters, and  $\varepsilon_i$  are i.i.d. mean zero variance  $\sigma^2$  real random variables. Removing the dots on the original data means that their sample averages, denoted by  $\bar{y}$  and  $\bar{x}$ , has been removed. In vector notation, we have for the centered model:

$$y = X\beta + \varepsilon. \quad (2)$$

To avoid handling complicated generalized inverses, the design matrix  $X$  is assumed to be of full-rank. The non full-rank case will be discussed in section 7.

We write  $s = X'y$  and  $S = X'X$ . We denote by  $I_n$  the  $(n, n)$  identity matrix. Let  $A$  be a  $(r, s)$  matrix, and  $M$  be a  $(r, r)$  symmetric positive matrix, we denote  $P_A^M = A(A'MA)^{-1}A'M$  the projector matrix onto  $\text{Range}(A)$  w.r.t.  $M$ . We also write  $P_A = P_A^{I_r}$ .

Consider a new individual  $\dot{x}_o$  and assume that its response  $\dot{y}_o$  is given by

$$\dot{y}_o = \mu + \dot{x}'_o \beta + \varepsilon_o \quad (3)$$

where  $\varepsilon_o$  is independent of  $(\varepsilon_1, \dots, \varepsilon_n)$  with mean 0 and variance  $\sigma_o^2$ . We want to predict  $\dot{y}_o$  by  $\bar{y} + \dot{x}'_o \hat{\beta}$ , where  $x_o = \dot{x}_o - \bar{x}_o$  and  $\hat{\beta}$  is an estimate of  $\beta$ . It is well known that the best linear unbiased estimator of  $x'_o \beta$  is  $x'_o \hat{\beta}^{\text{OLS}}$ , where  $\hat{\beta}^{\text{OLS}} = (X'X)^{-1}X'y = S^{-1}s$ . It is also well known that  $\bar{y} + x'_o \hat{\beta}^{\text{OLS}}$  is the best linear unbiased predictor of  $\dot{y}_o$ . However, for some  $x_o$ ,  $x'_o \hat{\beta}^{\text{OLS}}$  gives very poor estimate of  $x'_o \beta$  or predictor of  $\dot{y}_o$ . There are many methods to improve the estimation of  $\beta$  by biased estimators.

For example, Ridge estimators are obtained by approximating  $S$  by  $S + \alpha I$  in the normal equation  $S\hat{\beta} = s$ , where the real  $\alpha$  is a control parameter. They

can also be obtained as constrained least squares estimators:

$$\hat{\beta}_c^{\text{Ridge}} = \underset{\beta: \beta' \beta \leq c}{\text{ArgMin}} \|y - X\beta\|^2 \quad (4)$$

where  $c$  depends of  $\alpha$ .

In a similar way, Lasso estimators (Tibshirani, 1996) are defined by

$$\hat{\beta}_c^{\text{Lasso}} = \underset{\beta: \sum |\beta_i| \leq c}{\text{ArgMin}} \|y - X\beta\|^2. \quad (5)$$

In PCR, it is postulated that bad estimation of  $x'_o \beta$  by OLS is due to the fact that  $\text{var}(x'_o \hat{\beta}^{\text{OLS}})$  can be very "large". To give a sense to "large", a scale constraint must be added: in PCR, the constraint is  $x'_o x_o = 1$ . Denote by  $w_1^{\text{PCR}}, \dots, w_p^{\text{PCR}}$  an orthogonal basis of unit eigenvectors of  $S$  corresponding to decreasingly ordered eigenvalues. Since  $\text{var}(x'_o \hat{\beta}^{\text{OLS}}) = \sigma^2 x'_o S^{-1} x_o$ , PCR forces the OLS to be zero onto the directions given by the last eigenvectors. More precisely, define  $\overline{W}_{(q)}^{\text{PCR}} = (w_{q+1}^{\text{PCR}} | \dots | w_p^{\text{PCR}})$ , then

$$\beta_{(q)}^{\text{PCR}} = \underset{\beta: \overline{W}_{(q)}^{\text{PCR}' } \beta = 0}{\text{ArgMin}} \|y - X\beta\|^2. \quad (6)$$

This expression shows how the OLS estimator is modified in a statistically meaningful way. On another way,  $\beta_{(q)}^{\text{PCR}}$  can be obtained after a regression on the latent variables  $t_1^{\text{PCR}} = X w_1^{\text{PCR}}$ ,  $t_2^{\text{PCR}} = X w_2^{\text{PCR}}$ ,  $\dots$ ,  $t_q^{\text{PCR}} = X w_q^{\text{PCR}}$ . In that case, the vectors  $w_1^{\text{PCR}}, w_2^{\text{PCR}}, \dots, w_q^{\text{PCR}}$  are weight vectors. The fact that the same vector  $w_i$  represent two different things, either a direction or a weight vector, comes from the fact that  $w_1^{\text{PCR}}, w_2^{\text{PCR}}, \dots, w_p^{\text{PCR}}$  are orthogonal (see Lemma 5 and Formula 10). The main criticism of PCR is that the constraints  $\overline{W}_{(q)}^{\text{PCR}' } \beta = 0$  are obtained by considering  $\text{var}(x'_o \hat{\beta}^{\text{OLS}})$  (the noise) but not the actual value of  $x'_o \hat{\beta}^{\text{OLS}}$  (the signal). In other words, PCR constructs the latent variables independently of their relationship with the response  $y$ .

In the usual approach, PLSR is a latent variables method that aims to overcome this undesirable feature of PCR. The first latent variables  $t_1^{\text{PLS}} = X w_1^{\text{PLS}}$  maximizes the empirical covariance:

$$\text{cov}_e(t, y) = \text{cov}_e(X w, y) = w' X y \quad (7)$$

under the scale constraint  $w'w = 1$ . The  $i^{\text{th}}$  latent variables  $t_i = X w_i^{\text{PLS}}$  maximizes (7) under the constraints  $w'w = 1$  and  $t_i \perp (t_1^{\text{PLS}}, \dots, t_{i-1}^{\text{PLS}})$ . Define  $T_{(q)}^{\text{PLS}} = (t_1^{\text{PLS}} | \dots | t_q^{\text{PLS}})$  and  $W_{(q)}^{\text{PLS}} = (w_1^{\text{PLS}} | \dots | w_q^{\text{PLS}})$ . The PLS estimator  $\hat{\beta}_{(q)}^{\text{PLS}}$  based on the  $q$  first latent variables is:

$$\hat{\beta}_{(q)}^{\text{PLS}} = W_{(q)}^{\text{PLS}} \left( T_{(q)}^{\text{PLS}' } T_{(q)}^{\text{PLS}} \right)^{-1} T_{(q)}^{\text{PLS}' } y, \quad (8)$$

$$= W_{(q)}^{\text{PLS}} \left( W_{(q)}^{\text{PLS}' } S W_{(q)}^{\text{PLS}} \right)^{-1} W_{(q)}^{\text{PLS}' } S \hat{\beta}^{\text{OLS}}, \quad (9)$$

$$= P_{W_{(q)}^{\text{PLS}}}^S \hat{\beta}^{\text{OLS}}. \quad (10)$$

Note that equations 8, 9 and 10 are general properties of a latent variables regressors.

Helland (1988) shows that  $\text{Range} \left( W_{(q)} \right) = K_q$  where  $K_q$  is the Krylov subspace defined by:

$$K_q = \text{span} \left( \mathbf{s}, S \mathbf{s}, S^2 \mathbf{s}, \dots, S^{q-1} \mathbf{s} \right). \quad (11)$$

He also shows that Martens' algorithm (see Martens and Næs, 1989) is equivalent to a Gram-Schmidt orthogonalization procedure of the nested Krylov subspaces  $(K_1, K_2, \dots)$  and gives the same estimator, as easily seen with Formula 10. Note that Formulas 8, 9 and 10 are general expressions of a latent variables estimator.

The goal of the paper is to find an analogue to Formula 6 for PLSR. Of course, any estimator obtained by latent variables, whatever they are, can be expressed as a least squares estimator under some linear constraints. However, we seek for PLSR constraints corresponding to direction that have a clear statistical meaning. Since PLSR is a shrinkage method, shrinkage factors onto directions are of importance. This will be the key point of our approach.

### 3 Shrinkage factors and signal-noise ratio

The signal-noise ratio arises naturally in univariate estimation to deal with shrinkage. First, we recall some known facts on unidimensional shrinkage. Then, we apply these notions to directions in linear model.

#### 3.1 Shrinking a one-dimensional estimator or predictor

Let  $\hat{\theta}$  be an unbiased estimator of a real parameter  $\theta$ . We denote by  $\sigma^2$  the variance of  $\hat{\theta}$ . We want to improve  $\hat{\theta}$  by shrinking, so we consider new estimators  $a \hat{\theta}$  for  $a \in \mathbb{R}$ . Let  $a^*$  be the scalar that minimizes the MSE, i.e. :

$$a^* = \underset{a \in \mathbb{R}}{\text{ArgMin}} \mathbb{E} \left( a \hat{\theta} - \theta \right)^2. \quad (12)$$

We have,

$$a^* = \frac{\rho^2}{1 + \rho^2} \quad \text{where} \quad \rho = \frac{|\theta|}{\sigma}. \quad (13)$$

When  $\rho$  is known, which arises in some scale models,  $a^* \hat{\theta}$  is uniformly better under the quadratic lost among all the estimators  $a \hat{\theta}$ . When  $\rho$  is unknown but  $\sigma^2$  is known,  $\rho$  can be estimated by

$$\hat{\rho} = \frac{|\hat{\theta}|}{\sigma} \quad \text{and} \quad \hat{a}^* = \frac{\hat{\rho}^2}{1 + \hat{\rho}^2} = \frac{\hat{\theta}^2}{\sigma^2 + \hat{\theta}^2}. \quad (14)$$

We see that the shrinkage factor  $\hat{a}^*$  is an increasing function of the signal-noise ratio  $\hat{\rho}$ . Bibby and Toutenburg (1977) study the properties of the estimator  $\hat{a}^* \hat{\theta}$ . They also show that, in the normal case, the probability that  $\hat{a}^* \hat{\theta}$  improve  $\hat{\theta}$  is very high: 77% or more depending on  $\rho$ . When  $\rho < 1$ , the probability is equal to 1.

Consider now the following prediction problem: we want to predict  $y_o = \theta + \varepsilon_o$  where  $\varepsilon_o \sim (0, \sigma_o^2)$  is independent of  $\theta$ . The predictor  $\tilde{y}_o = \hat{\theta}$  is an unbiased predictor of  $y_o$ , i.e.  $\mathbb{E}(\hat{\theta}) = \mathbb{E}(y_o) = \theta$ , and we want to improve the prediction by shrinking. The equivalent of the MSE in prediction is the mean-squared error of prediction (MSEP):

$$\text{MSEP} = \mathbb{E}(\tilde{y}_o - y_o)^2 = \text{MSE} + \sigma_o^2.$$

The predictor  $a \tilde{y}_o$  that minimizes the MSEP is  $a^* \tilde{y}_o$  with  $a^*$  given by equation 13. Thus in both estimation and prediction problem, the best shrinkage factor is the same.

### 3.2 Shrinkage factors in multiple linear regression

We apply the results of Section 3.1 to the linear model (2). For some  $\mathbf{x}_o \in \mathbb{R}^p$ , we are interested in the estimation of  $\mathbf{x}'_o \beta$ . The minimum-variance unbiased linear estimator of  $\mathbf{x}'_o \beta$  is  $\mathbf{x}'_o \hat{\beta}^{\text{OLS}}$ . To improve this estimator, we apply the shrinkage factor presented in Section 3.1. So, we have:

$$a_{\mathbf{x}_o}^* = \frac{\rho_{\mathbf{x}_o}^2}{1 + \rho_{\mathbf{x}_o}^2} \quad \text{where} \quad \rho_{\mathbf{x}_o} = \frac{|\mathbf{x}'_o \beta|}{\sigma \sqrt{\mathbf{x}'_o S^{-1} \mathbf{x}_o}}, \quad (15)$$

that can be estimated by

$$\hat{a}_{x_o}^* = \frac{\hat{\rho}_{x_o}^2}{1 + \hat{\rho}_{x_o}^2} \quad \text{where} \quad \hat{\rho}_{x_o} = \frac{|\mathbf{x}'_o \hat{\beta}^{\text{OLS}}|}{\sigma \sqrt{\mathbf{x}'_o S^{-1} \mathbf{x}_o}}. \quad (16)$$

If  $\sigma^2$  is not known, it can be replaced by an estimation. However,  $\sigma^2$  is a positive constant and will not play any role in the following. Note that  $\mathbf{x} \mapsto \hat{\rho}_{\mathbf{x}}$  is positively homogeneous, that is:

$$\forall \alpha \neq 0, \quad \hat{\rho}_{\alpha \mathbf{x}} = \hat{\rho}_{\mathbf{x}}. \quad (17)$$

Consequently, the signal-noise ratio  $\hat{\rho}_{x_o}$  depends only on the direction given by  $\mathbf{x}_o$ , not on the exact  $\mathbf{x}_o$ . More precisely:

$$\hat{\rho}_{x_o} = \frac{|\mathbf{x}'_o \hat{\beta}|}{\sigma \sqrt{\mathbf{x}'_o S^{-1} \mathbf{x}_o}} = \frac{|\mathbf{x}'_o S^{-1} \mathbf{s}|}{\sigma \sqrt{\mathbf{x}'_o S^{-1} \mathbf{x}_o}} = \frac{\sqrt{\mathbf{s}' S^{-1} \mathbf{s}}}{\sigma} |\cos_{S^{-1}}(\widehat{\mathbf{x}_o}, \mathbf{s})|, \quad (18)$$

where  $\cos_{S^{-1}}(\widehat{\mathbf{x}_o}, \mathbf{s})$  is the cosine of the angle  $(\widehat{\mathbf{x}_o}, \mathbf{s})$  w.r.t. the quadratic form  $S^{-1}$ . Thus,  $\hat{\rho}_{x_o}$  depends only on the angle between  $\mathbf{x}_o$  and  $\mathbf{s} = X' \mathbf{y}$  w.r.t. the quadratic form  $S^{-1}$ . The greater the angle w.r.t  $S^{-1}$  between  $\mathbf{x}$  and  $\mathbf{s}$ , the less its signal ratio. Thus, the space of null signal-noise, which coincides with the space of null signal, is the orthogonal complement *w.r.t* the quadratic form  $S^{-1}$  of  $\mathbf{s}$ . Moreover, the relationship between  $a_{x_o}^*$  and  $\mathbf{x}_o$  is intrinsically non linear and depends only on  $\rho_{x_o}$ .

## 4 PLSR and directional signal-noise algorithm

In the classical algorithms used to construct PLS estimators, two kinds of arbitrariness are involved: a scale constraint and an orthogonality constraint. In this section, we present an algorithm based on the signal-noise ratio. Since the signal-noise ratio is scale invariant, no scale constraint is necessary. The algorithm seeks a sequence of orthogonal directions corresponding to decreasing signal-noise ratios. Then, we show that PLS estimators as least squares estimators constrained to be null onto the last directions of this sequence.

### 4.1 The signal-noise maximization procedure

Let start to seek the direction  $\mathbf{w}_1$  that maximizes the signal-noise ratio, i.e.

$$w_1 = \underset{w \in \mathbb{R}^p}{\text{ArgMax}} (\hat{\rho}_w) = \underset{w \in \mathbb{R}^p}{\text{ArgMax}} \frac{|w' \hat{\beta}|}{\sigma \sqrt{w' S^{-1} w}}. \quad (19)$$

By a direct application of Cauchy-Schwarz's inequality and because  $w' \hat{\beta} = w' S^{-1} s$ , we have:

**Proposition 1** *The first direction is spanned by  $s$ , that is  $w_1$  is any vector proportional to  $s$ .*

It is interesting to see that the direction of maximal noise play a central role in PLSR.

Iteratively, at step  $i$ , the direction  $w_i$  maximizes the signal-noise ratio  $\hat{\rho}_w$  under the orthogonality constraints  $w \perp (w_1, \dots, w_{i-1})$ .

Define  $q^*$  the lowest integer  $q$  such that  $K_{q+1} = K_q$ , where  $K_q$  is the Krylov subspace given by (11). We have:

$$S K_{q^*} = K_{q^*}. \quad (20)$$

This invariance property implies that  $\hat{\beta}^{\text{OLS}}$  belongs to  $K_{q^*}$ .

**Lemma 2** *The orthogonal subspace of  $K_{q^*}$  w.r.t. the quadratic form  $S$  is the same as the orthogonal subspace w.r.t.  $I$ . Moreover,  $K_{q^*}^\perp$  is included in the null-signal space, that is*

$$\forall w \in K_{q^*}^\perp, w' \hat{\beta}^{\text{OLS}} = 0.$$

**PROOF.** Let  $w$  be in the orthogonal to  $K_{q^*}$  w.r.t.  $I$  and  $w^*$  be any vector in  $K_{q^*}$ , then  $w' S w^* = w' (S w^*) = 0$  since  $S w^* \in K_{q^*+1} = K_{q^*}$ . Hence  $w$  lies in the orthogonal to  $K_{q^*}$  w.r.t.  $S$ . Consider now  $w \in K_{q^*}^\perp$ . Then,  $\forall w^* \in K_{q^*}$ ,  $w' w^* = 0$ . Equation 20 gives  $s = S w^{**}$  for some  $w^{**} \in K_{q^*}$  and then  $w' \hat{\beta}^{\text{OLS}} = w' S^{-1} s = w' S^{-1} S w^{**} = w' w^{**} = 0$ .  $\square$

We can now formulate the theorem that characterizes the sequence  $w_1, \dots, w_p$  and shows its relationship with the Krylov subspaces  $K_1, \dots, K_{q^*}$ .

**Theorem 3** *For  $i = 1, \dots, p$ , the directions  $w_i$  can be obtained as follow:*

- (1) *If  $i \leq q^*$ , then  $w_i$  belongs to  $K_i$  and is orthogonal to  $K_{i-1}$ . Moreover,  $w_i' \hat{\beta}^{\text{OLS}} \neq 0$ .*
- (2) *If  $i > q^*$ , then  $w_i$  is any vector orthogonal to  $(w_1, \dots, w_{i-1})$ . In that case,  $w_i' \hat{\beta}^{\text{OLS}} = 0$ , i.e.  $w_i$  belongs to the null-signal space.*



**PROOF.**

We use a recurrence procedure. It has already been shown that  $w_1 = s$ . We assume that the theorem is true for  $k = 1, \dots, i - 1$ . At step  $i$ , we seek  $w_i$  that maximizes the signal-noise ratio and that is orthogonal to  $w_1, \dots, w_{i-1}$ . By Formula 17, we actually seek a direction. Thus, without loss of generality, we may impose a scale constraint. The choice of the constraint will change the actual  $w_i$  but not the direction given by  $w_i$ . So, we choose the simplest constraint which is  $w' S^{-1} w = 1$ . Under this constraint, the signal-noise ratio reduce to  $|w' \hat{\beta}^{\text{OLS}}|$ . If  $w$  maximizes  $|w' \hat{\beta}^{\text{OLS}}|$  then either  $w$  either  $-w$  maximize  $w' \hat{\beta}^{\text{OLS}} a$ . So, we want to maximize  $w' \hat{\beta}^{\text{OLS}}$  under the constraint  $w' S^{-1} w = 1$  and  $w' w_k = 0$  for  $k = 1, \dots, i - 1$ .

**Case 1:**  $i \leq q^*$

The maximizing  $w_i$  will be a solution to the Lagrange multiplier equation

$$S^{-1} s + \lambda S^{-1} w_i + \sum_{k=1}^{i-1} \mu_k w_k = 0. \quad (21)$$

By recurrency,  $w_k \in K_k$  for  $k = 1, \dots, i - 1$ , thus  $S w_k \in K_{k+1} \subseteq K_i$ . Multiplying Equation 21 by  $S$  gives:

$$\lambda w_i = -s - \sum_{k=1}^{i-1} \mu_k S w_k. \quad (22)$$

Since  $k < q^*$ ,  $K_{k+1} \neq K_k$  and  $w_1, S w_1, \dots, S w_{i-1}$  are linearly independent. Hence,  $\lambda$  cannot be equal to 0 and  $w_i$  belongs necessarily to  $K_i$ . Point 1 is established.

Multiplying Equation 21 by  $w'_i$  gives

$$w'_i S^{-1} s + \lambda w'_i S^{-1} w_i = 0. \quad (23)$$

Thus,  $w'_i \hat{\beta}^{\text{OLS}} \neq 0$ .

**Case 2:**  $i > q^*$

Since  $w_i$  is orthogonal to  $K_{q^*} = \text{span}\{w_1, \dots, w_{q^*}\}$ , the signal  $w'_i \hat{\beta}^{\text{OLS}} = 0$  according to Lemma 2. Hence, the function to be maximized is null and every  $w_i$  in the orthogonal to  $w_1, \dots, w_{i-1}$  is a solution of the maximization problem. Point 2 is established.  $\square$

Theorem 3 shows that the sequence  $(w_1, \dots, w_{q^*})$  can be found by a Gram-Schmidt orthogonalization procedure from the nested Krylov subspaces  $K_1, \dots, K_{q^*}$ .

Therefore, it corresponds to Martens' weight vectors sequence. However, in our context,  $w_i$  corresponds to a direction for the parameter  $\beta$  or its estimate rather than weight vectors used to define latent variables. Denote

$$W_{(q)} = (w_1 | \dots | w_q) \text{ and } \overline{W}_{(q)} = (w_{q+1} | \dots | w_p). \quad (24)$$

We obviously have:

$$\text{Range}(\overline{W}_{(q)}) = \text{Range}(W_{(q)})^\perp. \quad (25)$$

Therefore, signal-noise ratios and shrinkage factors of directions in  $\text{Range}(\overline{W}_{(q)})$  are bounded above, i.e.

**Proposition 4**

$$\forall x \in \text{Range}(\overline{W}_{(q)}), \quad \hat{\rho}_x \leq \hat{\rho}_{w_{q+1}} \text{ and } \hat{a}_x^* \leq \hat{a}_{w_{q+1}}^*. \quad (26)$$

*4.2 PLS estimators as constrained least squares estimators*

We show that PLS estimators are least squares estimators constrained to be null onto the last directions of the sequence exhibited in Section 4.1. By Formula 26, we know that these directions are associated to the lowest shrinkage factors or the lowest signal-noise ratios of the sequence.

**Lemma 5** *Let  $G$  be a  $p \times g$  matrix of rank  $g$  ( $g \leq p$ ) and  $F$  a  $p \times f$  matrix ( $f = p - g$ ) of rank  $f$  such that  $G'F = 0$ , then:*

$$\text{ArgMin}_{\beta: G'\beta=0} \|y - X\beta\|^2 = S^{-1} (I - P_G^{S^{-1}}) (s) = P_F^S \hat{\beta}.$$

**PROOF.** The proof of this Lemma can be found in Seber (1984).

**Theorem 6** *For  $q \leq q^*$ , the PLS estimator  $\hat{\beta}_{(q)}^{\text{PLS}}$  given by (8) is the least squares estimator constrained to have null signal onto the space  $\overline{W}_{(q)}$ :*

$$\hat{\beta}_{(q)}^{\text{PLS}} = \text{ArgMin}_{\beta: \overline{W}_{(q)}'\beta=0} \|y - X\beta\|^2. \quad (27)$$

**PROOF.** By Theorem 3,  $\overline{W}_{(q)}'K_q = 0$ ,  $\text{Rank}(K_q) = q$  and  $\text{Rank}(\overline{W}_{(q)}) = p - q$ , Lemma 5 with  $G = \overline{W}_{(q)}$  and  $F = K_q$  gives the result.  $\square$

Actually, the constraints  $\overline{W}'_{(q)}\beta = 0$  do not need to take into account the vectors  $w_{q^*+1}, \dots, w_p$  since they belong to the null-signal space. More precisely, let  $\widetilde{W}_{(q)} = (w_{q+1}, \dots, w_{q^*})$  for  $q < q^*$  and  $\widetilde{W}_{(q)} = \emptyset$  for  $q \geq q^*$ , then, forcing the least squares estimator to be zero either on  $\overline{W}_{(q)}$  or on  $\widetilde{W}_{(q)}$  both gives  $\widehat{\beta}_{(q)}^{\text{PLS}}$ . First, we give two technical results:

**Lemma 7** *Let  $G_1$  be a  $p \times g_1$  matrix of rank  $g_1$  and  $G_2$  be a  $p \times g_2$  matrix of rank  $g_2$  such that  $G_2'S^{-1}G_1 = 0$ . Let  $G = (G_1|G_2)$ . Then,*

$$G_2'\widehat{\beta} = 0 \implies \underset{\beta:G'\beta=0}{\text{ArgMin}} \|y - X\beta\|^2 = \underset{\beta:G_1'\beta=0}{\text{ArgMin}} \|y - X\beta\|^2.$$

**PROOF.**  $G_2'S^{-1}G_1 = 0$  and  $G_2'\widehat{\beta} = 0$  imply respectively that  $P_G^{S^{-1}} = P_{G_1}^{S^{-1}} + P_{G_2}^{S^{-1}}$  and  $P_{G_2}^{S^{-1}}(s) = 0$ . Whence, from Lemma 5,

$$\begin{aligned} \underset{\beta:G'\beta=0}{\text{ArgMin}} \|y - X\beta\|^2 &= S^{-1} \left( I - P_G^{S^{-1}} \right) (s) \\ &= S^{-1} \left( I - P_{G_1}^{S^{-1}} \right) (s) \\ &= \underset{\beta:G_1'\beta=0}{\text{ArgMin}} \|y - X\beta\|^2. \end{aligned}$$

□

**Corollary 8** *If  $G'\widehat{\beta} = 0$ , then  $P_G^{S^{-1}}(s) = 0$  and*

$$\underset{\beta:G'\beta=0}{\text{ArgMin}} \|y - X\beta\|^2 = \widehat{\beta}^{\text{OLS}}.$$

**Theorem 9** *We have the following characterizations for  $\widehat{\beta}_{(q)}^{\text{PLS}}$ :*

1)  $\forall q < q^*$ ,

$$\widehat{\beta}_{(q)}^{\text{PLS}} = \underset{\beta:\widetilde{W}'_{(q)}\beta=0}{\text{ArgMin}} \|y - X\beta\|^2.$$

2)  $\forall q \geq q^*$ ,

$$\widehat{\beta}_{(q)}^{\text{PLS}} = \underset{\beta:\widetilde{W}'_{(q)}\beta=0}{\text{ArgMin}} \|y - X\beta\|^2 = \widehat{\beta}^{\text{OLS}}.$$

**PROOF.** From Theorem 3, we have  $\overline{W}'_{(q^*)}\widehat{\beta} = 0$ . Moreover, if  $q < q^*$ ,

$$\widetilde{W}_{(q)} \subset K_{q^*} \perp \overline{W}_{(q^*)} \implies K_{q^*} = SK_{q^*} \perp_{S^{-1}} \overline{W}_{(q^*)}.$$

The conditions of Lemma 7 are fulfilled with  $G_1 = \widetilde{W}_{(q)}$ ,  $G_2 = \overline{W}_{(q^*)}$  and  $G = \overline{W}_{(q)}$ . Then, Theorem 6 and Lemma 7 give the first part of the theorem.

If  $q \geq q^*$ ,  $\overline{W}_{(q)} \subset \overline{W}_{(q^*)}$ . Thus,  $\overline{W}_{(q)}' \widehat{\beta} = 0$ . Theorem 6 and Corollary 8 give  $\widehat{\beta}_{(q)}^{\text{PLS}} = \widehat{\beta}_{ols}$ .  $\square$

Thus, when the underlying space of an estimator on latent variables contains  $K_{q^*}$ , this latter estimator coincides with  $\widehat{\beta}^{\text{OLS}}$ . This follows from the fact that any estimator on latent variables can be obtained by projecting  $\widehat{\beta}^{\text{OLS}}$  on the space spanned by its weight vectors. Hence,  $K_{q^*}$  indicates the end of the procedure of shrinkage in the sequence of nested Krylov subspaces.

## 5 A constraint free algorithm for PLSR

In the classical approach of PLSR and the one presented in Section 4, the maximization algorithm involves some arbitrary constraints. In this section, we give another procedure for PLSR with no such arbitrariness. The main idea is to rebuild iteratively  $\widehat{\beta}^{\text{OLS}}$  according to a new criterion which is the difference between signals over the noise. At each step, we attempt to bring closer the current estimator to  $\widehat{\beta}^{\text{OLS}}$  by adding the direction that maximizes this criterion. This direction is the one where the two estimators differ the most.

The algorithm starts with the null estimator. At step one, we put  $\widehat{\beta}_o^{\text{PLS}} = 0$  and we define

$$\Delta\text{SN}_o(\mathbf{x}) = \frac{|\mathbf{x}' \widehat{\beta}^{\text{OLS}} - \mathbf{x}' \widehat{\beta}_o^{\text{PLS}}|}{\sigma \sqrt{\mathbf{x}' S^{-1} \mathbf{x}}}.$$

Then, we seek the direction that maximizes  $\Delta\text{SN}_o(\mathbf{x})$ . Since in that case,  $\Delta\text{SN}_o(\mathbf{x}) = \widehat{\rho}_x$ , this direction is spanned by  $\mathbf{w}_1 = \mathbf{s}$ . Now we want to rebuild from 0 the least squares estimator on that direction. So, we have:

$$\underset{\beta \in \text{span}\{\mathbf{w}_1\}}{\text{ArgMin}} \ \|y - X\beta\|^2 = \widehat{\beta}_{(1)}^{\text{PLS}}. \quad (28)$$

Iteratively, at step  $q + 1 \leq q^*$ , the current estimator is  $\widehat{\beta}_q^{\text{PLS}} = P_{W_{(q)}}^S \widehat{\beta}^{\text{OLS}}$  where  $W_{(q)}$  is defined in Section 4.1. To construct the subspace of dimension  $q + 1$  which gives  $\widehat{\beta}_{q+1}^{\text{PLS}}$ , we add to  $W_{(q)}$  the direction that maximizes  $\Delta\text{SN}_q(\mathbf{x})$ , where

$$\Delta\text{SN}_q(\mathbf{x}) = \frac{|\mathbf{x}' \widehat{\beta}^{\text{OLS}} - \mathbf{x}' P_{W_{(q)}}^S \widehat{\beta}^{\text{OLS}}|}{\sigma \sqrt{\mathbf{x}' S^{-1} \mathbf{x}}} = \frac{|\mathbf{x}' \widehat{\beta}^{\text{OLS}} - \mathbf{x}' \widehat{\beta}_q^{\text{PLS}}|}{\sigma \sqrt{\mathbf{x}' S^{-1} \mathbf{x}}}. \quad (29)$$

**Theorem 10**  $\forall q \leq q^* - 1$ ,

$$\text{ArgMax}_{x \in \mathbb{R}^p} \Delta \text{SN}_q(x) = w_{q+1} \quad (30)$$

where  $w_{q+1}$  is the  $(q+1)^{\text{th}}$  element of the sequence given by Theorem 3. Moreover,

$$\max_{x \in \mathbb{R}^p} \Delta \text{SN}_q(x) = \hat{\rho}_{w_{q+1}}. \quad (31)$$

**PROOF.** We have:  $P_{W_q}^S S^{-1} = S^{-1} P_{S W_q}^{S^{-1}}$ . By Cauchy-Schwarz's inequality:

$$\Delta \text{SN}_q(x) = \frac{|x' S^{-1} (s - P_{S W_q}^{S^{-1}} s)|}{\sqrt{x' S^{-1} x}} \leq \|s - P_{S W_q}^{S^{-1}} s\|_{S^{-1}}.$$

Equality holds iff  $x \propto (I - P_{S W_q}^{S^{-1}})s$ . Since  $\text{Range}(\overline{W}_q)$  is the orthogonal to  $\text{Range}(W_q)$  w.r.t.  $I$ ,  $\overline{W}_q$  is the orthogonal to  $\text{Range}(S W_q)$  w.r.t.  $S^{-1}$ . So,  $(I - P_{S W_q}^{S^{-1}})s = P_{\overline{W}_q}^{S^{-1}} s \in \text{Range}(\overline{W}_q)$ . Thus,  $\Delta \text{SN}_q(x)$  attains its maximum on  $\text{Range}(\overline{W}_q)$ . Since  $\Delta \text{SN}_q(x)$  and  $\hat{\rho}_x$  coincide on  $\text{Range}(\overline{W}_q)$ ,

$$\max_{x \in \mathbb{R}^p} \Delta \text{SN}_q(x) = \max_{x \in \text{Range}(\overline{W}_q)} \Delta \text{SN}_q(x) = \max_{x \in \text{Range}(\overline{W}_q)} \hat{\rho}_x = \hat{\rho}_{w_{q+1}}.$$

The last equality follows from Theorem 3.  $\square$

Thus,  $\text{Range}(W_{(q+1)})$  is the subspace that we are looking for, where  $W_{q+1} = (w_1 | \dots | w_{q+1})$ . The corresponding rebuilt least squares estimator is then  $\hat{\beta}_{(q+1)}^{\text{PLS}}$  since

$$\hat{\beta}_{(q+1)}^{\text{PLS}} = \underset{\beta \in \text{Range}(W_{(q+1)})}{\text{ArgMin}} \|y - X\beta\|^2. \quad (32)$$

Formulas 28 and 32 follows from Formula 27 since  $\beta \in \text{Range}(W_{(q+1)})$  is equivalent to  $\overline{W}'_{(q+1)}\beta = 0$ .

At step  $q^*$ ,  $\hat{\beta}_{q^*}^{\text{PLS}} = \hat{\beta}^{\text{OLS}}$  and  $\hat{\rho}_{w_{q^*+1}} = 0$ .

## 6 Some comments:

The two algorithms presented in Sections 4 and 5 focus on the way PLSR modify the OLS estimator from a directional point of view. The first one is based on the signal-noise ratio  $\hat{\rho}_x$ . This ratio appears naturally to improve the estimate of  $x'\beta$  by applying a shrinkage factor. It plays the same role for PLSR as the variance (or noise) for PCR. The second one is based on  $\Delta \text{SN}_q(x)$ , which is a measurement of the difference between the current estimator and the OLS

estimator according to  $x$ . The interesting feature of these two criteria is that there are scale invariant. Therefore, unlike the latent variables approach, no arbitrary scale constraints are needed .

Both algorithms differ in their structure. In the first one, the maximization procedure of  $\hat{\rho}_x$  is carried out to get the whole sequence of directions. Then, the latest directions are used as constraints in the least squares minimization. In the second approach, the alternation of  $\Delta\text{SN}_q(x)$  maximization and least squares minimization at each step gives the PLS estimators. This alternation allows to release the orthogonality constraints used in the first algorithm.

Both algorithms differ also in their interpretation. The first one shows PLS estimators as least squares estimators forcing to be null onto the smallest signal-noise ratio directions of the sequence. The second algorithm shed the light on how PLSR reduce the difference in terms of  $\Delta\text{SN}_q(x)$  between the current estimator and the OLS estimator.

## 7 The non full-rank case

We discuss briefly the non full-rank case. If  $\text{Rank}(S) = r < p$ , the same procedure describe in Section 4 can be applied but need to be adjusted. The OLS estimator  $\hat{\beta}^{\text{OLS}}$  is non unique, but if  $x'_o\beta$  is estimable, i.e. if  $x_o \in \text{Range } S$ ,  $x'\hat{\beta}^{\text{OLS}}$  and  $\text{var}(x'\hat{\beta}^{\text{OLS}})$  do not depend on the choice of  $\hat{\beta}^{\text{OLS}}$ . So, we seek directions  $w_1, \dots, w_r$  in  $\text{Range}(S)$  rather than in the whole space. With this restriction, part (a) of Theorem 3 and Theorem 10 are still valid and  $w_1, \dots, w_{q^*}$  are identical to Martens' weight vectors. Denote by  $w_{r+1}, \dots, w_p$  an orthogonal sequence in  $\text{Range}(S)^\perp$ . For  $q \leq r$ , we define  $\overline{W}_{(q)} = (w_{q+1} | \dots | w_p)$ . It can be shown that for  $q < q^*$ ,

$$\hat{\beta}_{(q)}^{\text{PLS}} = \underset{\beta: \overline{W}_{(q)}'\beta=0}{\text{ArgMin}} \|y - X\beta\|^2 \quad (33)$$

and for  $q^* \leq q \leq r$  :

$$\underset{\beta: \overline{W}_{(q)}'\beta=0}{\text{ArgMin}} \|y - X\beta\|^2 = \hat{\beta}^{\text{OLS}*} \quad (34)$$

where  $\hat{\beta}^{\text{OLS}*} = S^+s$  with  $S^+$  the Moore-Penrose inverse of  $S$ . The constraints  $w'_{r+1}\beta = 0, \dots, w'_p\beta = 0$  ensures that the estimator belongs to  $\text{Range}(S)$ . The interpretation of estimators in term of projectors can be easily adapted by using generalized projectors (see Rao and Mitra, 1974) and choosing  $\hat{\beta}^{\text{OLS}*}$  as initial OLS estimator.

## References

- [1] Bibby, J., Toutenburg, H., 1977. Prediction and improved estimation in linear models. John Wiley & Sons.
- [2] Breiman, L., Spector, P., 1992. Submodel selection and evaluation in regression. The  $X$ -random case. *International Statistical Review* 60, 291–319.
- [3] Brown, P. J., 1993. Measurement, regression, and calibration. Oxford University Press.
- [4] Butler, N. A., Denham, M. C., 2000. The peculiar shrinkage properties of partial least squares regression. *Journal of the Royal Statistical Society, Series B, Methodologica* 62, 585–593.
- [5] De Jong, S., 1995. PLS shrinks. *Journal of Chemometrics* 9, 323–326.
- [6] Frank, I. E., Friedman, J. H., 1993. A statistical view of some chemometrics regression tools (Disc: p136-148). *Technometrics* 35, 109–135.
- [7] Goutis, C., 1996. Partial least squares algorithm yields shrinkage estimators. *The Annals of Statistics* 24, 816–824.
- [8] Helland, I. S., 1988. On the structure of partial least squares regression. *Communications in Statistics, Part B – Simulation and Computation* [Split from: @J(CommStat)] 17, 581–607.
- [9] Helland, I. S., 1990. Partial least squares regression and statistical models. *Scandinavian Journal of Statistics* 17, 97–114.
- [10] Helland, I. S., 1992. Maximum likelihood regression on relevant components. *Journal of the Royal Statistical Society, Series B, Methodological* 54, 637–647.
- [11] Helland, I. S., Almøy, T., 1994. Comparison of prediction methods when only a few components are relevant. *Journal of the American Statistical Association* 89, 583–591.
- [12] Lingjærde, O. C., Christophersen, N., 2000. Shrinkage structure of partial least squares. *Scandinavian Journal of Statistics* 27 (3), 459–473.
- [13] Martens, H., Næs, T., 1989. *Multivariate calibration*. John Wiley & Sons.
- [14] Rao, C., Mitra, S., 1974. Projections under semi-norms and generalized inverse of matrices. *Linear Algebra and Applications* 9, 155–167.
- [15] Seber, G. A. F., 1984. *Multivariate observations*. John Wiley & Sons.
- [16] Stone, M., Brooks, R. J., 1990. Continuum regression: Cross-validated sequentially constructed prediction embracing ordinary least squares, partial least squares and principal components regression (Corr: V54 p906-07). *Journal of the Royal Statistical Society, Series B, Methodological* 52.
- [17] Tibshirani, R., 1996. Regression shrinkage and selection via the lasso. *Journal of the Royal Statistical Society, Series B, Methodological* 58, 267–288.