



LUDWIG-  
MAXIMILIANS-  
UNIVERSITÄT  
MÜNCHEN

INSTITUT FÜR STATISTIK  
SONDERFORSCHUNGSBEREICH 386



Krämer, Boulesteix, Tutz:

## Penalized Partial Least Squares Based on B-Splines Transformations

Sonderforschungsbereich 386, Paper 485 (2006)

Online unter: <http://epub.ub.uni-muenchen.de/>

Projektpartner



# Penalized Partial Least Squares Based on B-Splines Transformations

**Nicole Krämer**

`nkraemer@cs.tu-berlin.de`

Department of Electrical Engineering and Computer Science  
Technical University of Berlin  
Franklinstr. 28/29, 10587 Berlin, Germany

**Anne-Laure Boulesteix**

`anne-laure.boulesteix@tum.de`

Department of Medical Statistics and Epidemiology  
Technical University of Munich  
Ismaningerstr. 22, 81675 Munich, Germany

**Gerhard Tutz**

`tutz@stat.uni-muenchen.de`

Department of Statistics  
University of Munich  
Akademiestr. 1, 80799 Munich, Germany

## Abstract

We propose a novel method to model nonlinear regression problems by adapting the principle of penalization to Partial Least Squares (PLS). Starting with a generalized additive model, we expand the additive component of each variable in terms of a generous amount of B-Splines basis functions. In order to prevent overfitting and to obtain smooth functions, we estimate the regression model by applying a penalized version of PLS. Although our motivation for penalized PLS stems from its use for B-Splines transformed data, the proposed approach is very general and can be applied to other penalty terms or to other dimension reduction techniques. It turns out that penalized PLS can be computed virtually as fast as PLS. We prove a close connection of penalized PLS to the solutions of preconditioned linear systems. In the case of high-dimensional data, the new method is shown to be an attractive competitor to other techniques for estimating generalized additive models. If the number of predictor variables is high compared to the number of examples, traditional techniques often suffer from overfitting. We illustrate that penalized PLS performs well in these situations.

**Keywords:** generalized additive model, dimension reduction, nonlinear regression, conjugate gradient

# 1 Introduction

Nonlinear regression effects may be modeled via additive regression models of the form

$$Y = \beta_0 + f_1(X_1) + \cdots + f_p(X_p) + \varepsilon. \quad (1)$$

where the functions  $f_1, \dots, f_p$  have unspecified functional form. An approach which allows flexible representation of the functions  $f_1, \dots, f_p$  is the expansion in basis functions (Hastie & Tibshirani 1990). To prevent overfitting, there are two general approaches. In the first approach, each function  $f_j$  is the sum of only a small set of basis functions,

$$f_j(x) = \sum_{k=1}^{K_j} \beta_{kj} B_{kj}(x). \quad (2)$$

The basis functions  $B_{kj}$  are chosen adaptively by a selection procedure. The second approach (that is outlined in Section 3) circumvents the problem of basis function selection. Instead, we allow a generous amount  $K_j \gg 1$  of basis functions in the expansion (2). As this usually leads to high-dimensional and highly correlated data, we penalize the coefficients  $\beta_{jk}$  in the estimation process (Eilers & Marx 1996).

Quite generally, a different approach to deal with high dimensionality is to use dimension reduction techniques such as Partial Least Squares (PLS) (Wold 1975, Wold et al. 1984). The main idea is to build a few components from the predictor variables and to regress  $\mathbf{y}$  onto these components. A short overview on PLS can be found in Section 2.

As a linear approach, PLS probably fails to yield high prediction accuracy in the case of nonlinear relationships between predictors and responses as in (1). In order to incorporate nonlinear structures, it might be advisable to transform the original predictors preliminarily to a PLS regression. This approach has been proposed by Durand & Sabatier (1997) and Durand (2001) in different variants. The method proposed by Durand & Sabatier (1997) is based on a variant of PLS that may be computed via an iterative algorithm. They suggest an approach that incorporates splines transformations of the predictors within each iteration of the iterative algorithm. In contrast, the method proposed by Durand (2001) is global. The predictors are first transformed using splines basis functions as a preliminary step, then PLS regression is performed on the transformed data matrix. The choice of the degree  $d$  of the polynomial pieces and of the number of knots is performed by an either ascending or descending search procedure that is not automatic.

For large numbers of variables, this search procedure is computationally intensive and might overfit the training data. In the present article, we suggest an alternative approach based on the penalty strategy of Eilers & Marx (1996). As described in Section 3, we transform the initial data matrix nonlinearly using B-splines basis functions. Our new method, which we call penalized PLS, is based on the following principle. The equivalent of penalizing the (higher order) differences of adjacent B-splines coefficients is, in the

framework of dimension reduction, the penalization of (higher order) differences of adjacent weights.

In Section 4, we introduce an adaptation of the principle of penalization to PLS. More precisely, we present a penalized version of the optimization problem attached to PLS. Although the motivation stems from its use for B-splines transformed data, the proposed approach is very general and can be adapted to other penalty terms or to other dimension reduction techniques such as Principal Components Analysis. It turns out that the new method shares a lot of properties of PLS and that its computation requires virtually no extra costs. We highlight the close connection between penalized PLS and preconditioned linear systems. It is already known that PLS is equivalent to the conjugate gradient method (Hestenes & Stiefel 1952) applied to the set of normal equations associated to a linear regression problem. We prove that penalized PLS corresponds to a conjugate gradient method for a preconditioned set of normal equations, where the preconditioner depends on the penalty term. Furthermore, we show that this new technique is closely related to the so-called kernel trick. More precisely, we prove that penalized PLS is equivalent to ordinary PLS using a generalized inner product that is defined by the penalty term. In Sections 5 and 6, we illustrate our method on different data sets.

In the rest of the paper, we restrict ourselves to a univariate response. In Section 7, we stress that the extension of our method to a multivariate response is straightforward.

## 2 Partial Least Squares Regression

Let us consider the general linear regression problem. We want to predict a univariate response variable  $Y$  using  $p$  predictor variables  $X_1, \dots, X_p$  based on a finite set

$$\{(y_i, \mathbf{x}_i) = (y_i, x_{i1}, \dots, x_{ip}), i = 1, \dots, n\}$$

of observations. We set

$$\mathbf{X} = \begin{pmatrix} \mathbf{x}_1^T \\ \dots \\ \mathbf{x}_n^T \end{pmatrix} \in \mathbb{R}^{n \times p}, \quad \mathbf{y} = \begin{pmatrix} y_1 \\ \dots \\ y_n \end{pmatrix} \in \mathbb{R}^n,$$

and require for simplicity of notation that both  $\mathbf{X}$  and  $\mathbf{y}$  are centered. If we assume that the relationship between predictors and response is linear, this relationship can be represented in compact form by

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}.$$

Here,  $\boldsymbol{\beta}$  is the  $p$ -dimensional vector of regression coefficients and  $\boldsymbol{\epsilon}$  is the vector of residuals.

When  $n < p$ , the usual regression tools such as ordinary least squares (OLS) regression

cannot be applied to estimate  $\beta$  since the  $p \times p$  covariance matrix  $(1/n)\mathbf{X}^T\mathbf{X}$  (which has rank at most  $n - 1$ ) is singular. From a technical point of view, this may be solved by replacing the inverse of the covariance matrix by a generalized inverse. However, for  $n < p$ , OLS usually fits the training data perfectly and one cannot expect the method to perform well on a new data set. Partial Least Squares (PLS) (Wold 1975, Wold et al. 1984) is an alternative regression tool which is more appropriate in the case of highly correlated predictors and high-dimensional data. PLS is a standard tool for analyzing chemical data (Martens & Naes 1989), and in recent years, the success of PLS has led to applications in other scientific fields such as physiology (Rosipal et al. 2003) or bioinformatics (Boulesteix & Strimmer 2006).

The main idea of PLS is to build orthogonal components  $\mathbf{t}_1, \dots, \mathbf{t}_m$  from the original predictors  $\mathbf{X}$  and to use them as predictors in a least squares regression. There are different PLS techniques to extract these components, and each of them gives rise to a different variant of PLS. It is not our aim to explain all variants and we focus on two of them. An overview on different forms of PLS can be found in Rosipal & Krämer (2006). A component is a linear combination of the original predictors that hopefully reflects the relevant structure of the data. PLS is similar to Principal Components Regression (PCR). The difference is that PCR extracts components that explain the variance in the predictor variables whereas PLS extracts components that have a large covariance with  $\mathbf{y}$ . We now formalize this concept. A latent component  $\mathbf{t}$  is a linear combination  $\mathbf{t} = \mathbf{X}\mathbf{w}$  of the predictor variables. The vector  $\mathbf{w}$  is usually called the weight vector. We want to find a component with maximal covariance to  $\mathbf{y}$ , that is we want to maximize the empirical squared covariance

$$\text{cov}^2(\mathbf{X}\mathbf{w}, \mathbf{y}) = \mathbf{w}^T \mathbf{X}^T \mathbf{y} \mathbf{y}^T \mathbf{X} \mathbf{w}.$$

We have to constrain  $\mathbf{w}$  in order to obtain identifiability, choosing

$$\max \quad \mathbf{w}^T \mathbf{X}^T \mathbf{y} \mathbf{y}^T \mathbf{X} \mathbf{w}, \quad (3)$$

$$\text{subject to} \quad \|\mathbf{w}\| = 1. \quad (4)$$

Using Lagrangian multipliers, we conclude that the solution  $\mathbf{w}_1$  is – up to a scaling factor – equal to  $\mathbf{X}^T \mathbf{y}$ .

Let us remark that (3) and (4) are equivalent to

$$\max \quad \frac{\mathbf{w}^T \mathbf{X}^T \mathbf{y} \mathbf{y}^T \mathbf{X} \mathbf{w}}{\mathbf{w}^T \mathbf{w}}. \quad (5)$$

The solution of (5) is only unique up to a scalar. The normalization of the weight vectors  $\mathbf{w}$  to length 1 is not essential for the PLS algorithm and PLS algorithms differ in the way they scale the weight vectors and components. In this paper, we present all algorithms without the scaling of the vectors, in order to keep the notation as simple as possible.

Subsequent components  $\mathbf{t}_2, \mathbf{t}_3, \dots$  are chosen such that they maximize (3) and that all components  $\mathbf{t}_i$  are mutually orthogonal. In PLS, there are different techniques to extract subsequent components, and each technique gives rise to a variant of PLS. We briefly introduce two of them. In the method called SIMPLS (de Jong 1993), one computes for the  $i$ th component,

$$\begin{aligned} \max \quad & \mathbf{w}^T \mathbf{X}^T \mathbf{y} \mathbf{y}^T \mathbf{X} \mathbf{w}, \\ \text{subject to} \quad & \|\mathbf{w}\| = 1 \text{ and } \mathbf{X} \mathbf{w} \perp \mathbf{t}_j, j < i. \end{aligned}$$

Alternatively, one can deflate the original predictor variables  $\mathbf{X}$ . That is, we only consider the part of  $\mathbf{X}$  that is orthogonal onto all components  $\mathbf{t}_j, j < i$ . For any matrix  $\mathbf{V}$ , let us denote by  $\mathcal{P}_{\mathbf{V}}$  the orthogonal projection onto the space that is spanned by the columns of  $\mathbf{V}$ . In matrix notation, we have

$$\mathcal{P}_{\mathbf{V}} = \mathbf{V} (\mathbf{V}^T \mathbf{V})^{-1} \mathbf{V}^T. \quad (6)$$

The deflation of  $\mathbf{X}$  with respect to the components  $\mathbf{t}_1, \dots, \mathbf{t}_{i-1}$  is defined as

$$\mathbf{X}_i = \mathbf{X} - \mathcal{P}_{\mathbf{t}_1, \dots, \mathbf{t}_{i-1}} \mathbf{X} = \mathbf{X}_{i-1} - \mathcal{P}_{\mathbf{t}_{i-1}} \mathbf{X}_{i-1}. \quad (7)$$

For the computation of the  $i$ th component,  $\mathbf{X}$  is replaced by  $\mathbf{X}_i$  in (3). This method is called the NIPALS algorithm (Wold 1975). The two methods are equivalent if  $\mathbf{y}$  is univariate in the sense that we end up with the same components  $\mathbf{t}_i$  (de Jong 1993). In this paper, we use the NIPALS algorithm. In summary, the PLS algorithm is described in algorithm 1.

---

**Algorithm 1** NIPALS algorithm

---

Input:  $X_1 = X$ ,  $\mathbf{y}$ , number of components  $m$

**for**  $i=1, \dots, m$  **do**

$\mathbf{w}_i = \mathbf{X}_i^T \mathbf{y}$  (weight vector)

$\mathbf{t}_i = \mathbf{X}_i \mathbf{w}_i$  (component)

$\mathbf{X}_{i+1} = \mathbf{X}_i - \mathcal{P}_{\mathbf{t}_i} \mathbf{X}_i$  (deflation)

**end for**

---

PLS used to be overlooked by statisticians and was considered an algorithm rather than a sound statistical model. This attitude is in parts understandable, as in the early literature on the subject, PLS was explained solely in terms of formulas as in algorithm 1. Due to its success in applications, the interest in the statistical properties of PLS has risen. It can be related to other dimension reduction techniques such as Principal Components Regression and Ridge Regression and these methods can be cast under a unifying framework (Stone & Brooks 1990). The shrinkage properties of PLS have been studied extensively (Frank & Friedman 1993, de Jong 1995, Goutis 1996, Butler & Denham 2000). Furthermore, it can be shown that PLS is closely connected to Krylov subspaces and the conjugate gradient method (Helland 1988, Phatak & de Hoog 2003). We discuss this method in more detail

in Section 4.

Let us return to the PLS algorithm. With

$$\mathbf{T} = (\mathbf{t}_1, \dots, \mathbf{t}_m) .$$

denoting the collection of components, the fitted response is given by

$$\hat{\mathbf{y}} = \mathbf{T}(\mathbf{T}^T\mathbf{T})^{-1}\mathbf{T}^T\mathbf{y} = \mathcal{P}_{\mathbf{T}}\mathbf{y} . \quad (8)$$

In order to obtain the response for new observations, we have to determine the vector of regression coefficients  $\hat{\mathbf{y}} = \mathbf{X}\hat{\boldsymbol{\beta}}$ . Therefore, a representation of the components  $\mathbf{t}_i = \mathbf{X}_i\mathbf{w}_i$  as a linear combination of the original predictors  $\mathbf{X}$  is needed. In other words, we have to derive weight vectors  $\tilde{\mathbf{w}}_i$  with

$$\mathbf{X}\tilde{\mathbf{w}}_i = \mathbf{X}_i\mathbf{w}_i .$$

They are in general different from the “pseudo” weight vectors  $\mathbf{w}_i$  that are computed by the NIPALS algorithm. In order to avoid redundancy, the derivation of these weight vectors is deferred until Section 4.

It should be noted that the number  $m$  of PLS components is an additional model parameter that has to be estimated. One way of determining  $m$  is by cross-validation.

### 3 Penalized Regression Splines

The fitting of generalized additive models by use of penalized regression splines has become a widely used tool in statistics. Starting with the seminal paper by Eilers & Marx (1996), the approach has been extended and applied in various publications (Ruppert 2002, Wood 2000, Wood 2006). The basic concept is to expand the additive component of each variable  $X_j$  in basis functions as in (2) and to estimate the coefficients by penalization techniques. As suggested in Eilers & Marx (1996), B-splines are used as basis functions yielding so-called P-splines (for penalized B-splines). Splines are one-dimensional piecewise polynomial functions. The points at which the pieces are connected are called knots or breakpoints. We say that a spline is of order  $d$  if all polynomials are of degree  $\leq d$  and if the spline is  $(d - 1)$  times continuously differentiable at the breakpoints. A particular efficient set of basis functions are B-splines (de Boor 1978). The number of basis functions depends on the order of the splines and the number of breakpoints. For a given variable  $X_j$ , we consider a set of corresponding B-splines basis functions  $B_{1j}, \dots, B_{Kj}$ . These basis functions define a nonlinear map

$$\Phi_j(x) = (B_{1j}(x), \dots, B_{Kj}(x))^T .$$

By performing such a transformation on each of the variables  $X_1, \dots, X_p$ , the observation vector  $\mathbf{x}_i$  turns into a vector

$$\begin{aligned} \mathbf{z}_i &= (B_{11}(x_{i1}), \dots, B_{m1}(x_{i1}), \dots, B_{1p}(x_{ip}), \dots, B_{mp}(x_{ip}))^T \\ &= \Phi(\mathbf{x}_i) \end{aligned} \quad (9)$$

of length  $pK$ . Here  $\Phi$  is the function defined by the B-splines. The resulting data matrix obtained by the transformation of  $\mathbf{X}$  has dimensions  $n \times pK$  and will be denoted by  $\mathbf{Z}$  in the rest of the paper. In the examples in Sections 5 and 6, we consider the most widely used cubic B-splines, i.e. we choose  $d = 3$ .

The estimation of (1) is transformed into the estimation of the  $pK$ -dimensional vector that consists of the coefficients  $\beta_{jk}$ :

$$\boldsymbol{\beta}^T = (\beta_{11}, \dots, \beta_{K1}, \dots, \beta_{12}, \dots, \beta_{Kp}) = (\boldsymbol{\beta}_{(1)}^T, \dots, \boldsymbol{\beta}_{(p)}^T).$$

As explained above, the vector  $\boldsymbol{\beta}$  determines a nonlinear, additive function

$$f(\mathbf{x}) = \beta_0 + \sum_{j=1}^p f_j(x_j) = \beta_0 + \sum_{j=1}^p \sum_{k=1}^K \beta_{kj} B_{kj}(x_j) = \beta_0 + \Phi(\mathbf{x})^T \boldsymbol{\beta}.$$

As  $\mathbf{Z}$  is usually high-dimensional, the estimation of  $\boldsymbol{\beta}$  by minimizing the squared error

$$\frac{1}{n} \sum_{i=1}^n (y_i - f(\mathbf{x}_i))^2 = \frac{1}{n} \|\mathbf{y} - \beta_0 - \mathbf{Z}\boldsymbol{\beta}\|^2$$

usually leads to overfitting. Following Eilers & Marx (1996), we use for each variable many basis functions, say  $K \approx 20$ , and estimate by penalization. The idea is to penalize the second derivative of the function  $f$ . Eilers & Marx (1996) show that the following difference penalty term is a good approximation of the penalty on the second derivative of  $f$ ,

$$P(\boldsymbol{\beta}) = \sum_{j=1}^p \sum_{k=3}^m \lambda_j (\Delta^2 \beta_{kj})^2.$$

These are also called the second-order differences of adjacent parameters. The difference operator  $\Delta^2 \beta_{kj}$  has the form

$$\begin{aligned} \Delta^2 \beta_{kj} &= (\beta_{kj} - \beta_{k-1,j}) - (\beta_{k-1,j} - \beta_{k-2,j}) \\ &= \beta_{kj} - 2\beta_{k-1,j} + \beta_{k-2,j}. \end{aligned}$$

The coefficients  $\lambda_j \geq 0$  control the amount of penalization. This penalty term can be



expressed in terms of a penalty matrix  $\mathbf{P}$ . We denote by  $\mathbf{D}_K$  the  $(K - 1) \times K$  matrix

$$\mathbf{D}_K = \begin{pmatrix} 1 & -1 & . & . & . \\ . & 1 & -1 & . & . \\ . & . & . & . & . \\ . & . & . & 1 & -1 \end{pmatrix}$$

that defines the first order difference operator. Setting

$$\mathbf{K}_2 = (\mathbf{D}_{K-1}\mathbf{D}_K)^T \mathbf{D}_{K-1}\mathbf{D}_K,$$

we conclude that the penalty term equals

$$P(\boldsymbol{\beta}) = \sum_{j=1}^p \lambda_j \boldsymbol{\beta}_{(j)}^T \mathbf{K}_2 \boldsymbol{\beta}_{(j)} = \boldsymbol{\beta}^T (\boldsymbol{\Delta}_\lambda \otimes \mathbf{K}_2) \boldsymbol{\beta}.$$

Here  $\boldsymbol{\Delta}_\lambda$  is the  $p \times p$  diagonal matrix containing  $\lambda_1, \dots, \lambda_p$  on its diagonal and  $\otimes$  is the Kronecker product. The generalization of this method to higher-order differences of the coefficients of adjacent B-splines is straightforward. We simply replace  $\mathbf{K}_2$  by

$$\mathbf{K}_q = (\mathbf{D}_{K-q+1} \dots \mathbf{D}_K)^T (\mathbf{D}_{K-q+1} \dots \mathbf{D}_K).$$

To summarize, the penalized least squares criterion has the form

$$\widehat{R}_P(\boldsymbol{\beta}) = \frac{1}{n} \|\mathbf{y} - \beta_0 - \mathbf{Z}\boldsymbol{\beta}\|^2 + \boldsymbol{\beta}^T \mathbf{P} \boldsymbol{\beta} \quad (10)$$

with the penalty matrix  $\mathbf{P}$  defined as

$$\mathbf{P} = \boldsymbol{\Delta}_\lambda \otimes \mathbf{K}_q. \quad (11)$$

This is a symmetric matrix that is positive semidefinite.

## 4 Penalized Partial Least Squares Regression

We now introduce a general framework to combine PLS with penalization terms. We remark that this is not limited to spline transformed variables or to the special shape of the penalty matrix  $\mathbf{P}$  that is defined in (11). For this reason, we present the new method in terms of the original data matrix  $\mathbf{X}$  and only demand that  $\mathbf{P}$  is a symmetric matrix such that  $\mathbf{I}_p + \mathbf{P}$  is positive definite.

Again, we restrict ourselves to univariate responses  $\mathbf{y}$ . Penalized PLS for multivariate responses is briefly discussed in Section 7. We modify the optimization criterion (5) of PLS in the following way. The first component  $\mathbf{t}_1 = \mathbf{X}\mathbf{w}_1$  is defined by the solution of the

problem

$$\arg \max_{\mathbf{w}} \frac{\mathbf{w}^T \mathbf{X}^T \mathbf{y} \mathbf{y}^T \mathbf{X} \mathbf{w}}{\mathbf{w}^T \mathbf{w} + \mathbf{w}^T \mathbf{P} \mathbf{w}}. \quad (12)$$

Using Lagrangian multipliers, we obtain the solution

$$\mathbf{w}_1 = \mathbf{M} \mathbf{X}^T \mathbf{y} \quad (13)$$

with  $\mathbf{M} = (\mathbf{I}_p + \mathbf{P})^{-1}$ . Subsequent weight vectors and components are computed by deflating  $\mathbf{X}$  as described in (7) and then maximizing (12) with  $\mathbf{X}$  replaced by  $\mathbf{X}_i$ . In particular, we can compute the weight vectors and components of penalized PLS by simply replacing  $\mathbf{w}_i = \mathbf{X}_i^T \mathbf{y}$  by (13) in algorithm 1.

We now present results on penalized PLS that allow us to compute its regression vectors efficiently. Note that all results on penalized PLS also hold for ordinary PLS if we choose  $\mathbf{P} = \mathbf{0}$ . Let

$$\mathbf{T} = (\mathbf{t}_1, \dots, \mathbf{t}_m) \quad , \quad \mathbf{W} = (\mathbf{w}_1, \dots, \mathbf{w}_m) \quad ,$$

denote the matrices of components and weight vectors respectively.

**Lemma 1.** *The matrix*

$$\mathbf{R} = \mathbf{T}^T \mathbf{X} \mathbf{W} \in \mathbb{R}^{m \times m}$$

*is upper bidiagonal, that is*

$$r_{ij} = \mathbf{t}_i^T \mathbf{X} \mathbf{w}_j = 0$$

*if  $i < j$  or  $i + 1 > j$ . The matrix  $\mathbf{R}$  is invertible. Furthermore, the columns of  $\mathbf{T}$  and the columns of  $\mathbf{X} \mathbf{W}$  span the same space.*

This is an extension of a result for ordinary PLS that can be found e.g. in Manne (1987). The proof can be found in the appendix. We can now determine the regression coefficients for penalized PLS.

**Proposition 2.** *The Penalized PLS regression vector obtained after  $m$  steps is*

$$\widehat{\boldsymbol{\beta}}_{PLS}^{(m)} = \mathbf{W} (\mathbf{W}^T \mathbf{X}^T \mathbf{X} \mathbf{W})^{-1} \mathbf{W}^T \mathbf{X}^T \mathbf{y}. \quad (14)$$

*In particular, the penalized PLS estimator is the solution of the constrained minimization problem*

$$\begin{aligned} \min_{\boldsymbol{\beta}} \quad & \|\mathbf{y} - \mathbf{X} \boldsymbol{\beta}\|^2 \\ \text{subject to} \quad & \boldsymbol{\beta} \in \text{span}\{\mathbf{w}_1, \dots, \mathbf{w}_m\}. \end{aligned} \quad (15)$$

*Proof.* We deduce from lemma 1 that the columns of  $\mathbf{X} \mathbf{w}$  span the same space as the columns of  $\mathbf{T}$ . As PLS is ordinary least squares regression with predictors  $\mathbf{t}_1, \dots, \mathbf{t}_m$ , we

have

$$\hat{\mathbf{y}} = \mathcal{P}_{\mathbf{T}}\mathbf{y} = \mathcal{P}_{\mathbf{X}\mathbf{W}}\mathbf{y} = \mathbf{X}\mathbf{W} (\mathbf{W}^T \mathbf{X}^T \mathbf{X}\mathbf{W})^{-1} \mathbf{W}^T \mathbf{X}^T \mathbf{y}.$$

The second statement can be proven by noting that the OLS minimization problem with constraints (15) is equivalent to an unconstrained minimization problem for  $\boldsymbol{\beta} = \mathbf{W}\boldsymbol{\alpha}$  with  $\boldsymbol{\alpha} \in \mathbb{R}^m$ . If we plug this into the formula for the OLS estimator, we obtain (14).  $\square$

Formula (14) is beneficial for theoretical purposes but it is computationally inefficient. We now show how the calculation can be done in a recursive and faster way. The key point is to find “effective” weight vectors  $\tilde{\mathbf{w}}_i$  such that for every  $i$

$$\mathbf{t}_i = \mathbf{X}_i \mathbf{w}_i = \mathbf{X} \tilde{\mathbf{w}}_i. \quad (16)$$

This can be done by exploiting the fact that  $\mathbf{R}$  is bidiagonal.

**Proposition 3.** *The effective weight vectors  $\tilde{\mathbf{w}}_i$  defined in (16) and the regression vectors of penalized PLS are determined by setting  $\tilde{\mathbf{w}}_0 = \mathbf{0}$  and  $\hat{\boldsymbol{\beta}}^{(0)} = \mathbf{0}$  and computing iteratively*

$$\begin{aligned} \tilde{\mathbf{w}}_i &= \mathbf{w}_i - \frac{\tilde{\mathbf{w}}_{i-1}^T \mathbf{X}^T \mathbf{X} \mathbf{w}_i}{\tilde{\mathbf{w}}_{i-1}^T \mathbf{X}^T \mathbf{X} \tilde{\mathbf{w}}_{i-1}} \tilde{\mathbf{w}}_{i-1}, \\ \hat{\boldsymbol{\beta}}^{(i)} &= \hat{\boldsymbol{\beta}}^{(i-1)} + \frac{\tilde{\mathbf{w}}_i^T \mathbf{X}^T \mathbf{Y}}{\tilde{\mathbf{w}}_i^T \mathbf{X}^T \mathbf{X} \tilde{\mathbf{w}}_i} \tilde{\mathbf{w}}_i. \end{aligned}$$

The proof can be found in the appendix. Combining this result with the PLS algorithm 1, we obtain the penalized PLS algorithm 2.

---

**Algorithm 2** Penalized PLS

---

$\mathbf{X}_1 = \mathbf{X}$ ,  $\mathbf{y}$ , number of components  $m$ , penalty matrix  $\mathbf{P}$  (input)  
 $\mathbf{M} = (\mathbf{I}_p + \mathbf{P})^{-1}$ ,  $\tilde{\mathbf{w}}_0 = \mathbf{0}$ ,  $\hat{\boldsymbol{\beta}}^{(0)} = \mathbf{0}$  (initialization)  
**for**  $i=1, \dots, m$  **do**  
 $\mathbf{w}_i = \mathbf{M} \mathbf{X}_i^T \mathbf{y}$  (weight vector)  
 $\tilde{\mathbf{w}}_i = \mathbf{w}_i - \frac{\tilde{\mathbf{w}}_{i-1}^T \mathbf{X}^T \mathbf{X} \mathbf{w}_i}{\tilde{\mathbf{w}}_{i-1}^T \mathbf{X}^T \mathbf{X} \tilde{\mathbf{w}}_{i-1}} \tilde{\mathbf{w}}_{i-1}$  (effective weight vector)  
 $\hat{\boldsymbol{\beta}}^{(i)} = \hat{\boldsymbol{\beta}}^{(i-1)} + \frac{\tilde{\mathbf{w}}_i^T \mathbf{X}^T \mathbf{Y}}{\tilde{\mathbf{w}}_i^T \mathbf{X}^T \mathbf{X} \tilde{\mathbf{w}}_i} \tilde{\mathbf{w}}_i$  (regression vector)  
 $\mathbf{t}_i = \mathbf{X}_i \mathbf{w}_i$  (component)  
 $\mathbf{X}_{i+1} = \mathbf{X}_i - \mathcal{P}_{\mathbf{t}_i} \mathbf{X}_i$  (deflation)  
**end for**

---

#### 4.1 Partial Least Squares and Krylov Subspaces

It is well-known that PLS is closely connected to Krylov subspaces and conjugate gradient methods. Quite generally, linear regression problems can be transformed into algebraic problems in the following way. The OLS estimator is the solution of the minimization

problem

$$\min_{\boldsymbol{\beta}} \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|^2. \quad (17)$$

This is equivalent to finding the solution of the associated normal equation

$$\mathbf{A}\boldsymbol{\beta} = \mathbf{b} \quad (18)$$

with  $\mathbf{b} = \mathbf{X}^T \mathbf{y}$  and  $\mathbf{A} = \mathbf{X}^T \mathbf{X}$ . If the matrix  $\mathbf{A}$  is invertible, the solution of the normal equations is the OLS estimator  $\hat{\boldsymbol{\beta}} = \mathbf{A}^{-1} \mathbf{b}$ . If  $\mathbf{A}$  is singular, the solution of (18) with minimal Euclidean norm is  $\mathbf{A}^+ \mathbf{b}$ . We already mentioned in Section 2 that in the case of high dimensional data, the matrix  $\mathbf{A}$  is often (almost) singular and that the OLS estimator performs poorly on new data sets. A popular strategy is to regularize the least squares criterion (17) in the hope of improving the performance of the estimator. This corresponds to finding approximate solutions of (18). For example, Ridge Regression corresponds to the solution of the modified normal equations

$$(\mathbf{A} + \lambda \mathbf{I}_p) \boldsymbol{\beta} = \mathbf{b}.$$

Here  $\lambda > 0$  is the Ridge parameter. Principal Components Regression uses the eigen decomposition of  $\mathbf{A}$

$$\mathbf{A} = \mathbf{U} \boldsymbol{\Lambda} \mathbf{U}^T = \sum_{i=1}^p \lambda_i \mathbf{u}_i \mathbf{u}_i^T$$

and approximates  $\mathbf{A}$  and  $\mathbf{b}$  via the first  $m$  eigenvectors

$$\mathbf{A} \approx \sum_{i=1}^m \lambda_i \mathbf{u}_i \mathbf{u}_i^T, \quad \mathbf{b} \approx \sum_{i=1}^m (\mathbf{u}_i^T \mathbf{b}) \mathbf{u}_i.$$

It can be shown that the PLS estimators are equal to the approximate solutions of the conjugate gradient method (Hestenes & Stiefel 1952). This is a procedure that iteratively computes approximate solutions of (18) by minimizing the quadratic function

$$\phi(\boldsymbol{\beta}) = \frac{1}{2} \boldsymbol{\beta}^T \mathbf{A} \boldsymbol{\beta} - \boldsymbol{\beta}^T \mathbf{b} = \frac{1}{2} \langle \boldsymbol{\beta}, \mathbf{A} \boldsymbol{\beta} \rangle - \langle \boldsymbol{\beta}, \mathbf{b} \rangle \quad (19)$$

along directions that are  $\mathbf{A}$ -orthogonal. The approximate solution obtained after  $m$  steps is equal to the PLS estimator obtained after  $m$  iterations.

The conjugate gradient algorithm is in turn closely related to Krylov subspaces and the Lanczos algorithm (Lanczos 1950). The latter is a method for approximating eigenvalues. The connection between PLS and these methods is well-elaborated in Phatak & de Hoog (2003). We now establish a similar connection between penalized PLS and the above

mentioned methods. Set

$$\mathbf{A}_M = M\mathbf{A} \quad \text{and} \quad \mathbf{b}_M = M\mathbf{b}.$$

Recall that  $M$  is a symmetric and positive definite matrix that is determined by the penalty term  $P$ . We now illustrate that penalized PLS finds approximate solutions of the preconditioned normal equation

$$\mathbf{A}_M\boldsymbol{\beta} = \mathbf{b}_M. \quad (20)$$

**Lemma 4.** *The space spanned by the weight vectors  $\mathbf{w}_1, \dots, \mathbf{w}_m$  of penalized PLS is the same as the space spanned by the Krylov sequence*

$$\mathbf{b}_M, \mathbf{A}_M\mathbf{b}_M, \dots, \mathbf{A}_M^{m-1}\mathbf{b}_M. \quad (21)$$

This is the generalization of a result for ordinary PLS and can be proven via induction. Details are given in the appendix. We denote by

$$\mathcal{K}^{(m)} = \mathcal{K}^{(m)}(\mathbf{A}_M, \mathbf{b}_M)$$

the space that is spanned by the Krylov sequence (21). This space is called a Krylov space.

**Corollary 5.** *The penalized PLS estimator is the solution of the optimization problem*

$$\begin{aligned} \min \quad & \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|^2 \\ \text{subject to} \quad & \boldsymbol{\beta} \in \mathcal{K}^{(m)}. \end{aligned}$$

*Proof.* This follows immediately from proposition 2 and the fact that the weight vectors span the Krylov space  $\mathcal{K}^{(m)}$ .  $\square$

We now present the conjugate gradient method for the equation

$$\mathbf{A}_M\boldsymbol{\beta} = \mathbf{b}_M. \quad (22)$$

The Conjugate gradient method is normally applied if the involved matrix is symmetric. Note that in general, the matrix  $\mathbf{A}_M$  is not symmetric with respect to the canonical inner product, but with respect to the inner product

$$\langle \mathbf{x}, \tilde{\mathbf{x}} \rangle_{M^{-1}} = \mathbf{x}^T M^{-1} \tilde{\mathbf{x}}$$

defined by  $M^{-1}$ . We can rewrite the quadratic function  $\phi$  defined in (19) as

$$\phi(\boldsymbol{\beta}) = \frac{1}{2} \langle \boldsymbol{\beta}, \mathbf{A}_M\boldsymbol{\beta} \rangle_{M^{-1}} - \langle \boldsymbol{\beta}, \mathbf{b}_M \rangle_{M^{-1}}.$$

We replace the canonical inner product by the inner product defined by  $M^{-1}$  and minimize

this function iteratively along directions that are  $\mathbf{A}_M$ -orthogonal.

We start with an initial guess  $\beta_0 = \mathbf{0}$  and define  $\mathbf{d}_0 = \mathbf{r}_0 = \mathbf{b}_M - \mathbf{A}_M\beta_0 = \mathbf{b}_M$ . The quantity  $\mathbf{d}_m$  is the search direction and  $\mathbf{r}_m$  is the residual. For a given direction  $\mathbf{d}_m$ , we have to determine the optimal step size, that is we have to find

$$a_m = \arg \min_a \phi(\beta_m + a\mathbf{d}_m).$$

It is straightforward to check that

$$a_m = \frac{\langle \mathbf{d}_m, \mathbf{r}_m \rangle_{M^{-1}}}{\langle \mathbf{d}_m, \mathbf{A}_M \mathbf{d}_m \rangle_{M^{-1}}}.$$

The new approximate solution is then

$$\beta_{m+1} = \beta_m + a_m \mathbf{d}_m.$$

After updating the residuals via

$$\mathbf{r}_{m+1} = \mathbf{b}_M - \mathbf{A}_M \beta_{m+1},$$

we define a new search direction  $\mathbf{d}_{m+1}$  that is  $\mathbf{A}_M$ -orthogonal to the previous search directions. This is ensured by projecting the residual  $\mathbf{r}_m$  onto the space that is  $\mathbf{A}_M$ -orthogonal to  $\mathbf{d}_0, \dots, \mathbf{d}_m$ . We obtain

$$\mathbf{d}_{m+1} = \mathbf{r}_{m+1} - \sum_{i=0}^m \frac{\langle \mathbf{r}_{m+1}, \mathbf{A}_M \mathbf{d}_i \rangle_{M^{-1}}}{\langle \mathbf{d}_i, \mathbf{A}_M \mathbf{d}_i \rangle_{M^{-1}}} \mathbf{d}_i.$$

**Theorem 6.** *The penalized PLS algorithm is equal to the conjugate gradient algorithm for the preconditioned system (20).*

The presentation of the conjugate gradient method above and the proof of its equivalence to penalized PLS are an extension of the corresponding results for PLS that is given in Phatak & de Hoog (2003). The proof can be found in the appendix. Note that there is a different notion of conjugate gradients for preconditioned systems (Golub & van Loan 1983). We transform the preconditioned equation (19) by postmultiplying with  $\mathbf{M}$ :

$$\mathbf{MAM}\tilde{\beta} = \mathbf{Mb} \quad \text{with} \quad \tilde{\beta} = \mathbf{M}^{-1}\beta.$$

As the matrix  $\mathbf{MAM}$  is symmetric, we can apply the ordinary conjugate gradient algorithm to this equation. This approach differs from the one described above.

**Proposition 7.** *Suppose that  $\mathbf{A} = \mathbf{X}^T \mathbf{X}$  is regular. After at most  $p$  iterations, the penalized PLS estimator equals the OLS estimator.*

*Proof.* Using (15), the above statement is equivalent to showing that

$$\widehat{\boldsymbol{\beta}}_{OLS} \in \mathcal{K}^{(p)}.$$

Hence, we have to show that there is a polynomial  $\pi$  of degree  $\leq p - 1$  such that  $\widehat{\boldsymbol{\beta}}_{OLS} = \pi(\mathbf{A}_M) \mathbf{b}_M$ . As  $\mathbf{M}$  is invertible, the OLS estimator is

$$\widehat{\boldsymbol{\beta}}_{OLS} = \mathbf{A}^{-1} \mathbf{b} = \mathbf{A}^{-1} \mathbf{M}^{-1} \cdot \mathbf{M} \mathbf{b} = (\mathbf{M} \mathbf{A})^{-1} \mathbf{M} \mathbf{b} = \mathbf{A}_M^{-1} \mathbf{b}_M.$$

As  $\mathbf{A}_M$  is the product of two symmetric matrices and  $\mathbf{M}$  is positive definite,  $\mathbf{A}_M$  has a real eigendecomposition,

$$\mathbf{A}_M = \mathbf{U} \boldsymbol{\Gamma} \mathbf{U}^{-1}.$$

We define the polynomial  $\pi$  via the at most  $p$  equations

$$\pi(\gamma_i) = \frac{1}{\gamma_i}.$$

It follows immediately that  $\pi(\mathbf{A}_M) = \mathbf{A}_M^{-1}$ . This concludes the proof.  $\square$

## 4.2 Kernel Penalized Partial Least Squares

The computation of the penalized PLS estimator as presented in algorithm 2 involves matrices and vectors of dimension  $p \times p$  and  $p$  respectively. If the number of predictors  $p$  is very large, this leads to high computational costs. In this subsection, we show that we can represent this algorithm in terms of matrices and vectors of dimension  $n \times n$  and  $n$  respectively. Let us define the  $n \times n$  matrix  $\mathbf{K}_M$  via

$$\mathbf{K}_M = (\langle \mathbf{x}_i, \mathbf{x}_j \rangle_M) = \mathbf{X} \mathbf{M} \mathbf{X}^T.$$

This matrix is called the Gram matrix or the kernel matrix of  $\mathbf{X}$ . We conclude from corollary 5 that the penalized PLS estimator obtained after  $m$  steps is an element of the Krylov space  $\mathcal{K}^{(m)}(\mathbf{A}_M, \mathbf{b}_M)$ . It follows that we can represent the penalized PLS estimator as

$$\widehat{\boldsymbol{\beta}}^{(m)} = \mathbf{M} \mathbf{X}^T \boldsymbol{\alpha}^{(m)}, \quad \boldsymbol{\alpha}^{(m)} \in \mathcal{K}^{(m)}(\mathbf{K}_M, \mathbf{y}).$$

Here, the Krylov space  $\mathcal{K}^{(m)}(\mathbf{K}_M, \mathbf{y})$  is the space spanned by the vectors

$$\mathbf{y}, \mathbf{K}_M \mathbf{y}, \dots, \mathbf{K}_M^{m-1} \mathbf{y}.$$

Analogously, we can represent the effective weight vectors by

$$\tilde{\mathbf{w}}_m = \mathbf{M} \mathbf{X}^T \tilde{\boldsymbol{\alpha}}_m, \quad \tilde{\boldsymbol{\alpha}}_m \in \mathcal{K}^{(m)}(\mathbf{K}_M, \mathbf{y}).$$

It follows from the definition of the deflation step that

$$\mathbf{X}_m^T \mathbf{y} = \mathbf{X}^T (\mathbf{I}_n - \mathcal{P}_{\mathbf{t}_1, \dots, \mathbf{t}_{m-1}}) \mathbf{y} = \mathbf{X}^T (\mathbf{y} - \hat{\mathbf{y}}^{(m-1)}) .$$

We conclude that the weight vector  $\mathbf{w}_i$  is simply

$$\mathbf{w}_m = \mathbf{M} \mathbf{X}^T \mathbf{y}_{res}^{(m)} , \quad \mathbf{y}_{res}^{(m)} = \mathbf{y} - \hat{\mathbf{y}}^{(m-1)} .$$

If we plug in these representations into the penalized PLS algorithm 2, we obtain algorithm 3 that depends only on  $\mathbf{K}_M$  and  $\mathbf{y}$ .

---

**Algorithm 3** Kernel penalized PLS

---

$\mathbf{X}$ ,  $\mathbf{y}$ , number of components  $m$ , penalty term  $\mathbf{P}$  (input)  
 $\mathbf{M} = (\mathbf{I}_p + \mathbf{P})^{-1}$ ,  $\mathbf{K}_M = \mathbf{X} \mathbf{M} \mathbf{X}^T$ ,  $\boldsymbol{\alpha}^{(0)} = \tilde{\boldsymbol{\alpha}}^{(m)} = \mathbf{0}$  (initialization)  
**for**  $i=1, \dots, m$  **do**  
 $\mathbf{y}_{res}^{(m)} = \mathbf{y} - \hat{\mathbf{y}}^{(m-1)}$  (residuals)  
 $\tilde{\boldsymbol{\alpha}}_m = \mathbf{y}_{res}^{(m)} - \frac{\tilde{\boldsymbol{\alpha}}_{m-1}^T \mathbf{K}_M^2 \mathbf{y}_{res}^{(m)}}{\tilde{\boldsymbol{\alpha}}_{m-1}^T \mathbf{K}_M^2 \tilde{\boldsymbol{\alpha}}_{m-1}} \tilde{\boldsymbol{\alpha}}_{m-1}$  (effective weight vector)  
 $\boldsymbol{\alpha}^{(m)} = \boldsymbol{\alpha}^{(m-1)} + \frac{\tilde{\boldsymbol{\alpha}}_m^T \mathbf{K}_M \mathbf{y}}{\tilde{\boldsymbol{\alpha}}_m^T \mathbf{K}_M^2 \tilde{\boldsymbol{\alpha}}_m} \tilde{\boldsymbol{\alpha}}_m$  (regression vector)  
 $\mathbf{t}_i = \mathbf{K}_M \tilde{\boldsymbol{\alpha}}_m$  (component)  
 $\hat{\mathbf{y}}^{(m+1)} = \hat{\mathbf{y}}^{(m)} + \mathcal{P}_{\mathbf{t}_i} \mathbf{y}$  (estimation of  $\mathbf{y}$ )  
**end for**

---

A kernel version of PLS has already been defined in Rännar et al. (1994) in order to speed up the computation of PLS. We repeat that the speed of the kernel version of penalized PLS does not depend on the number of predictor variables at all but on the number of observations. This implies that – from an algorithmic point of view – there are no restrictions in terms of the number of predictor variables. The importance of this so-called “dual” representation also becomes apparent if we want to extend PLS to nonlinear problems by using the kernel trick. In this paper, the kernel trick appears in two different versions.

Let us only consider the case of ordinary PLS on B-Splines transformed variables. Recall that in (9), we transform the original data  $\mathbf{X}$  using a nonlinear function  $\Phi$  defined by the B-Splines. As algorithm 3 only relies on inner products between observations, the nonlinear transformation does not increase the computational costs. We only have to compute the kernel matrix of inner products

$$\mathbf{K} = (\langle \Phi(\mathbf{x}_i), \Phi(\mathbf{x}_j) \rangle)_{i,j=1, \dots, n} .$$

This implies that we do not have to map the data points explicitly using a function  $\Phi$ . It suffices to compute the function

$$k(\mathbf{x}, \tilde{\mathbf{x}}) = \langle \Phi(\mathbf{x}), \Phi(\tilde{\mathbf{x}}) \rangle . \quad (23)$$



The function  $k$  is called a kernel. The replacement of the usual inner product by kernel is known as the kernel trick and has turned up to be very popular in the machine learning community. Instead of defining a nonlinear map  $\Phi$ , we define a “valid” kernel function  $k(x, z)$ . E.g., polynomial relationships can be modeled via kernels of the form

$$k_d(x, z) = (1 + \langle x, z \rangle)^d, d \in \mathbb{N}.$$

Furthermore, it is possible to define kernels for complex data structures as graphs or text. Literature on the kernel trick and its applications is abundant. A detailed treatise of the subject can be found in Schölkopf & Smola (2002). A nonlinear version of PLS using the kernel trick is presented in Rosipal & Trejo (2001).

If we represent penalized PLS in terms of the kernel matrix  $\mathbf{K}_M$ , we realize that penalized PLS is closely connected to the kernel trick in other respects. Using algorithm 3 or the definition of the kernel matrix  $\mathbf{K}_M$ , we realize that penalized PLS equals ordinary PLS with the canonical inner product replaced by the inner product

$$\langle \mathbf{x}, \mathbf{z} \rangle_M = \mathbf{x}^T \mathbf{M} \mathbf{z}.$$

This function is called a linear kernel. Why is this a sensible inner product? Let us consider the eigendecomposition of the penalty matrix,  $\mathbf{P} = \mathbf{S} \mathbf{\Theta} \mathbf{S}^T$ . We prefer direction  $\mathbf{s}$  such that  $\mathbf{s}^T \mathbf{P} \mathbf{s}$  is small, that is we prefer directions that are defined by eigenvectors  $\mathbf{s}_i$  of  $\mathbf{P}$  with a small corresponding eigenvalue  $\theta_i$ . If we represent the vectors  $\mathbf{x}$  and  $\mathbf{z}$  in terms of the eigenvectors of  $\mathbf{P}$ ,

$$\tilde{\mathbf{x}} = \mathbf{S}^T \mathbf{x} \quad , \quad \tilde{\mathbf{z}} = \mathbf{S}^T \mathbf{z},$$

we conclude that

$$\langle \mathbf{x}, \mathbf{z} \rangle_M = \tilde{\mathbf{x}}^T (\mathbf{I}_p + \mathbf{\Theta})^{-1} \tilde{\mathbf{z}} = \sum_{i=1}^p \frac{1}{1 + \theta_i} \tilde{\mathbf{x}}_i \tilde{\mathbf{z}}_i$$

This implies that directions  $\mathbf{s}_i$  with a small eigenvalue  $\theta_i$  receive a higher weighting than directions with a large eigenvalue.

## 5 Example: Birth Data

In this section, we analyze a real data set describing pregnancy and delivery for 42 infants who are sent to a neonatal intensive care unit after birth. The data are taken from the R (R Development Core Team 2005) software package `exactmaxsel` and are introduced in Boulesteix (2006). Our goal is to predict the number of days spent in the neonatal intensive care unit ( $y$ ) based on the following predictors: birth weight (in g), birth height (in cm), head circumference (in cm), term (in week), age of the mother (in year), weight of the mother before pregnancy (in kg), weight of the mother before delivery (in kg), height of the mother (in cm), time (in month). Some of the predictors are expected to be strongly associated with the response (e.g., birth weight, term), in contrast to poor predictors like

time or height of the mother.

The parameter settings are as follows. We make the simplifying assumption that  $\lambda = \lambda_1 = \dots = \lambda_p$ , which reduces the problem of selecting the optimal smoothing parameter to a one-dimensional problem. As already mentioned above, we use cubic splines. Furthermore, the order of difference of adjacent weights is set to 2. The shape of the fitted functions

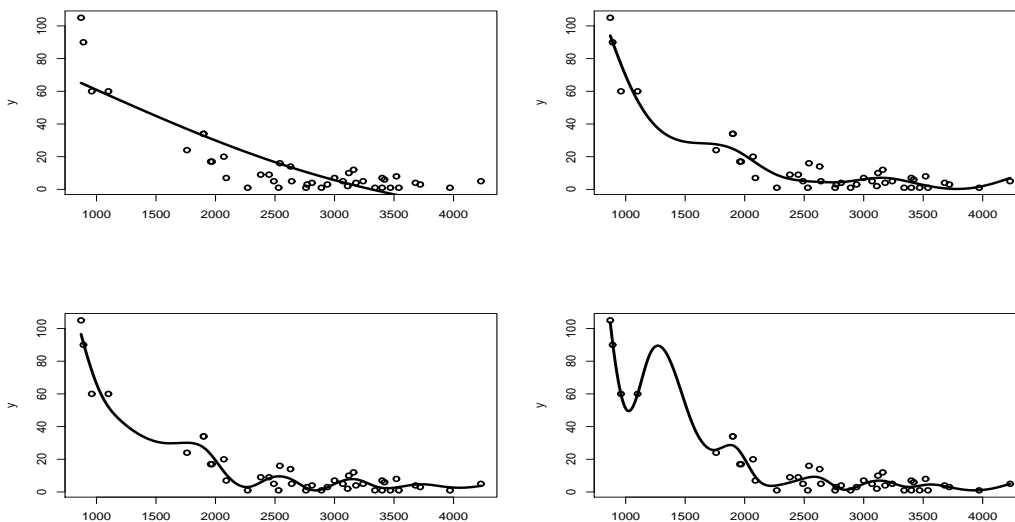


Figure 1: Fitted function for the predictor variable “weight” using penalized PLS. The value of  $\lambda$  is 2000 and the numbers of components are 1, 5 (top) and 9, 13 (bottom).

$f_j$  depends on the two model parameters  $\lambda$  and  $m$ . We first illustrate that the number  $m$  of penalized PLS components controls the smoothness of the estimated functions. For this purpose, we only consider the predictor variable “weight”. Figure 1 displays the fitted functions obtained by penalized PLS for  $\lambda = 2000$  and 4 different numbers of components  $m = 1, 5, 9, 13$ . For small values of  $m$ , the obtained functions are smooth. For higher values of  $m$ , the functions adapt themselves more and more to the data which leads to overfitting for high values of  $m$ .

We compare our novel method to PLS without penalization as described in (Durand 2001) and the `gam()` package in R. This is the implementation of an adaptive selection procedure for the basis functions in (2). More details can be found in Wood (2000) and Wood (2006). This is the standard tool for estimating generalized additive models. The optimal parameter values of (penalized) PLS are determined by computing the leave-one-out squared error. We remark that the split into training and test set is done before transforming the original predictors using B-splines. In order to have comparable results, we normalize the response such that  $\text{var}(\mathbf{y}) = 1$ . The results are summarized in Table 1. Penalized PLS is the best out of the three method. In particular, it receives a considerably lower error than PLS without penalization.

	leave-one-out-error	$m_{opt}$	$\lambda_{opt}$
PLS	0.159	8	–
penalized PLS	<b>0.090</b>	2	330
GAM	0.115	–	–

Table 1: Optimal model parameters and leave-one-out error for the birth data set and normalized response.

## 6 Example:Polymer Data

This data set consists of  $p = 10$  predictor variables and four response variables. The number of observations is  $n = 61$ . The data are taken from a polymer test plant. It can be downloaded from <ftp://ftp.cis.upenn.edu/pub/ungar/chemdata/>. The predictor variables are measurements of controlled variables in a polymer processing plant (e.g. temperatures, feed rates ...). No more details on the variables are given due to confidentiality reasons. As in the last section, we first scale each response variable to have a variance equal to 1. Again, we compare penalized PLS to PLS and `gam()`. The results are summarized in Table 2. For all four response variables, penalized PLS is better than PLS

	1 <sup>st</sup> response	2 <sup>nd</sup> response	3 <sup>rd</sup> response	4 <sup>th</sup> response
PLS	0.672	0.863	0.254	0.204
penalized PLS	0.607	<b>0.801</b>	<b>0.206</b>	<b>0.164</b>
GAM	<b>0.599</b>	0.881	0.218	0.182

Table 2: Leave-one-out error for the polymer data set and normalized response.

without penalization. Penalized PLS is also better than GAM for three out of the four response variables, although the difference is considerably smaller.

## 7 Concluding Remarks

In this work, we proposed an extension of Partial Least Squares Regression using penalization techniques. Apart from its computational efficiency (it is virtually as fast as PLS), it also shares a lot of mathematical properties of PLS. Our novel method obtains good results in applications. In the two examples that are discussed, penalized PLS clearly outperforms PLS without penalization. Furthermore, the results indicate that it is a competitor of `gam()` in the case of very high-dimensional data.

We might think of other penalty terms. Kondylis & Whittaker (2006) consider a preconditioned version of PLS by giving weights to the predictor variables. Higher weights are given to those predictor variables that are highly correlated to the response. These weights can be expressed in terms of a penalty term. Goutis & Fearn (1996) combine PLS with an additive penalty term to data derived from near infra red spectroscopy. The penalty term controls the smoothness of the regression vector.

The introduction of a penalty term can easily be adapted to other dimension reduction techniques. For example for Principal Components Analysis, the penalized optimization criterion is

$$\max_{\mathbf{w}} \frac{\text{var}(\mathbf{X}\mathbf{w})}{\mathbf{w}^T\mathbf{w} + \mathbf{w}^T\mathbf{P}\mathbf{w}}.$$

PLS can handle multivariate responses  $\mathbf{Y}$ . The natural extension of criterion (3) is the following.

$$\max_{\mathbf{w}} \frac{\|\text{cov}(\mathbf{X}\mathbf{w}, \mathbf{Y})\|^2}{\mathbf{w}^T\mathbf{w}} = \max_{\mathbf{w}} \frac{\mathbf{w}^T \mathbf{X}^T \mathbf{Y} \mathbf{Y}^T \mathbf{X} \mathbf{w}}{\mathbf{w}^T \mathbf{w}}.$$

Using Lagrangian multipliers, we deduce that the solution is the eigenvector of the matrix

$$\mathbf{B} = \mathbf{X}^T \mathbf{Y} \mathbf{Y}^T \mathbf{X}$$

that corresponds to the largest eigenvalue of  $\mathbf{B}$ . This eigenvector is usually computed in an iterative fashion. If we want to apply penalized PLS for multivariate responses, we compute

$$\max_{\mathbf{w}} \frac{\mathbf{w}^T \mathbf{X}^T \mathbf{Y} \mathbf{Y}^T \mathbf{X} \mathbf{w}}{\mathbf{w}^T \mathbf{w} + \mathbf{w}^T \mathbf{P} \mathbf{w}}.$$

The solution fulfills

$$\mathbf{B}\mathbf{w} = \gamma (\mathbf{I}_p + \mathbf{P}) \mathbf{w}, \gamma \in \mathbb{R}.$$

This is called a generalized eigenvalue problem or a matrix pencil. Note that for multivariate  $\mathbf{Y}$ , the equivalence of SIMPLS and NIPALS does not hold, so we expect the penalized versions of these methods to be different as well. There are kernel versions for PLS with multivariate  $\mathbf{Y}$  (Rännar et al. 1994, Rosipal & Trejo 2001), hence we can also represent multivariate penalized PLS in terms of kernel matrices.

## Acknowledgement

This research was supported by the Deutsche Forschungsgemeinschaft (SFB 386, ‘‘Statistical Analysis of Discrete Structures’’).

## References

- Boulesteix, A.-L. (2006), ‘Maximally Selected Chi-Square Statistics for Ordinal Variables’, *Biometrical Journal* **48**, 451–462.
- Boulesteix, A.-L. & Strimmer, K. (2006), ‘Partial Least Squares: A Versatile Tool for the Analysis of High-Dimensional Genomic Data’, *Briefings in Bioinformatics*. to appear.
- Butler, N. & Denham, M. (2000), ‘The Peculiar Shrinkage Properties of Partial Least Squares Regression’, *Journal of the Royal Statistical Society B* **62**, 585–593.
- de Boor, C. (1978), *A Practical Guide to Splines*, Springer.

- de Jong, S. (1993), ‘SIMPLS: An Alternative Approach to Partial Least Squares Regression’, *Chemometrics and Intelligent Laboratory Systems* **18**, 251 – 263.
- de Jong, S. (1995), ‘PLS Shrinks’, *Journal of Chemometrics* **9**, 323–326.
- Durand, J. F. (1993), ‘Generalized Principal Component Analysis with Respect to Instrumental Variables via Univariate Spline Transformations’, *Computational Statistics and Data Analysis* **16**, 423–440.
- Durand, J. F. (2001), ‘Local Polynomial Additive Regression Through PLS and Splines: PLSS’, *Chemometrics and Intelligent Laboratory Systems* **58**, 235–246.
- Durand, J. F. & Sabatier, R. (1997), ‘Additive Splines for Partial Least Squares Regression’, *Journal of the American Statistical Association* **92**, 1546–1554.
- Eilers, P. & Marx, B. (1996), ‘Flexible Smoothing with B-Splines and Penalties’, *Statistical Science* **11**, 89–121.
- Frank, I. & Friedman, J. (1993), ‘A Statistical View of some Chemometrics Regression Tools’, *Technometrics* **35**, 109–135.
- Golub, G. & van Loan, C. (1983), *Matrix Computation*, John Hopkins University Press, Baltimore.
- Goutis, C. (1996), ‘Partial Least Squares yields Shrinkage Estimators’, *The Annals of Statistics* **24**, 816–824.
- Goutis, C. & Fearn, T. (1996), ‘Partial Least Squares Regression on Smooth Factors’, *Journal of the American Statistical Association* **91**, 627–632.
- Hastie, T. & Tibshirani, R. (1990), *Generalized Additive Models*, Chapman and Hall.
- Helland, I. (1988), ‘On the Structure of Partial Least Squares Regression’, *Communications in Statistics, Simulation and Computation* **17**(2), 581–607.
- Hestenes, M. & Stiefel, E. (1952), ‘Methods for Conjugate Gradients for Solving Linear Systems’, *Journal of Research of the National Bureau of Standards* **49**, 409–436.
- Kondylis, A. & Whittaker, J. (2006), ‘Preconditioning Krylov Spaces and Variable Selection in PLSR’, *preprint* .
- Lanczos, C. (1950), ‘An Iteration Method for the Solution of the Eigenvalue Problem of Linear Differential and Integral Operators’, *Journal of Research of the National Bureau of Standards* **45**, 225–280.
- Manne, R. (1987), ‘Analysis of Two Partial-Least-Squares Algorithms for Multivariate Calibration’, *Chemometrics and Intelligent Laboratory Systems* **2**, 187–197.
- Martens, H. & Naes, T. (1989), *Multivariate Calibration*, Wiley, New York.

- Phatak, A. & de Hoog, F. (2003), ‘Exploiting the Connection between PLS, Lanczos, and Conjugate Gradients: Alternative Proofs of some Properties of PLS’, *Journal of Chemometrics* **16**, 361–367.
- R Development Core Team (2005), *R: A Language and Environment for Statistical Computing*, R Foundation for Statistical Computing, Vienna, Austria. ISBN 3-900051-07-0.  
**URL:** <http://www.R-project.org>
- Rännar, S., Lindgren, F., Geladi, P. & Wold, S. (1994), ‘A PLS Kernel Algorithm for Data Sets with many Variables and Fewer Objects, Part I: Theory and Applications’, *Journal of Chemometrics* **8**, 111–125.
- Rosipal, R. & Krämer, N. (2006), Overview and Recent Advances in Partial Least Squares, *in* ‘Subspace, Latent Structure and Feature Selection Techniques’, Lecture Notes in Computer Science, Springer, pp. 34–51.
- Rosipal, R. & Trejo, L. (2001), ‘Kernel Partial Least Squares Regression in Reproducing Kernel Hilbert Spaces’, *Journal of Machine Learning Research* **2**, 97–123.
- Rosipal, R., Trejo, L. & Matthews, B. (2003), Kernel PLS-SVC for Linear and Nonlinear Classification, *in* ‘Proceedings of the Twentieth International Conference on Machine Learning’, Washington, DC, pp. 640–647.
- Ruppert, D. (2002), ‘Selecting the Number of Knots for Penalized Splines’, *Journal of Computational and Graphical Statistics* **11**, 735–757.
- Schölkopf, B. & Smola, A. (2002), *Learning with Kernels. Support Vector Machines, Regularization, Optimization, and Beyond.*, The MIT Press.
- Stone, M. & Brooks, R. (1990), ‘Continuum Regression: Cross-Validated Sequentially Constructed Prediction Embracing Ordinary Least Squares, Partial Least Squares and Principal Components Regression (with Discussion)’, *Journal of the Royal Statistical Society B* **52**, 237–269.
- Wold, H. (1975), Path Models with Latent Variables: The NIPALS Approach, *in* H. B. et al., ed., ‘Quantitative Sociology: International Perspectives on Mathematical and Statistical Model Building’, Academic Press, pp. 307–357.
- Wold, S., Ruhe, H., Wold, H. & III, W. D. (1984), ‘The Collinearity Problem in Linear Regression. The Partial Least Squares (PLS) Approach to Generalized Inverses’, *SIAM Journal of Scientific and Statistical Computations* **5**, 735–743.
- Wood, S. (2006), *Generalized Additive Models: An Introduction with R*, Chapman and Hall.
- Wood, S. N. (2000), ‘Modelling and Smoothing Parameter Estimation with Multiple Quadratic Penalties’, *Journal of the Royal Statistical Society B* **62**(2), 413–428.

## A Proofs

We recall that for  $k < i$

$$\mathbf{X}_i = \prod_{j=k}^{i-1} (\mathbf{I}_n - \mathcal{P}_{\mathbf{t}_j}) \mathbf{X}_k = (\mathbf{I}_n - \mathcal{P}_{\mathbf{t}_k, \dots, \mathbf{t}_{i-1}}) \mathbf{X}_k \quad (24)$$

The last equality follows from the fact that the components  $\mathbf{t}_i$  are mutually orthogonal. In particular, we obtain

$$\mathbf{X}_i = (\mathbf{I}_n - \mathcal{P}_{\mathbf{t}_1, \dots, \mathbf{t}_{i-1}}) \mathbf{X}. \quad (25)$$

*Proof of lemma 1.* First note that (25) is equivalent to  $\mathbf{X} = \mathbf{X}_j + \mathcal{P}_{\mathbf{t}_1, \dots, \mathbf{t}_{j-1}} \mathbf{X}$ . It follows that

$$\mathbf{X} \mathbf{w}_j = \mathbf{X}_j \mathbf{w}_j + \mathcal{P}_{\mathbf{t}_1, \dots, \mathbf{t}_{j-1}} \mathbf{X} \mathbf{w}_j = \mathbf{t}_j + \sum_{i=1}^{j-1} \frac{\mathbf{t}_i^T \mathbf{X} \mathbf{w}_j}{\mathbf{t}_i^T \mathbf{t}_i} \mathbf{t}_i. \quad (26)$$

As all components  $\mathbf{t}_i$  are mutually orthogonal,

$$\mathbf{t}_i^T \mathbf{X} \mathbf{w}_j = \begin{cases} \mathbf{t}_i^T \mathbf{t}_i \neq 0 & , i = j \\ 0 & , i > j \\ * & , \text{otherwise} \end{cases}.$$

We conclude that  $\mathbf{R}$  is an upper triangular matrix with all diagonal elements  $\neq 0$ . Furthermore, it follows from (26) that all vectors  $\mathbf{X} \mathbf{w}_j$  are linear combinations of the components  $\mathbf{t}_1, \dots, \mathbf{t}_j$ . This implies that the columns of  $\mathbf{X} \mathbf{W}$  and the columns of  $\mathbf{T}$  span the same space. Finally, we have to show that  $\mathbf{R}$  is bidiagonal. To prove this, we show that  $\mathbf{X}_i \mathbf{w}_j = 0$  for  $j < i$ . The condition  $i > j$  implies (recall (24)) that  $\mathbf{X}_i = \mathbf{X}_j - \mathcal{P}_{\mathbf{t}_j, \dots, \mathbf{t}_{i-1}} \mathbf{X}_j$  and consequently

$$\mathbf{X}_i \mathbf{w}_j = \mathbf{X}_j \mathbf{w}_j - \mathcal{P}_{\mathbf{t}_1, \dots, \mathbf{t}_{i-1}} \mathbf{X}_j \mathbf{w}_j = \mathbf{t}_j - \mathcal{P}_{\mathbf{t}_1, \dots, \mathbf{t}_{i-1}} \mathbf{t}_j \stackrel{j \leq i-1}{=} \mathbf{t}_j - \mathbf{t}_j = \mathbf{0}.$$

This implies that for  $i - 1 > j$

$$\begin{aligned} \mathbf{t}_i^T \mathbf{X} \mathbf{w}_j &= \mathbf{t}_i^T (\mathbf{X}_i + \mathcal{P}_{\mathbf{t}_1, \dots, \mathbf{t}_{i-1}} \mathbf{X}) \mathbf{w}_j \\ &= \mathbf{t}_i^T (\mathbf{X}_i \mathbf{w}_j + \mathcal{P}_{\mathbf{t}_1, \dots, \mathbf{t}_{i-1}} \mathbf{X} \mathbf{w}_j) \\ &= \mathbf{t}_i^T (\mathcal{P}_{\mathbf{t}_1, \dots, \mathbf{t}_{i-1}} \mathbf{X} \mathbf{w}_j) = \mathbf{0}. \end{aligned}$$

□

*Proof of proposition 3.* For  $i = 1$ , we have  $\tilde{\mathbf{w}}_1 = \mathbf{w}_1$  as  $\mathbf{X}_1 = \mathbf{X}$ . For a general  $i$ , we have

$$\mathbf{t}_{i+1} = \mathbf{X}_{i+1} \mathbf{w}_{i+1} = (\mathbf{X} - \mathcal{P}_{\mathbf{t}_1, \dots, \mathbf{t}_i} \mathbf{X}) \mathbf{w}_{i+1} = \mathbf{X} \mathbf{w}_{i+1} - \mathcal{P}_{\mathbf{t}_i} \mathbf{X} \mathbf{w}_{i+1}.$$

The last equality holds as  $\mathbf{R} = \mathbf{T}^T \mathbf{X} \mathbf{W}$  is bidiagonal. Using formula (6) for the projection operator, it follows that

$$\mathbf{t}_{i+1} = \mathbf{X} \mathbf{w}_{i+1} - \mathbf{X} \frac{\tilde{\mathbf{w}}_i \tilde{\mathbf{w}}_i^T}{\tilde{\mathbf{w}}_i^T \mathbf{X}^T \mathbf{X} \tilde{\mathbf{w}}_i} \mathbf{X}^T \mathbf{X} \mathbf{w}_{i+1}.$$

We conclude that

$$\tilde{\mathbf{w}}_{i+1} = \mathbf{w}_{i+1} - \frac{\tilde{\mathbf{w}}_i^T \mathbf{X}^T \mathbf{X} \mathbf{w}_{i+1}}{\tilde{\mathbf{w}}_i^T \mathbf{X}^T \mathbf{X} \tilde{\mathbf{w}}_i} \tilde{\mathbf{w}}_i.$$

The regression estimate after  $i + 1$  steps is

$$\begin{aligned} \mathbf{X} \hat{\boldsymbol{\beta}}^{(i+1)} &= \mathcal{P}_{\mathbf{t}_1, \dots, \mathbf{t}_{i+1}} \mathbf{Y} \\ &= \mathbf{X} \hat{\boldsymbol{\beta}}^{(i)} + \mathcal{P}_{\mathbf{t}_{i+1}} \mathbf{Y} \\ &= \mathbf{X} \hat{\boldsymbol{\beta}}^{(i)} + \mathbf{X} \frac{\tilde{\mathbf{w}}_{i+1} \tilde{\mathbf{w}}_{i+1}^T}{\tilde{\mathbf{w}}_{i+1}^T \mathbf{X}^T \mathbf{X} \tilde{\mathbf{w}}_{i+1}} \mathbf{X}^T \mathbf{Y}. \end{aligned}$$

This concludes the proof.  $\square$

*Proof of lemma 4.* We use induction. For  $m = 1$  we know that  $\mathbf{w}_1 = \mathbf{b}_M$ . For a fixed  $m > 1$ , we conclude from the induction hypothesis and lemma 1 that every vector  $\mathbf{s}$  that lies in the span of  $\mathbf{t}_1, \dots, \mathbf{t}_m$  is of the form

$$\mathbf{s} = \mathbf{X} \mathbf{v} \quad , \quad \mathbf{v} \in \text{span}\{\mathbf{w}_1, \dots, \mathbf{w}_m\} = \mathcal{K}^{(m)}. \quad (27)$$

We conclude that

$$\mathbf{X}_{m+1}^T \mathbf{y} = (\mathbf{X} - \mathcal{P}_{\mathbf{t}_1, \dots, \mathbf{t}_m} \mathbf{X})^T \mathbf{y} = \mathbf{X}^T \mathbf{y} - \mathbf{X}^T \mathcal{P}_{\mathbf{t}_1, \dots, \mathbf{t}_m} \mathbf{y} \stackrel{(27)}{=} \mathbf{b} - \mathbf{X}^T \mathbf{X} \mathbf{s}$$

and that

$$\mathbf{w}_{m+1} = \mathbf{M} \mathbf{X}_{m+1}^T \mathbf{y} = \mathbf{M} \mathbf{b} - \mathbf{M} \mathbf{A} \mathbf{s} = \mathbf{b}_M - \mathbf{A}_M \mathbf{s} \in \mathcal{K}^{(m+1)}.$$

$\square$

In the rest of the appendix, we show the equivalence of penalized PLS and the preconditioned conjugate gradient method.

**Lemma 8.** *We have*

$$\text{span}\{\mathbf{d}_0, \dots, \mathbf{d}_{m-1}\} = \text{span}\{\mathbf{r}_0, \dots, \mathbf{r}_{m-1}\} = \text{span}\{\mathbf{x}_1, \dots, \mathbf{x}_m\} = \mathcal{K}^{(m)}.$$

This can be proven via induction.

**Lemma 9.** *We have*

$$\boldsymbol{\beta}_m = \sum_{i=0}^{m-1} \frac{\langle \mathbf{d}_i, \mathbf{b}_M \rangle_{M^{-1}}}{\langle \mathbf{d}_i, \mathbf{A}_M \mathbf{d}_i \rangle_{M^{-1}}} \mathbf{d}_i$$



*Proof.* This corresponds to the iterative definition of  $\boldsymbol{\beta}_{m+1}$ . We only have to show that

$$\langle \mathbf{d}_i, \mathbf{r}_i \rangle_{M^{-1}} = \langle \mathbf{d}_i, \mathbf{b}_M \rangle_{M^{-1}} .$$

Note that

$$\mathbf{r}_i = \mathbf{b} - \sum_{j=0}^{i-1} a_j \mathbf{A}_M \mathbf{d}_j .$$

As  $\mathbf{d}_i$  is  $\mathbf{A}_M$ -orthogonal onto all directions  $\mathbf{d}_j$ ,  $j < i$ , the proof is complete.  $\square$

Now we are able to proof the equivalence of penalized PLS and the conjugate gradient method.

*Proof of theorem 6.* As the search directions  $\mathbf{d}_i$  span the Krylov space  $\mathcal{K}^{(m)}$ , we can replace the matrix  $\mathbf{W}$  in (14) by the matrix  $\mathbf{D} = (\mathbf{d}_0, \dots, \mathbf{d}_{m-1})$ . As the search directions are  $\mathbf{A}_M$ -orthogonal, we have

$$\begin{aligned} \widehat{\boldsymbol{\beta}}_{PPLS} &= \mathbf{D} (\mathbf{D}^T \mathbf{A} \mathbf{D})^{-1} \mathbf{D}^T \mathbf{b} \\ &= \mathbf{D} (\mathbf{D}^T \mathbf{M}^{-1} \mathbf{A}_M \mathbf{D})^{-1} \mathbf{D}^T \mathbf{M}^{-1} \mathbf{b}_M \\ &= \sum_{i=0}^{m-1} \frac{\langle \mathbf{d}_i, \mathbf{b}_M \rangle_{M^{-1}}}{\langle \mathbf{d}_i, \mathbf{A}_M \mathbf{d}_i \rangle_{M^{-1}}} \mathbf{d}_i \end{aligned}$$

and this equals the formula in lemma 9.  $\square$