
Kernelizing PLS, Degrees of Freedom, and Efficient Model Selection

Nicole Krämer
Mikio L. Braun

NKRAEMER@CS.TU-BERLIN.DE
MIKIO@CS.TU-BERLIN.DE

Technische Universität Berlin, Franklinstr. 28/29, 10587 Berlin, Germany

Abstract

Kernelizing partial least squares (PLS), an algorithm which has been particularly popular in chemometrics, leads to kernel PLS which has several interesting properties, including a sub-cubic runtime for learning, and an iterative construction of directions which are relevant for predicting the outputs. We show that the kernelization of PLS introduces interesting properties not found in ordinary PLS, giving novel insights into the workings of kernel PLS and the connections to kernel ridge regression and conjugate gradient descent methods. Furthermore, we show how to correctly define the degrees of freedom for kernel PLS and how to efficiently compute an unbiased estimate. Finally, we address the practical problem of model selection. We demonstrate how to use the degrees of freedom estimate to perform effective model selection, and discuss how to implement cross-validation schemes efficiently.

1. Introduction

Kernel methods have proven to work extremely well for a large range of machine learning applications, and a multitude of different variants exist, ranging from support vector machines (SVMs) to kernel ridge regression (KRR). Each of these algorithms comes with its own strengths and weaknesses. For example, SVMs produce sparse results permitting fast predictions, while the complexity of implementing SVMs is significantly larger than, for example, KRR. On the other hand, KRR produces non-sparse solutions leading to scalability problems.

On the algorithmic-complexity-vs.-scalability axis, kernel PLS inhabits an interesting middle position: Kernel PLS can be implemented in a simple iterative

scheme whose runtime scales sub-cubic in the number of training examples. In each iteration, only a matrix-vector multiplication is required. In this respect, kernel PLS is similar to using conjugate gradient type methods for solving KRRs, but as we argue below, kernel PLS usually converges to a good solution using fewer steps.

Kernel PLS therefore seems to be a promising addition to the toolbox of existing kernel methods. However, some ingredients for making kernel PLS readily available are still missing, and the goal of this paper is to address some of these issues.

On the theoretical side, we discuss several properties of kernel PLS not found in ordinary PLS, which help to explain how kernel PLS works. We illustrate that for kernels which can be interpreted as smoothing operators (e.g., Gaussian kernels), kernel PLS can be interpreted as an iterative refinement using smoothed residuals.

For ordinary PLS, there is a close relationship to conjugate gradient (CG) methods, namely that the PLS solutions are equal to the conjugate gradient approximations applied to the normal equation. We show that this connection is less strong between kernel PLS and KRR. In fact, it turns out that kernel PLS and CG for KRR optimize on the same subspaces, but the objective function used by kernel PLS is more suited to regression problems leading to faster convergence of the solutions.

Finally, for regression methods, the notion of degrees of freedom characterizes the complexity of the model well. We discuss how to define the degrees of freedom correctly for kernel PLS and propose an efficient algorithm for computing an unbiased estimate.

On the practical side, the issue of efficient and effective model selection is of prime importance for applications. Using the degrees of freedom estimate, we perform model selection in a minimum description length framework. For cross-validation schemes, we discuss how to efficiently implement testing across different number of components.

Appearing in *Proceedings of the 24th International Conference on Machine Learning*, Corvallis, OR, 2007. Copyright 2007 by the author(s)/owner(s).

2. Kernelizing PLS

We briefly introduce the ordinary PLS algorithm and derive a variant of kernelized PLS. We start with defining some notation. As usual, we consider a learning problem with inputs (or predictors) \mathbf{x}_i and outputs y_i . The data matrix \mathbf{X} is the matrix whose *rows* are the centered \mathbf{x}_i . The vector \mathbf{y} consists of the centered outputs y_i . The number of observations is denoted by n , the number of variables is denoted by d . The number of PLS components is m . For any set of linearly independent vectors $\mathbf{S} = \mathbf{s}_1, \dots, \mathbf{s}_k$, we denote by $\mathcal{P}_{\mathbf{S}}$ the projection onto the space that is spanned by the vectors \mathbf{s}_i . In matrix notation, we have $\mathcal{P} = \mathbf{S}(\mathbf{S}^T \mathbf{S})^{-1} \mathbf{S}^T$. The kernel matrix $\mathbf{X} \mathbf{X}^T$ is denoted by \mathbf{K} . Furthermore, we set $\mathbf{A} = \mathbf{X}^T \mathbf{X}$ and $\mathbf{b} = \mathbf{X}^T \mathbf{y}$. In the linear case, we consider the linear regression model $\mathbf{y} = \mathbf{X} \boldsymbol{\beta} + \mathbf{e}$. We only discuss PLS for a one-dimensional response, although PLS can also handle multivariate responses. As we illustrate below, PLS is closely related to conjugate gradients and Krylov methods. On this account, we recall the definition of Krylov subspaces. (For more details on Krylov spaces see e.g. Golub & Van Loan, 1996) For a matrix $\mathbf{C} \in \mathbb{R}^{c \times c}$ and $\mathbf{c} \in \mathbb{R}^c$, we call the set of vectors $\mathbf{c}, \mathbf{C}\mathbf{c}, \dots, \mathbf{C}^{m-1}\mathbf{c}$ the Krylov sequence of length m . The space spanned by these vectors is called a Krylov space and is denoted by

$$\mathcal{K}_m(\mathbf{C}, \mathbf{c}) = \text{span} \{ \mathbf{c}, \mathbf{C}\mathbf{c}, \dots, \mathbf{C}^{m-1}\mathbf{c} \}.$$

2.1. Ordinary Partial Least Squares

The main idea of PLS is to build a few orthogonal components $\mathbf{T} = (\mathbf{t}_1, \dots, \mathbf{t}_m)$ from the original predictors \mathbf{X} and to use them in a least squares regression in place of \mathbf{X} . 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} .

A latent component \mathbf{t} is a linear combination $\mathbf{t} = \mathbf{X}\mathbf{w}$ of the predictor variables. The vector \mathbf{w} is called the weight vector. We want to find a component with maximal covariance to \mathbf{y} , that is, for the first component $\mathbf{t}_1 = \mathbf{X}\mathbf{w}_1$, we maximize

$$\mathbf{w}_1 = \underset{\mathbf{w} \in \mathbb{R}^d}{\text{argmax}} \frac{\mathbf{w}^T \mathbf{X}^T \mathbf{y} \mathbf{y}^T \mathbf{X} \mathbf{w}}{\mathbf{w}^T \mathbf{w}} = \mathbf{X}^T \mathbf{y}. \quad (1)$$

Subsequent components $\mathbf{t}_2, \mathbf{t}_3, \dots$ are chosen such that they maximize the squared covariance to \mathbf{y} and that all components are mutually orthogonal. Orthogonality can be ensured by deflating the original predictor variables \mathbf{X} . That is, we only consider the part of \mathbf{X}

that is orthogonal on all components $\mathbf{t}_j, j < i$:

$$\mathbf{X}_i = \mathbf{X} - \mathcal{P}_{\mathbf{t}_1, \dots, \mathbf{t}_{i-1}} \mathbf{X}. \quad (2)$$

We then replace \mathbf{X} by \mathbf{X}_i in (1). This is called the NIPALS method (Wold, 1975). We remark that there are different techniques to extract subsequent components. More information on different versions of PLS can be found in Rosipal and Krämer (2006).

In order to predict the output for new observations, we have to determine the vector of regression coefficients $\hat{\boldsymbol{\beta}}_m$ defined by

$$\hat{\mathbf{y}}_m = \mathcal{P}_{\mathbf{t}_1, \dots, \mathbf{t}_m} \mathbf{y} = \mathbf{X} \hat{\boldsymbol{\beta}}_m.$$

Therefore, we need a representation of the components in terms of the original data, i.e.

$$\mathbf{t}_i = \mathbf{X}_i \mathbf{w}_i = \mathbf{X} \tilde{\mathbf{w}}_i. \quad (3)$$

This can be done by exploiting the fact (Manne, 1987) that the matrix $\tilde{\mathbf{L}} = \mathbf{T}^T \mathbf{X} \mathbf{W}$ is upper bidiagonal, i.e. $\tilde{l}_{ij} = 0$ for $i > j$ or $i < j - 1$. Furthermore, setting $\mathbf{D} = \text{diag}(1/\|\mathbf{t}_1\|, \dots, 1/\|\mathbf{t}_m\|)$, we have

$$\mathbf{X} \mathbf{W} = (\mathbf{T} \mathbf{D}) (\mathbf{D} \tilde{\mathbf{L}}). \quad (4)$$

In particular, the columns of \mathbf{T} and the columns of $\mathbf{X} \mathbf{W}$ span the same space. Note that (4) is in fact the QR-factorization of $\mathbf{X} \mathbf{W}$. Below, we derive an efficient algorithm for both kernel PLS and its degrees of freedom based on this relationship. We conclude this section by remarking that (Helland, 1988)

$$\begin{aligned} \mathcal{K}_m(\mathbf{A}, \mathbf{b}) &= \text{span}(\mathbf{w}_1, \dots, \mathbf{w}_m), \quad (5) \\ \mathcal{K}_m(\mathbf{K}, \mathbf{K} \mathbf{y}) &= \text{span}(\mathbf{t}_1, \dots, \mathbf{t}_m). \quad (6) \end{aligned}$$

which already points at the close relationship of PLS and Krylov spaces.

2.2. Kernel Partial Least Squares

A derivation of PLS in terms of the kernel matrix $\mathbf{K} = \mathbf{X} \mathbf{X}^T$ and \mathbf{y} is already defined in Rännar et al. (1994) in order to speed up the computation of PLS. The extension of PLS to nonlinear regression and classification using the kernel trick is proposed in Rosipal and Trejo (2001) and Rosipal et al. (2003). There, a more general setting with multivariate output is considered. For the univariate case, the derivation can however be simplified. There are different equivalent possibilities to define the PLS solution in terms of the Kernel matrix. All of them are based on (6) which implies that

$$\hat{\mathbf{y}}_m = \mathcal{P}_{\mathbf{T}} \mathbf{y} = \mathcal{P}_{\mathcal{K}_m(\mathbf{K}, \mathbf{K} \mathbf{y})} \mathbf{y} = \mathbf{K} \boldsymbol{\alpha}_m. \quad (7)$$

Here, α_m are the coefficients in the kernel expansion

$$f(\mathbf{x}) = \sum_{i=1}^n \alpha_i k(\mathbf{x}, \mathbf{x}_i)$$

for PLS with m components. In order to compute α_m , we start by defining an ordinary basis $\tilde{\mathbf{T}} = (\tilde{\mathbf{t}}_1, \dots, \tilde{\mathbf{t}}_m)$ of $\mathcal{K}_m(\mathbf{K}, \mathbf{K}\mathbf{y})$. E.g. we might simply choose $\tilde{\mathbf{t}}_i = \mathbf{K}^i \mathbf{y}$. In order to determine α_m we first represent the basis $\tilde{\mathbf{T}}$ as $\tilde{\mathbf{T}} = \mathbf{K}\mathbf{R}$. Let us denote the QR-factorization of $\tilde{\mathbf{T}}$ by $\tilde{\mathbf{T}} = \mathbf{T}\mathbf{L}$. It follows that

$$\mathbf{T} = \tilde{\mathbf{T}}\mathbf{L}^{-1} = \mathbf{K}\mathbf{R}\mathbf{L}^{-1}. \quad (8)$$

Combining (7) and (8), we obtain

$$\hat{\mathbf{y}}_m = \mathbf{T}\mathbf{T}^T \mathbf{y} = \mathbf{K}\mathbf{R}\mathbf{L}^{-1} \mathbf{T}^T \mathbf{y}.$$

It follows that

$$\alpha_m = \mathbf{R}\mathbf{L}^{-1} \mathbf{T}^T \mathbf{y}.$$

Although a priori any basis $\tilde{\mathbf{T}}$ will do, it is crucial to derive an algorithm that is computationally efficient. This can be achieved by exploiting (4). More precisely, we choose $\tilde{\mathbf{T}} = \mathbf{X}\mathbf{W}$ and derive the QR-factorization iteratively by exploiting that $\tilde{\mathbf{L}}$ is bidiagonal. First, a kernel representation of $\tilde{\mathbf{T}} = \mathbf{K}\mathbf{R}$ can be derived by exploiting (2), which implies that $\mathbf{w}_i = \mathbf{X}^T \mathbf{r}_i$ with $\mathbf{r}_i = \mathbf{y} - \hat{\mathbf{y}}_i$. We set $\tilde{\mathbf{t}}_i = \mathbf{X}\mathbf{w}_i = \mathbf{K}\mathbf{r}_i$, and derive a recursive formula for the expansion (recall (3))

$$\mathbf{t}_i = \mathbf{X}\tilde{\mathbf{w}}_i =: \mathbf{K}\gamma_i.$$

Note that by definition, $\mathbf{t}_i = (\mathbf{X} - \mathcal{P}_{\mathbf{t}_1, \dots, \mathbf{t}_{i-1}}) \mathbf{w}_i$. As $\tilde{\mathbf{L}} = \mathbf{T}^T \tilde{\mathbf{T}}$ is upper bidiagonal, we have

$$\mathbf{t}_i = \tilde{\mathbf{t}}_i - \mathcal{P}_{\mathbf{t}_{i-1}} \tilde{\mathbf{t}}_i,$$

from which we conclude that

$$\gamma_i = \mathbf{r}_i - \frac{\mathbf{t}_{i-1}^T \tilde{\mathbf{t}}_i}{\mathbf{t}_{i-1}^T \mathbf{t}_{i-1}} \gamma_{i-1}.$$

Finally, we have

$$\hat{\mathbf{y}}_i = \hat{\mathbf{y}}_{i-1} + \mathcal{P}_{\mathbf{t}_i} \mathbf{y} = \mathbf{K}\alpha_{i-1} + \frac{\mathbf{t}_i^T \mathbf{y}}{\mathbf{t}_i^T \mathbf{t}_i} \mathbf{K}\gamma_i.$$

From this, we derive a recursive formula for α_i . The results are summarized in algorithm 1. To ensure numeric stability, we include the stopping criterium in line 13.

Note that the formulas in steps 9, 11 and 15 in algorithm 1 are redundant. E.g. for the computation of \mathbf{t}_i in step 11, it is computationally more efficient to use the relationship $\mathbf{t}_i = \mathbf{K}\gamma_i$. However, for the interpretation of Kernel PLS as an iterative refinement (see next subsection) and for the estimation of its degrees of freedom, the recursive formula $\mathbf{t}_i = \mathbf{K}\mathbf{r}_i - \mathcal{P}_{\mathbf{t}_{i-1}} \mathbf{K}\mathbf{r}_i$ is more convenient.

Algorithm 1 Iterative Computation of kernel PLS

```

1: Input:  $\mathbf{X}$ ,  $\mathbf{y}$ , number of components  $m$ 
2: Initialize:  $\mathbf{K} = \mathbf{X}\mathbf{X}^T$ ,  $\alpha_0 = \hat{\mathbf{y}}_0 = \mathbf{0}$ 
3: for  $i = 1, \dots, m$  do
4:   if  $i = 1$  then
5:      $\mathbf{r}_i = \gamma_i = \mathbf{y}$ 
6:      $\tilde{\mathbf{t}}_i = \mathbf{t}_i = \mathbf{K}\mathbf{y}$ 
7:   else
8:      $\mathbf{r}_i = \mathbf{r}_{i-1} - \mathcal{P}_{\mathbf{t}_{i-1}} \mathbf{y}$ 
9:      $\tilde{\mathbf{t}}_i = \mathbf{K}\mathbf{r}_i [= \tilde{\mathbf{t}}_{i-1} - \mathbf{K}\mathcal{P}_{\mathbf{t}_{i-1}} \mathbf{y}]$ 
10:     $\gamma_i = \mathbf{r}_i - \frac{\mathbf{t}_{i-1}^T \tilde{\mathbf{t}}_i}{\mathbf{t}_{i-1}^T \mathbf{t}_{i-1}} \gamma_{i-1}$ 
11:     $\mathbf{t}_i = \mathbf{K}\gamma_i [= \mathbf{K}\mathbf{r}_i - \mathcal{P}_{\mathbf{t}_{i-1}} \mathbf{K}\mathbf{r}_i]$ 
12:  end if
13:  exit if  $|\angle(\mathbf{t}_i, \mathbf{t}_j)| > 0.1$  for any  $1 \leq j < i$ 
14:   $\alpha_i = \alpha_{i-1} + \frac{\mathbf{t}_i^T \mathbf{y}}{\mathbf{t}_i^T \mathbf{t}_i} \gamma_i$ 
15:   $\hat{\mathbf{y}}_i = \hat{\mathbf{y}}_{i-1} + \mathcal{P}_{\mathbf{t}_i} \mathbf{y} [= \mathbf{K}\alpha_i]$ 
16: end for
    
```

2.3. Kernel PLS as Iterative Refinement

Kernel PLS can also be interpreted as an iterative refinement of $\hat{\mathbf{y}}$ using smoothed residuals. In order to compute $\hat{\mathbf{y}}_i$, the i th residual $\mathbf{r}_i = \mathbf{y} - \hat{\mathbf{y}}_{i-1}$ is computed first (line 8). The residual \mathbf{r}_i is then mapped through \mathbf{K} , and \mathbf{t}_i is constructed by orthogonalizing $\mathbf{K}\mathbf{r}_i$ with respect to \mathbf{t}_{i-1} (line 11, due to the bidiagonality, \mathbf{t}_i is then also orthogonal to all $\mathbf{t}_1, \dots, \mathbf{t}_{i-1}$). Finally, $\hat{\mathbf{y}}_i$ is updated by adding the projection of \mathbf{y} to \mathbf{t}_i (line 15).

The key ingredient here is the application of \mathbf{K} to \mathbf{r}_i . A component of \mathbf{r}_i along the eigenvector \mathbf{u}_i of \mathbf{K} is scaled by the corresponding eigenvalue λ_i . Since smooth kernels, as are typically employed in machine learning, have rapidly decaying eigenvalues, most of these components will be effectively suppressed, and \mathbf{t}_i will reflect only the part of the residual corresponding to directions with large eigenvalues. It turns out that this is a very effective strategy for learning. As discussed by Braun et al. (2007), in a kernel setting, the relevant information about the outputs \mathbf{y} is contained in a number of leading eigendirections. Kernel PLS iteratively refines $\hat{\mathbf{y}}$ to minimize the residual on this subspace, thereby extracting the relevant information and discarding the noise.

2.4. Kernel PLS vs. Conjugate Gradients

It is well-known that PLS is closely connected to Krylov subspaces and conjugate gradient (CG) methods (Lingjærde & Christophersen, 2000; Phatak & de Hoog, 2002). The connection is given as follows. The ordinary least-squares problem $\min_{\beta} \|\mathbf{y} - \mathbf{X}\beta\|^2$

is solved by finding the solution of the associated normal equation (recall that $\mathbf{A} = \mathbf{X}\mathbf{X}^T$, $\mathbf{b} = \mathbf{X}^T\mathbf{y}$)

$$\mathbf{A}\boldsymbol{\beta} = \mathbf{b}. \quad (9)$$

One can show that applying the CG method to this equation starting with the initialization $\boldsymbol{\beta}_0 = 0$ produces the same solution as PLS (Helland, 1988; Phatak et al., 2002). The CG method iteratively computes approximate solutions of (9) by minimizing the quadratic function

$$\phi(\boldsymbol{\beta}) = \frac{1}{2}\boldsymbol{\beta}^T \mathbf{A}\boldsymbol{\beta} - \boldsymbol{\beta}^T \mathbf{b}. \quad (10)$$

along directions that are \mathbf{A} -orthogonal. The CG algorithm is also closely related to Krylov subspaces and the Lanczos algorithm (Lanczos, 1950). One can show that the CG algorithm minimizes (10) over the Krylov spaces. It follows from (4) that the PLS estimate $\hat{\boldsymbol{\beta}}_m$ obtained after m steps is obtained by minimizing ϕ over $\mathcal{K}_m(\mathbf{A}, \mathbf{b})$

$$\begin{aligned} \hat{\boldsymbol{\beta}}_m &= \underset{\boldsymbol{\beta} \in \mathcal{K}_m(\mathbf{A}, \mathbf{b})}{\operatorname{argmin}} \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|^2 \\ &= \underset{\boldsymbol{\beta} \in \mathcal{K}_m(\mathbf{A}, \mathbf{b})}{\operatorname{argmin}} \left(\frac{1}{2}\boldsymbol{\beta}^T \mathbf{A}\boldsymbol{\beta} - \boldsymbol{\beta}^T \mathbf{b} \right). \end{aligned} \quad (11)$$

This connection is remarkable since Krylov spaces are also known to allow good approximations of leading eigenvalues and eigenspaces in the context of Lanczos methods, such that PLS can be interpreted of making use of low-rank approximations to \mathbf{K} in a similar spirit as principal component regression, which approximates \mathbf{K} by restricting \mathbf{K} to the exact leading eigenspaces.

In the context of kernel PLS, a similar result concerning the Krylov spaces holds. Using (6), we conclude that the m th kernel PLS solution can be written as

$$\begin{aligned} \boldsymbol{\alpha}_m &= \underset{\boldsymbol{\alpha} \in \mathcal{K}_m(\mathbf{K}, \mathbf{y})}{\operatorname{argmin}} \|\mathbf{y} - \mathbf{K}\boldsymbol{\alpha}\|^2 \\ &= \underset{\boldsymbol{\alpha} \in \mathcal{K}_m(\mathbf{K}, \mathbf{y})}{\operatorname{argmin}} \left(\frac{1}{2}\boldsymbol{\alpha}^T \mathbf{K}^2 \boldsymbol{\alpha} - \boldsymbol{\alpha}^T \mathbf{K}\mathbf{y} \right) \end{aligned} \quad (12)$$

While this term looks very similar to conjugate gradient applied to the equation $\mathbf{K}^2\boldsymbol{\alpha} = \mathbf{K}\mathbf{y}$, note that the Krylov-space is different, since it should read $\mathcal{K}_m(\mathbf{K}^2, \mathbf{K}\mathbf{y})$ for CG.

Using the conjugate gradients approach to the solution of kernel ridge regression (KRR) leads to optimization over the same Krylov spaces, however with a different objective function. In KRR, we have to solve

$$(\mathbf{K} + \lambda\mathbf{I})\boldsymbol{\alpha} = \mathbf{y},$$

for some $\lambda > 0$. Since explicitly inverting the matrix is impractical for large numbers of examples, one alternative consists in solving this equation iteratively using the CG method. This often leads to a significant reduction in computation time since one iteration only requires a matrix-vector operation (and additional vector-vector operations). Furthermore, the CG method is guaranteed to converge to the optimal solution after n steps, but usually converges much faster.

Now if we compare the CG cost function for the KRR problem, we get that (with $\mathbf{K}_\lambda = \mathbf{K} + \lambda\mathbf{I}$)

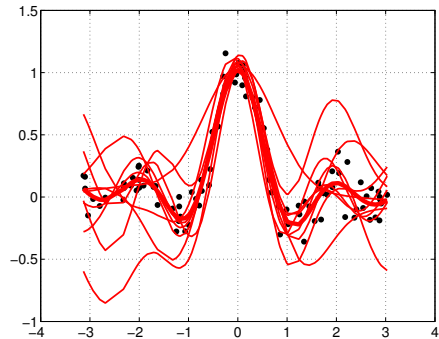
$$\begin{aligned} \phi_{\text{KRR}}(\boldsymbol{\alpha}) &= \frac{1}{2}\boldsymbol{\alpha}^T \mathbf{K}_\lambda \boldsymbol{\alpha} - \boldsymbol{\alpha}^T \mathbf{y} \\ &= \frac{1}{2}(\mathbf{K}_\lambda \boldsymbol{\alpha} - \mathbf{y})^T \mathbf{K}_\lambda^{-1} (\mathbf{K}_\lambda \boldsymbol{\alpha} - \mathbf{y}) + \text{const.} \end{aligned}$$

Comparing ϕ_{KRR} with (12), we see that both methods minimize the sum of squared residual, however, CG applied to KRR (CG-KRR) rescales the errors by \mathbf{K}_λ^{-1} , leading to an amplification of the error along directions of small eigenvalues of \mathbf{K} of up to $1/\lambda$, emphasizing some components of the residual. In other words, the intermediate solutions of CG applied to KRR tries to find a solution with higher accuracy than actually specified by the optimization problem.

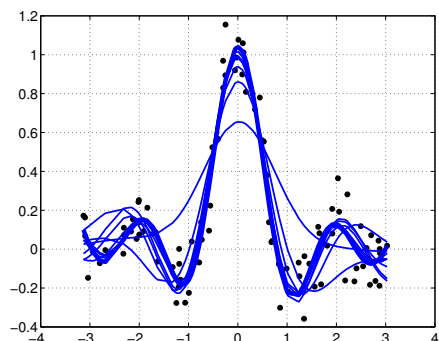
Practically, this results in larger errors and more iterations necessary to obtain a good solution compared to kernel PLS. Consider the following example. We draw 100 inputs X_i uniformly from $[-\pi, \pi]$, and $Y_i = \text{sinc}(X_i) + \varepsilon_i$ with ε_i being normally distributed and having standard deviation $1/10$. We choose an rbf-kernel of width $1/2$ and set $\lambda = 10^{-3}$. Both kernel PLS and CG-KRR find good solutions after 10 iterations. Comparing the resulting fit functions for CG-KRR and kernel PLS (Figures 1(a) and 1(b)), one can see that the kernel PLS solutions converge faster. This is also reflected by the sum of squared residuals (Figure 1(c)). The observations on the different objective functions is also reflected by Figure 1(d), which plots the difference between $\boldsymbol{\alpha}_5$ and $\boldsymbol{\alpha}_{10}$ decomposed by eigenvector directions of \mathbf{K} . One can see that kernel PLS is more accurate in the leading eigenvectors, and less accurate on the directions corresponding to small eigenvectors, while the CG-KRR solutions distributes the error more evenly. By mapping $\boldsymbol{\alpha}$ through \mathbf{K} , any information contained in directions corresponding to small eigenvalues is discarded anyway, such that kernel PLS leads to better fits.

3. Degrees of Freedom

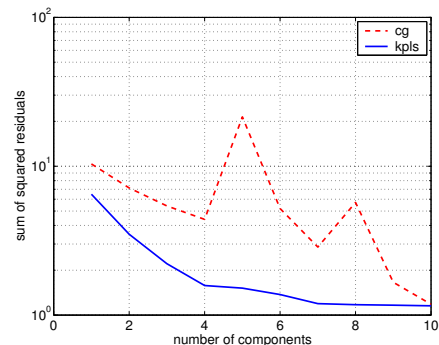
In kernel PLS, both the number of PLS components and the kernel parameters have to be estimated.



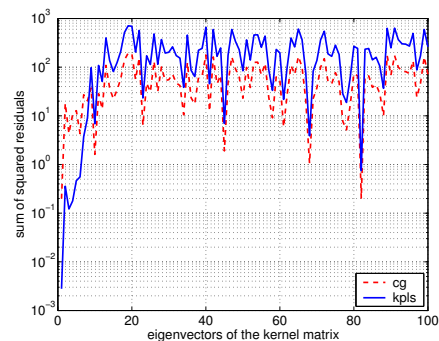
(a) CGs applied to KRR



(b) Kernel PLS



(c) sum of squared residuals



(d) residual by eigendirections

Figure 1. Comparing kernel PLS to conjugate gradients (CG) applied to kernel ridge regression (KRR).

Roughly, we can distinguish between two approaches for model selection.

The k -fold cross-validation method is based on a repeated random splitting of the sample into training and test data. As kernel PLS has a sub-cubic runtime for learning, cross-validation is very efficient as long as k is not too large. The runtime of leave-one-out-cross-validation (i.e. $k = n$) is however cubic in the number of samples, and in contrast to, e.g., kernel ridge regression, there is no efficient computation in terms of generalized cross-validation (Wahba, 1990). This is mainly due to the fact that kernel PLS is not linear in the sense that $\hat{\mathbf{y}}$ depends nonlinearly on \mathbf{y} .

Information criteria rely on the fact that the true error of a model can be estimated in terms of its training error and its complexity. In regression problems, the complexity is defined in terms of degrees of freedom. For any regression method that fits the response data \mathbf{y} by a vector $\hat{\mathbf{y}}$, Efron (2004) defines

$$\text{df} = E_{\mathbf{y}} \left[\text{trace} \left(\frac{\partial \hat{\mathbf{y}}}{\partial \mathbf{y}} \right) \right].$$

Note that if $\hat{\mathbf{y}}$ is linear in \mathbf{y} , i.e. there is a $n \times n$ matrix \mathbf{H} (called the hat-matrix) such that $\hat{\mathbf{y}} = \mathbf{H}\mathbf{y}$, we retrieve $\text{df} = \text{trace}(\mathbf{H})$, which coincides with the well-known definition given in Hastie and Tibshirani (1990). We obtain an unbiased estimate of the degrees of freedom via

$$\hat{\text{df}} = \text{trace} \left(\frac{\partial \hat{\mathbf{y}}}{\partial \mathbf{y}} \right).$$

Degrees of freedom are not only relevant for model selection, they also allow us to quantify the intrinsic complexity of a regression method. As kernel PLS depends nonlinearly on \mathbf{y} (via \mathbf{T}), we derive the estimated degrees of freedom by an explicit computation of the first derivative of $\hat{\mathbf{y}}_m$. In contrast to Phatak et al. (2002), who compute the first derivative of $\hat{\beta}_m$ based on relationship (11) (which turns out to be extremely time consuming and numerically unstable), we follow the iterative approach of Serneels et al. (2004) for the derivative of $\hat{\beta}_m$, and make explicit use of the recursive definition of kernel PLS of Algorithm 1. To this end, we have to determine the first derivative of the projection operator $\mathcal{P}_{\mathbf{v}}\mathbf{z} = \mathbf{v}(\mathbf{v}^T\mathbf{v})^{-1}\mathbf{v}^T\mathbf{z}$, with both vectors \mathbf{v} and \mathbf{z} depending on \mathbf{y} .

Theorem 1. *The first derivative of $\mathcal{P}_{\mathbf{v}}\mathbf{z}$ is*

$$\frac{\partial \mathcal{P}_{\mathbf{v}}\mathbf{z}}{\partial \mathbf{y}} = \frac{1}{\mathbf{v}^T\mathbf{v}} (\mathbf{v}\mathbf{z}^T + \mathbf{v}^T\mathbf{z}\mathbf{I}_n) (\mathbf{I}_n - \mathcal{P}_{\mathbf{v}}) \frac{\partial \mathbf{v}}{\partial \mathbf{y}} + \mathcal{P}_{\mathbf{v}} \frac{\partial \mathbf{z}}{\partial \mathbf{y}}.$$

Proof. We split the computation of the derivative into two parts by using the fact that after normalization of

the vector \mathbf{v} via $\mathbf{v} \leftarrow \mathbf{v}/\|\mathbf{v}\|$, the projection operator is simply $\mathbf{v}\mathbf{v}^T\mathbf{z}$. Using elementary calculus, we see that

$$\frac{\partial(\mathbf{v}/\|\mathbf{v}\|)}{\partial\mathbf{y}} = \frac{1}{\|\mathbf{v}\|} (\mathbf{I}_n - \mathcal{P}_{\mathbf{v}}) \frac{\partial\mathbf{v}}{\partial\mathbf{y}}.$$

Furthermore, using the product rule, we have

$$\frac{\partial(\mathbf{v}\mathbf{v}^T\mathbf{z})}{\partial\mathbf{y}} = (\mathbf{v}\mathbf{z}^T + \mathbf{v}^T\mathbf{z}\mathbf{I}_n) \frac{\partial\mathbf{v}}{\partial\mathbf{y}} + \mathbf{v}\mathbf{v}^T \frac{\partial\mathbf{z}}{\partial\mathbf{y}}.$$

Combining these two results via the chain rule, we obtain the desired result. \square

We can now derive the degrees of freedom of kernel PLS by using Theorem 1 and by differentiating the recursive formulas in Algorithm 1. This is displayed in Algorithm 2. Figure 2 illustrates the behavior of

Algorithm 2 Kernel PLS with first derivative of PLS

Input: \mathbf{K} , \mathbf{y} , m

Initialize: $\hat{\mathbf{y}}_0 = \boldsymbol{\alpha}_0 = \mathbf{0}_n$, $\frac{\partial\hat{\mathbf{y}}_0}{\partial\mathbf{y}} = \mathbf{0}_{n \times n}$

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

 if $i = 1$ then

$$\mathbf{r}_i = \boldsymbol{\gamma}_i = \mathbf{y}$$

$$\tilde{\mathbf{t}}_i = \mathbf{t}_i = \mathbf{K}\mathbf{y}$$

$$\frac{\partial\tilde{\mathbf{t}}_i}{\partial\mathbf{y}} = \frac{\partial\mathbf{t}_i}{\partial\mathbf{y}} = \mathbf{K}$$

 else

$$\mathbf{r}_i = \mathbf{r}_{i-1} - \mathcal{P}_{\tilde{\mathbf{t}}_{i-1}}\mathbf{y}$$

$$\tilde{\mathbf{t}}_i = \mathbf{K}\mathbf{r}_i = \tilde{\mathbf{t}}_{i-1} - \mathbf{K}\mathcal{P}_{\tilde{\mathbf{t}}_{i-1}}\mathbf{y}$$

$$\frac{\partial\tilde{\mathbf{t}}_i}{\partial\mathbf{y}} = \frac{\partial\tilde{\mathbf{t}}_{i-1}}{\partial\mathbf{y}} - \mathbf{K} \left(\frac{\partial\mathcal{P}_{\tilde{\mathbf{t}}_{i-1}}\mathbf{y}}{\partial\mathbf{y}} \right)$$

$$\boldsymbol{\gamma}_i = \mathbf{r}_i - \frac{\mathbf{t}_{i-1}^T \tilde{\mathbf{t}}_i}{\mathbf{t}_{i-1}^T \mathbf{t}_{i-1}} \boldsymbol{\gamma}_{i-1}$$

$$\mathbf{t}_i = \mathbf{K}\boldsymbol{\gamma}_i = \tilde{\mathbf{t}}_i - \mathcal{P}_{\tilde{\mathbf{t}}_{i-1}}\tilde{\mathbf{t}}_i$$

$$\frac{\partial\mathbf{t}_i}{\partial\mathbf{y}} = \frac{\partial\tilde{\mathbf{t}}_i}{\partial\mathbf{y}} - \frac{\partial\mathcal{P}_{\tilde{\mathbf{t}}_{i-1}}\tilde{\mathbf{t}}_i}{\partial\mathbf{y}}$$

 end if

 exit if $|\angle(\mathbf{t}_i, \mathbf{t}_j)| > 0.1$ for any $1 \leq j < i$

$$\boldsymbol{\alpha}_i = \boldsymbol{\alpha}_{i-1} + \frac{\mathbf{t}_i^T \mathbf{y}}{\mathbf{t}_i^T \mathbf{t}_i} \boldsymbol{\gamma}_i$$

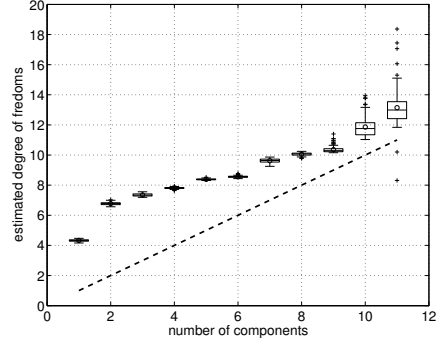
$$\hat{\mathbf{y}}_i = \mathbf{K}\boldsymbol{\alpha}_i = \hat{\mathbf{y}}_{i-1} + \mathcal{P}_{\mathbf{t}_i}\mathbf{y}$$

$$\frac{\partial\hat{\mathbf{y}}_i}{\partial\mathbf{y}} = \frac{\partial\hat{\mathbf{y}}_{i-1}}{\partial\mathbf{y}} + \frac{\partial\mathcal{P}_{\mathbf{t}_i}\mathbf{y}}{\partial\mathbf{y}}$$

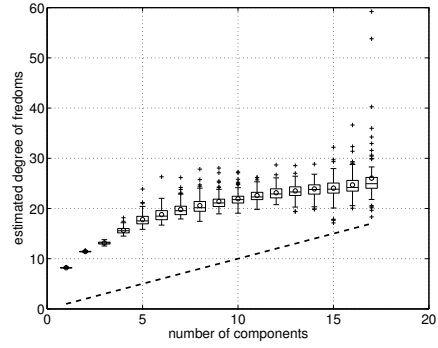
end for

the degrees of freedom of kernel PLS. We use the sinc-data introduced in Section 2.4. We use rbf kernels with bandwidths 1/2 and 1/10 respectively. Keeping the input data \mathbf{X} fixed, we generate the output \mathbf{y} 100 different times (using a standard deviation of ε_i of 1/10) and estimate the degrees of freedom of kernel PLS. The graphs in Figure 2 depict the mean degrees of freedom and their boxplots.

In both cases the degrees of freedom is a concave function, i.e., the degrees of freedom of kernel PLS exceed the number of components. To our experience from



(a) Rbf kernel with width 1/2



(b) Rbf kernel with width 1/10

Figure 2. Estimated degrees of freedom for Kernel PLS.

simulated and real world data, the concavity of the graphs is typical for kernel PLS. This supports the conjecture voiced by Frank and Friedman (1993) that $\text{df} \geq m$. In particular, applying the naïve approach $\text{df}(\text{PLS}) = m$ tends to underestimate the complexity of kernel PLS and is expected to choose model parameters that overfit. In the next section, we compare the performance of kernel PLS with three different information criteria for model selection. We use the Akaike information criterion (AIC), the Bayesian information criterion (BIC) and generalized minimum description length (Hansen & Yu, 2001), which is defined as

$$\text{gMDL} = \frac{n}{2} \log S + \frac{d}{2} \log F + \frac{1}{2} \log n,$$

$$S = \text{RSS}/(n - d), \quad F = (\mathbf{y}^T \mathbf{y} - \text{RSS})/dS,$$

where RSS stands for the residual sums of squares.

4. Experiments

We now present experimental results for the proposed model selection schemes. We will also discuss how to efficiently implement these procedures practically. Note that we do not want to prove the competitiveness of kernel PLS. This has been done before, e.g. by Rosipal and Trejo (2001).

We compare three sets of model selection schemes: (1) the AIC, BIC, and gMDL methods based on our degrees of freedom estimate, (2) leave-one-out cross-validation as a baseline, and (3) 5-fold cross-validation as a computationally faster variant.

In order to implement cross-validation schemes efficiently for kernel PLS, note that the iterative computation of kernel PLS from Algorithm 1 not only computes the final parameter vector for m components, but also computes all α_i for $1 \leq i \leq m$. Therefore, one computation allows to compute the test errors for all choices up to m components. Since the main algorithm requires $O(mn^2)$ computations, it follows that the leave-one-out cross-validation error can be computed in $O(mn^3)$ for a given kernel and all numbers of components up to m . This computational complexity is similar to that of computing the leave-one-out cross-validation error for KRR in closed form (Wahba, 1990). For k -fold cross-validation, the computational cost become significantly better, as the runtime is of order $O(kmn^2)$.

The degrees of freedom estimate also computes all intermediate quantities, such that all degrees of freedom can be computed at once. However, the algorithm involves one matrix-matrix multiplication, such that the overall complexity is practically cubic (although the theoretical run-time is sub-cubic). Note that matrix-matrix multiplications can be parallelized efficiently in contrast to matrix inversions. In summary, all m degrees of freedom estimates can be computed in $O(mn^3)$.

As data sets, we chose the “bank” and “kin” regression data sets from the delve repository¹. Both data sets are synthetic data sets generated from simulations. The bank data set models customers waiting in a queue, and the tasks consists in predicting the percentage of customers who loose their patience before arriving at the head of the queue. The kin data set is based on a model of a robotic arm, and the task consists in predicting the position of the arm based on the angles of its joints. The data sets come in four different flavors, “fh”, “fm”, “nh”, “nm”, where “f” indicates fairly linear connection between inputs and outputs, “n” non-linear connection, “h” “high noise”, and “m” medium levels of noise. Both data sets are eight-dimensional, and consist of 8192 data points.

For our experiments, the data sets are split into 20 realizations of 100 training examples and 8092 test examples. Furthermore, the inputs are normalized to lie in the interval $[-1, 1]$. We use a Gaussian ker-

nel with widths chosen from 20 logarithmically spaced points between 1 and 10^4 , and the maximum number of components are set to 50. Both kernel width and the number of components are chosen using the respective model selection criterion. Tables 1 and 2 show the resulting mean squared errors and standard deviations over the 20 realizations.

As we expected, leave-one-out cross-validation leads to the best performances. While the AIC and BIC perform poorly on the data sets, the gMDL criterion is on par with leave-one-out cross-validation. But also the computationally much more efficient 5-fold cross-validation still leads to comparably good results.

In summary, we see that the gMDL criterion leads to effective model selection giving further evidence for the practical relevance of our degrees of freedom estimator. On the cross-validation side, even the relatively coarse, but fast, 5-fold cross validation performs well. The additional benefit of the gMDL criterion is that it also produces a degrees of freedom estimate which can then be used to compare kernel PLS to other (linear) methods. Both methods lead to performances on par to those achieved by leave-one-out cross-validation.

5. Conclusions

We have discussed several topics related to kernel PLS relevant to the more widespread adoption of kernel PLS in a machine learning setting. It has turned out that kernel PLS has some interesting properties not shared with ordinary PLS, allowing for alternative interpretations of the algorithm beyond co-optimizing covariance. Furthermore, we have addressed practical issues including efficient implementations, estimation of degrees of freedom, and model selection, using either degrees of freedom or cross-validation based criteria. In practical experiments, we have shown how to perform efficient and effective model selection.

References

- Braun, M. L., Buhmann, J. M., & Müller, K.-R. (2007). Denoising and Dimension Reduction in Feature Space. *Advances in Neural Information Processing Systems 19*. Cambridge, MA: MIT Press.
- Efron, B. (2004). The Estimation of Prediction Error: Covariance Penalties and Cross-Validation. *Journal of the American Statistical Association*, 99, 619–633.
- Frank, I., & Friedman, J. (1993). A Statistical View of Some Chemometrics Regression Tools. *Technometrics*, 35, 109–135.

¹<http://www.cs.toronto.edu/~delve>

Table 1. Results on the “bank” data set (best result, on par = within 2× standard deviation of best result)

METHOD	FH	FM	NH	NM
AIC	$9.38 \pm 2.49 \times 10^{-3}$	$4.18 \pm 1.71 \times 10^{-3}$	$6.09 \pm 1.11 \times 10^{-3}$	$1.77 \pm 0.17 \times 10^{-3}$
BIC	$6.25 \pm 0.58 \times 10^{-3}$	$2.18 \pm 0.41 \times 10^{-3}$	$4.32 \pm 0.75 \times 10^{-3}$	$1.23 \pm 0.17 \times 10^{-3}$
GMDL	$6.25 \pm 0.58 \times 10^{-3}$	$2.16 \pm 0.38 \times 10^{-3}$	$4.00 \pm 0.43 \times 10^{-3}$	$1.23 \pm 0.17 \times 10^{-3}$
5-FOLD CV	$6.29 \pm 0.46 \times 10^{-3}$	$2.31 \pm 0.44 \times 10^{-3}$	$4.16 \pm 0.56 \times 10^{-3}$	$1.16 \pm 0.11 \times 10^{-3}$
LOO-CV	$6.30 \pm 0.53 \times 10^{-3}$	$2.28 \pm 0.35 \times 10^{-3}$	$4.12 \pm 0.65 \times 10^{-3}$	$1.18 \pm 0.16 \times 10^{-3}$

Table 2. Results on the “kin” data set (best result, on par = within 2× standard deviation of best result)

METHOD	FH	FM	NH	NM
AIC	$3.17 \pm 0.43 \times 10^{-2}$	$2.99 \pm 0.47 \times 10^{-2}$	$6.92 \pm 0.33 \times 10^{-2}$	$4.91 \pm 0.46 \times 10^{-2}$
BIC	$2.44 \pm 1.37 \times 10^{-2}$	$2.77 \pm 7.15 \times 10^{-3}$	$6.92 \pm 0.33 \times 10^{-2}$	$4.91 \pm 0.46 \times 10^{-2}$
GMDL	$2.31 \pm 0.15 \times 10^{-3}$	$4.84 \pm 0.51 \times 10^{-4}$	$5.73 \pm 1.07 \times 10^{-2}$	$4.84 \pm 0.91 \times 10^{-2}$
5-FOLD CV	$2.29 \pm 0.12 \times 10^{-3}$	$4.38 \pm 0.65 \times 10^{-4}$	$5.14 \pm 0.53 \times 10^{-2}$	$4.23 \pm 0.32 \times 10^{-2}$
LOO-CV	$2.30 \pm 0.17 \times 10^{-3}$	$4.60 \pm 0.36 \times 10^{-4}$	$5.26 \pm 0.63 \times 10^{-2}$	$4.19 \pm 0.30 \times 10^{-2}$

- Golub, G., & Van Loan, C. (1996). *Matrix Computations*. Johns Hopkins University Press.
- Hansen, M., & Yu, B. (2001). Model Selection and Minimum Description Length Principle. *Journal of the American Statistical Association*, 96, 746–774.
- Hastie, T., & Tibshirani, R. (1990). *Generalized Additive Models*. Chapman and Hall, London.
- Helland, I. (1988). On the Structure of Partial Least Squares Regression. *Communications in Statistics, Simulation and Computation*, 17, 581–607.
- 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.
- Lingjærde, O., & Christophersen, N. (2000). Shrinkage Structure of Partial Least Squares. *Scandinavian Journal of Statistics*, 27, 459–473.
- Manne, R. (1987). Analysis of Two Partial-Least-Squares Algorithms for Multivariate Calibration. *Chemometrics and Intelligent Laboratory Systems*, 2, 187–197.
- Phatak, A., & de Hoog, F. (2002). Exploiting the Connection between PLS, Lanczos Methods and Conjugate Gradients: Alternative Proofs of Some Properties of PLS. *Journal of Chemometrics*, 16, 361–367.
- Phatak, A., Riley, P., & Penlidis, A. (2002). The Asymptotic Variance of the Univariate PLS Estimator. *Linear Algebra and its Applications*, 354, 245–253.
- 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, 34–51. Springer.
- 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. *Proceedings of the Twentieth International Conference on Machine Learning* (pp. 640–647). Washington, DC.
- Serneels, S., Lemberge, P., & Espen, P. V. (2004). Calculation of PLS Prediction Intervals Using Efficient Recursive Relations for the Jacobian Matrix. *Journal of Chemometrics*, 18, 76–80.
- Wahba, G. (1990). *Spline Models For Observational Data*. Society for Industrial and Applied Mathematics.
- Wold, H. (1975). Path models with Latent Variables: The NIPALS Approach. In et H. B. al. (Ed.), *Quantitative Sociology: International Perspectives on Mathematical and Statistical Model Building*, 307–357. Academic Press.