

# Lanczos Approximations for the Speedup of Kernel Partial Least Squares Regression\*

Nicole Krämer

Machine Learning Group  
Berlin Institute of Technology  
nkraemer@cs.tu-berlin.de

Masashi Sugiyama

Department of Computer Science  
Tokyo Institute of Technology  
sugi@cs.titech.ac.jp

Mikio L. Braun

Machine Learning Group  
Berlin Institute of Technology  
mikio@cs.tu-berlin.de

February 19, 2009

## Abstract

The runtime for Kernel Partial Least Squares (KPLS) to compute the fit is quadratic in the number of examples. However, the necessity of obtaining sensitivity measures as degrees of freedom for model selection or confidence intervals for more detailed analysis requires cubic runtime, and thus constitutes a computational bottleneck in real-world data analysis. We propose a novel algorithm for KPLS which not only computes (a) the fit, but also (b) its approximate degrees of freedom and (c) error bars in quadratic runtime. The algorithm exploits a close connection between Kernel PLS and the Lanczos algorithm for approximating the eigenvalues of symmetric matrices, and uses this approximation to compute the trace of powers of the kernel matrix in quadratic runtime.

## 1 INTRODUCTION

Partial Least Squares (PLS) [28, 29] is a supervised dimensionality reduction technique. Given  $n$  observations  $(\mathbf{x}_i, y_i) \in \mathbb{R}^d \times \mathbb{R}$ , it iteratively constructs an orthogonal set  $\mathbf{T} = (\mathbf{t}_1, \dots, \mathbf{t}_m) \in \mathbb{R}^{n \times m}$  of  $m$  latent features which have maximal covariance with the target variable  $\mathbf{y} = (y_1, \dots, y_n)$ . For regression, these latent components are used as predictors in a least squares fit instead of the original data leading to fitted values

$$\hat{\mathbf{y}}_m = \mathbf{T} \left( \mathbf{T}^\top \mathbf{T} \right)^{-1} \mathbf{T}^\top \mathbf{y} = \mathcal{P}_{\mathbf{T}} \mathbf{y}, \quad (1)$$

where  $\mathcal{P}$  denotes the orthogonal projection operator. PLS is the standard tool e.g. in chemometrics [16], and has been successfully applied in various other scientific fields such as chemoinformatics, physiology or bioinformatics [25, 23, 1]. In combination with the kernel trick [20, 22], Kernel Partial Least Squares (KPLS) performs dimensionality reduction and

---

\*to appear in Proceedings of the 12th International Conference on Artificial Intelligence and Statistics (AISTATS 09)

regression in a non-linear fashion. KPLS has some appealing properties over existing kernel methods. Due to its iterative nature, it only relies on matrix-vector multiplications. Hence its runtime is quadratic in the number of training examples, as opposed to – for example – Kernel Ridge Regression, which requires the inversion of a large symmetric matrix, having time complexity  $\mathcal{O}(n^3)$ . Furthermore, it is possible to compute the derivative of the KPLS solution with respect to  $\mathbf{y}$  by differentiating the iterative formulation itself. Taking the trace of the derivative of the fitted values, one obtains an estimate of the degrees of freedom for KPLS, which can be used, for example, for effective model selection based on information criteria like AIC, BIC, or gMDL [13]. The first order Taylor approximation can also be used to construct confidence intervals for PLS [5, 19]. However, since we take the derivative of a vector (1), the derivative is a matrix, and the computation of the derivative involves a number of matrix-matrix multiplications which have time complexity  $\mathcal{O}(n^3)$  for all practical considerations.

In this work, we propose an algorithm which computes (a) the fit of KPLS as well as (b) its approximate degrees of freedom and (c) confidence intervals for the KPLS solutions, all in quadratic runtime. These results are based on the fact that PLS is equivalent to the Lanczos method for approximating the eigenvalues of the kernel matrix  $\mathbf{K}$  by the eigenvalues of a tridiagonal  $m \times m$  matrix  $\mathbf{D}$ . The main contribution is to compute these approximate eigenvalues using KPLS *itself*. Then, using a different formulation of the derivative of the fit in KPLS, one can approximate the trace of powers  $\mathbf{K}^j$  of the kernel matrix using the matrix  $\mathbf{D}$ . Since  $\mathbf{D}$  is typically small (as it scales with the number of components), the runtime for computing the eigenvalues is cubic in  $m$ , and therefore, unproblematic. Since the powers of the Kernel matrices  $\mathbf{K}^j$  are the only matrix-matrix multiplications of order  $n$  in the formula for the degrees of freedom, the approximation leads to quadratic runtime. Hence, we use the KPLS fit to approximate its degrees of freedom. In addition, using the alternative formulation of the derivative, one can perform a sensitivity analysis of KPLS resulting in confidence intervals on the estimates, also in quadratic runtime.

This paper is structured as follows. In Section 2, we review the connection between KPLS and Lanczos approximations, and summarize the state-of-the-art for computing the derivative of Kernel PLS. In Sections 3 and 4, we propose our novel formulation of the derivative together with the quadratic runtime algorithms for the degrees of freedom and the confidence intervals. We conclude with some practical examples.

PLS is closely related to Krylov methods. Therefore, we briefly recall the definition of Krylov subspaces. 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  $\mathcal{K}_m(\mathbf{C}, \mathbf{c})$ .

## 2 BACKGROUND: PLS, LANCZOS METHODS, AND SENSITIVITY ANALYSIS

In this paper, we focus on the NIPALS algorithm [28] for PLS. For different forms of PLS, see [21]. The  $n$  centered observations  $(\mathbf{x}_i, y_i)$  are pooled into a  $n \times d$  data matrix  $\mathbf{X}$  and a vector  $\mathbf{y} \in \mathbb{R}^n$ . PLS constructs  $m$  orthogonal latent components  $\mathbf{T} = (\mathbf{t}_1, \dots, \mathbf{t}_m) \in \mathbb{R}^{n \times m}$  in a greedy fashion. The first component  $\mathbf{t}_1 = \mathbf{X}\mathbf{w}_1$  fulfills

$$\mathbf{w}_1 = \arg \max_{\|\mathbf{w}\|=1} \text{cov}(\mathbf{X}\mathbf{w}, \mathbf{y})^2 = \frac{1}{\|\mathbf{X}^\top \mathbf{y}\|} \mathbf{X}^\top \mathbf{y}. \quad (2)$$

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 variables  $\mathbf{X}$

$$\mathbf{X}_i = \mathbf{X} - \mathcal{P}_{\mathbf{t}_1, \dots, \mathbf{t}_{i-1}} \mathbf{X},$$

and then replacing  $\mathbf{X}$  by  $\mathbf{X}_i$  in (2). The matrix  $\mathbf{W} = (\mathbf{w}_1, \dots, \mathbf{w}_m) \in \mathbb{R}^{d \times m}$  can be shown to be orthogonal as well (e.g. [9]). Note furthermore that the latent components are usually scaled to unit norm. Kernel PLS [20, 22] can be derived by noting that  $\mathbf{w}_i = \mathbf{X}^\top \mathbf{r}_i$  with

$$\mathbf{r}_i = (\mathbf{y} - \hat{\mathbf{y}}_{i-1}) / \|\mathbf{K}^{1/2}(\mathbf{y} - \hat{\mathbf{y}}_{i-1})\| \quad (3)$$

denoting the normalized residuals, and by deflating the kernel matrix  $\mathbf{K}$  instead of  $\mathbf{X}$ ,

$$\mathbf{K}_i = (\mathbf{I}_n - \mathcal{P}_{\mathbf{t}_{i-1}}) \mathbf{K}_{i-1} (\mathbf{I}_n - \mathcal{P}_{\mathbf{t}_{i-1}}).$$

In contrast to e.g. Principal Component Analysis, the latent components  $\mathbf{T}$  depend on the response, and hence the fitted values (1) are a nonlinear function of  $\mathbf{y}$ .

Recall that in the nonlinear case, KPLS depends on the kernel parameters (e.g. the width of an rbf-kernel) and the optimal number  $m$  of latent components. Thus, for model selection, one has to select the optimal combination on a grid of possible kernel parameters and components from 1 to a maximal amount  $m^*$  of components.

## 2.1 KRYLOV METHODS AND LANCZOS APPROXIMATION

To predict the output for a new observation, we have to derive the regression coefficients  $\hat{\boldsymbol{\beta}}_m$  (in the linear case) and kernel coefficients  $\hat{\boldsymbol{\alpha}}_m$  (in the nonlinear case), which are defined via  $\hat{\mathbf{y}}_m = \mathbf{X} \hat{\boldsymbol{\beta}}_m = \mathbf{K} \hat{\boldsymbol{\alpha}}_m$ . This can be done by using the fact that PLS is equivalent to the Lanczos bidiagonalization of  $\mathbf{X}$  [14]: The orthogonal matrices  $\mathbf{T}$  and  $\mathbf{W}$  represent a decomposition of  $\mathbf{X}$  into a bidiagonal matrix  $\mathbf{L} \in \mathbb{R}^{m \times m}$  via

$$\mathbf{XW} = \mathbf{TL} \quad (4)$$

with  $l_{ij} = 0$  for  $i > j$  or  $i < j - 1$ . This matrix is defined as  $\mathbf{L} = \mathbf{T}^\top \mathbf{XW}$ . This implies [15, 9, 22]

$$\hat{\boldsymbol{\beta}}_m = \mathbf{W} \mathbf{L}^{-1} \mathbf{T}^\top \mathbf{y} \quad \text{and} \quad \hat{\boldsymbol{\alpha}}_m = \mathbf{R} \mathbf{L}^{-1} \mathbf{T}^\top \mathbf{y}$$

with  $\mathbf{R} = (\mathbf{r}_1, \dots, \mathbf{r}_m) \in \mathbb{R}^{n \times m}$ . Furthermore, it can be shown [18] that PLS is equivalent to the conjugate gradient (CG) algorithm [8]. The latter is a procedure that iteratively computes approximate solutions of the normal equation  $\mathbf{A}\boldsymbol{\beta} = \mathbf{b}$  (with  $\mathbf{A} = \mathbf{X}^\top \mathbf{X}$  and  $\mathbf{b} = \mathbf{X}^\top \mathbf{y}$ ) by minimizing the quadratic function  $1/2 \boldsymbol{\beta}^\top \mathbf{A} \boldsymbol{\beta} - \boldsymbol{\beta}^\top \mathbf{b}$  along directions that are  $\mathbf{A}$ -orthogonal. These search directions span the Krylov space defined by  $\mathbf{A}$  and  $\mathbf{b}$ . The approximate solution of CG obtained after  $m$  steps is equal to the PLS estimate  $\hat{\boldsymbol{\beta}}_m$  with  $m$  components. Moreover, the weight vectors  $\mathbf{W}$  are an orthogonal basis of  $\mathcal{K}_m(\mathbf{A}, \mathbf{b})$ .

Krylov methods are also used to approximate eigenvalues of  $\mathbf{A}$  by “restricting”  $\mathbf{A}$  onto Krylov subspaces: In terms of the orthogonal basis  $\mathbf{W}$  of  $\mathcal{K}_m(\mathbf{A}, \mathbf{b})$ , the map

$$\mathbf{D} = \mathcal{P}_{\mathcal{K}_m(\mathbf{A}, \mathbf{b})} \mathbf{A} \mathcal{P}_{\mathcal{K}_m(\mathbf{A}, \mathbf{b})}$$

is represented by

$$\mathbf{D} = \mathbf{W}^\top \mathbf{A} \mathbf{W}. \quad (5)$$

$\mathbf{D}$  is shown to be tridiagonal, and the  $m$  distinct eigenvalues  $\mu_1 > \mu_2 > \dots > \mu_m$  of  $\mathbf{D}$  – called Ritz values – are good approximations of the eigenvalues of  $\mathbf{A}$  [24]. One immediate consequence of the connection between PLS and Krylov spaces is the fact that the latent components span the Krylov space defined by  $\mathbf{K}$  and  $\mathbf{K}\mathbf{y}$ . This implies that

$$\hat{\mathbf{y}}_m = \mathcal{P}_{\mathcal{K}_m(\mathbf{K}, \mathbf{K}\mathbf{y})} \mathbf{y}. \quad (6)$$

## 2.2 SENSITIVITY ANALYSIS FOR KPLS

Sensitivity measures are crucial in at least two important scenarios. On the one hand, they are needed to select the correct model (in terms of a suitable kernel and the number of components) when using information criteria. On the other hand, to assess the stability of the solution, one needs to measure the influence of small noise in the training points on the learned function. For example, areas with a high sensitivity require further data points to stabilize the solution in an ambiguous area. Furthermore, if for some regions, the prediction does not depend on the training points at all, this indicates that further data points are necessary.

Both of these questions – model selection and stability analysis – can be addressed by computing the derivatives of the KPLS solution with respect to  $\mathbf{y}$ , either of the fitted labels  $\hat{\mathbf{y}}_m$ , or of the learned kernel coefficients  $\hat{\boldsymbol{\alpha}}_m$ . Let us consider the regression model

$$Y_i = f(\mathbf{x}_i) + \varepsilon_i, \varepsilon_i \sim \mathcal{N}(0, \sigma^2). \quad (7)$$

For a general regression method with fitted values  $\hat{\mathbf{y}}$ , the degrees of freedom are defined as [30, 6]

$$\text{DoF} = E[\text{trace}(\partial \hat{\mathbf{y}} / \partial \mathbf{y})]$$

with the expectation  $E$  taken with respect to  $Y_1, \dots, Y_n$ . An unbiased plug-in estimate of the degrees of freedom is therefore given by

$$\widehat{\text{DoF}} = \text{trace}(\partial \hat{\mathbf{y}} / \partial \mathbf{y}). \quad (8)$$

Degrees of freedom in combination with information criteria can be used for model selection. As the KPLS solution depends nonlinearly on  $\mathbf{y}$ , the computation of the derivative is necessary. Krämer & Braun [13] derive an algorithm for the derivative of  $\hat{\mathbf{y}}_m$  by transforming the Lanczos decomposition (4) into a Kernel representation and by exploiting its sparsity. The resulting iterative algorithm for (8) is then used successfully for model selection. This method scales cubically in the number of examples.

For the construction of confidence intervals for a fitted kernel function

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

one needs to study the influence of an infinitesimal perturbation in the values of  $\mathbf{y}$ . If the kernel coefficients depended linearly on  $\mathbf{y}$  via  $\hat{\boldsymbol{\alpha}} = \mathbf{H}\mathbf{y}$ , the distribution of the prediction  $\hat{f}(\mathbf{x})$  at any point  $\mathbf{x}$  would be given by

$$\hat{f}(\mathbf{x}) \sim \mathcal{N}\left(\mathbf{k}(\mathbf{x})^\top E[\hat{\boldsymbol{\alpha}}], \sigma^2 \mathbf{k}(\mathbf{x})^\top \mathbf{H} \mathbf{H}^\top \mathbf{k}(\mathbf{x})\right) \quad (9)$$

with  $\mathbf{k}(\mathbf{x}) = (k(\mathbf{x}, \mathbf{x}_1), \dots, k(\mathbf{x}, \mathbf{x}_n)) \in \mathbb{R}^n$ . However, as KPLS depends nonlinearly on  $\mathbf{y}$ , the distribution of  $\hat{\alpha}_m$  can only be determined approximately by using a first order Taylor expansion, i.e. one uses

$$\mathbf{H}_m \approx (\partial \hat{\alpha}_m / \partial \mathbf{y}). \quad (10)$$

To the best of our knowledge, confidence intervals for PLS have only been constructed in the linear setting, but the results can easily be extended to the Kernel case. Phatak et. al. [19] use (6) to explicitly calculate the derivative of the PLS coefficients  $\hat{\beta}_m$ , and obtain an approximate distribution of  $\hat{\beta}_m$ . As the formula depends on matrix multiplications of order  $(nm) \times (nm)$ , this approach is computationally expensive. Furthermore, as the Krylov sequence  $\mathbf{K}\mathbf{y}, \dots, \mathbf{K}^m\mathbf{y}$  is nearly collinear, the formula is numerically unstable. In [5, 26], an iterative formulation of PLS is used to construct the derivative of  $\hat{\beta}_m$ . Finally, we remark that the approach by [13] using the Lanczos decomposition can be extended to the derivative of the kernel coefficients.

The drawback of all of these approaches is their poor scalability. All algorithms are cubic in the number of observations. In the following two sections we exploit that we do not need the derivative itself, but only the trace of the derivative for the degrees of freedom, and (9) for the construction of confidence intervals. The key advantage is that we can compute these approximation schemes in quadratic runtime.

### 3 APPROXIMATE DEGREES OF FREEDOM IN QUADRATIC RUNTIME

The key ingredients for the derivation of approximate degrees of freedom are (1) the identification of those terms that are cubic in  $n$ , and (2) the approximation of those terms using Lanczos approximations.

First, we extend the results of [19] to the computation of the derivative of  $\hat{\mathbf{y}}_m$ . We define the  $m \times m$  matrix  $\mathbf{B}$  via  $b_{ij} = \langle \mathbf{t}_i, \mathbf{K}^j \mathbf{y} \rangle$ . The matrix is regular and upper triangular, as the latent components  $\mathbf{T}$  are an orthogonal basis of the Krylov subspace  $\mathcal{K}_m(\mathbf{K}, \mathbf{K}\mathbf{y})$ .

**Proposition 1.** *Let  $\mathbf{c} = \mathbf{B}^{-1} \mathbf{T}^\top \mathbf{y}$  and  $\mathbf{V} = (\mathbf{v}_1, \dots, \mathbf{v}_m) = \mathbf{T} \mathbf{B}^{-\top}$ . We have*

$$\begin{aligned} \frac{\partial \hat{\mathbf{y}}_m}{\partial \mathbf{y}} &= \left[ \mathbf{c}^\top \otimes (\mathbf{I}_n - \mathcal{P}_{\mathbf{T}}) \right] \mathbf{Q}^\top + \left[ \mathbf{V} \otimes (\mathbf{y} - \hat{\mathbf{y}}_m)^\top \right] \mathbf{Q}^\top + \mathcal{P}_{\mathbf{T}} \\ &= \sum_{j=1}^m c_j \left( \mathbf{I}_n - \mathbf{T} \mathbf{T}^\top \right) \mathbf{K}^j + \sum_{j=1}^m \mathbf{v}_j (\mathbf{y} - \hat{\mathbf{y}}_m)^\top \mathbf{K}^j + \mathbf{T} \mathbf{T}^\top. \end{aligned}$$

Here,  $\otimes$  is the Kronecker product and  $\mathbf{Q} = (\mathbf{K}, \mathbf{K}^2, \dots, \mathbf{K}^m) \in \mathbb{R}^{n \times nm}$ .

*Proof.* The first line follows by computing the derivative of the projection operator (6) and by applying a basis transformation from the Krylov sequence  $\mathbf{K}\mathbf{y}, \dots, \mathbf{K}^m\mathbf{y}$  to the orthogonal basis  $\mathbf{t}_1, \dots, \mathbf{t}_m$ . The latter ensures that the formula is numerically more stable. For the second line, we represent the Kronecker product as a sum.  $\square$

As a consequence, we yield a formula for the degrees of freedom of KPLS.

**Corollary 2.** *An unbiased estimated of the degrees of freedom of KPLS with  $m$  components is given by*

$$\begin{aligned}\widehat{\text{DoF}}(m) &= \sum_{j=1}^m c_j \text{trace} \left[ \left( \mathbf{I}_n - \mathbf{T}\mathbf{T}^\top \right) \mathbf{K}^j \right] + \sum_{j=1}^m (\mathbf{y} - \widehat{\mathbf{y}}_m)^\top \mathbf{K}^j \mathbf{v}_j + m \\ &= \sum_{j=1}^m c_j \text{trace} (\mathbf{K}^j) + m \\ &\quad - \sum_{j=1}^m \left( \sum_{l=1}^m \mathbf{t}_l^\top \mathbf{K}^j \mathbf{t}_l \right) + (\mathbf{y} - \widehat{\mathbf{y}}_m)^\top \sum_{j=1}^m \mathbf{K}^j \mathbf{v}_j\end{aligned}$$

This representation of the DoF of KPLS reveals an interesting feature. The computation of last line is quadratic in  $n$ , as it only involves matrix-vector multiplications. The first line however is cubic in  $n$ , as we need to compute the trace of powers of the kernel matrix  $\mathbf{K}^j$  for  $j = 1, \dots, m$ .

### 3.1 APPROXIMATE DEGREES OF FREEDOM VIA RITZ VALUES

As explained above, PLS is equivalent to Lanczos approximations, and can be used to approximate the eigenvalues of  $\mathbf{X}^\top \mathbf{X}$  via the tridiagonal matrix  $\mathbf{D}$  defined in (5). Note that  $\mathbf{D}$  has a kernel representation

$$\mathbf{D} = \mathbf{R}^\top \mathbf{K}^2 \mathbf{R} = \mathbf{L}^\top \mathbf{L}. \quad (11)$$

with  $\mathbf{R}$  the matrix of normalized residuals defined in (3). The eigenvalues of  $\mathbf{D}$  are called Ritz values and constitute approximations of the eigenvalues of  $\mathbf{X}^\top \mathbf{X}$  [24]. The quality of the approximation increases with the number  $m$  of latent components, and as the computation of the Ritz values scales cubically only in  $m$ , an efficient strategy is to allow a generous amount of components for the computation of  $\mathbf{D}$ .

As the eigenvalues of  $\mathbf{X}^\top \mathbf{X}$  correspond to the eigenvalues of  $\mathbf{K}$  in the kernel setting, we can use Ritz values to derive an approximation of the trace of  $\mathbf{K}^j$ .

**Definition 3** (Approximate Degrees of Freedom). *We define the approximate degrees of freedom of KPLS with  $m$  components as*

$$\begin{aligned}\widehat{\text{DoF}}_{\text{appr}}(m) &= \sum_{j=1}^m c_j \text{trace} (\mathbf{D}_{m_{\max}}^j) + m \\ &\quad - \sum_{j=1}^m \left( \sum_{l=1}^m \mathbf{t}_l^\top \mathbf{K}^j \mathbf{t}_l \right) + (\mathbf{y} - \widehat{\mathbf{y}}_m)^\top \sum_{j=1}^m \mathbf{K}^j \mathbf{v}_j,\end{aligned}$$

where  $\mathbf{D}_{m_{\max}}$  is the tridiagonal matrix defined in (11) computed with  $m_{\max} \geq m$  latent components.

The computation of  $\mathbf{D}$  only requires one additional  $m_{\max} \times m_{\max}$  matrix multiplication  $\mathbf{L}^\top \mathbf{L}$ . As the matrix  $\mathbf{D}$  is of size  $m_{\max} \times m_{\max}$ , the runtime for the computation is cubic in the number of maximal components  $m_{\max}$  (which is typically small), and quadratic in the number  $n$  of examples.

### 3.2 QUALITY OF THE APPROXIMATION

Theoretically, the validity of this approximation can be justified in terms of a deviation bound.

**Proposition 4** (Saad [24]). *Denote by  $\mu_1, \dots, \mu_m$  the eigenvalues of  $\mathbf{D}$  and by  $\lambda_1 \geq \dots \geq \lambda_n$  the eigenvalues of the Kernel matrix  $\mathbf{K}$ . We have*

$$0 \leq \lambda_i - \mu_i \leq (\lambda_1 - \lambda_n) \left( \frac{\kappa_i \tan \theta_i}{C_{m-i}(1 + 2\gamma_i)} \right)^2$$

with  $\mathbf{u}_i$  the  $i$ th eigenvector of  $\mathbf{K}$ ,

$$\theta_i = \arccos \frac{\langle \mathbf{y}, \sqrt{\lambda_i} \mathbf{u}_i \rangle}{\|\mathbf{y}\|_K}$$

and

$$\kappa_i = \prod_{j=1}^{i-1} \frac{\mu_j - \lambda_n}{\mu_j - \lambda_i} \quad \gamma_i = \frac{\lambda_i - \lambda_{i-1}}{\lambda_{i+1} - \lambda_n}.$$

Here,  $C_l$  denotes the Chebychev polynomial of order  $l$ .

Note that  $\theta_i$  is the angle between  $\mathbf{b}$  and the  $i$ th eigenvector of  $\mathbf{X}^\top \mathbf{X}$  - computed in feature space. This inequality implies that the approximation for the  $i$ th eigenvalue is good under two different scenarios. Either  $\lambda_i$  is already close to zero, so  $\mu_i \leq \lambda_i$  is close to zero as well. For large eigenvalues  $\lambda_i$ , the approximation is good if (a) the eigenvalues of  $\mathbf{K}$  decay fast, (b) the angle  $\theta_i$  corresponding to the  $i$ th eigenvector is small, and (c) the index  $i$  is not too large compared to  $m$ . Property (a) is a feature of rbf-kernels, which we use throughout the rest of the paper. Condition (b) is typically fulfilled for the leading eigenvectors of  $\mathbf{K}$  [2, 3], and condition (c) can be fulfilled by using a sufficient large amount  $m_{\max}$  of components.

In practical applications, two important issues are the quality of the approximate degrees of freedom, and the quality of the model selection criteria based on these approximate degrees of freedom. In accordance with [13], we choose generalized minimum description length (gMDL) [7] as model selection criterion.

The simulation setting follows the regression model (7) with  $f(x) = \text{sinc}(x)$ . We draw  $n = 100$  inputs  $X_i$  uniformly from  $[-\pi, \pi]$  and set the standard deviation to  $\sigma = 0.1$ . We fit KPLS with three different rbf-kernels of width 0.01, 0.1, 1 and use different numbers  $m_{\max}$  of maximal components. In addition, we compute the DoF, the approximate DoF, the gMDL criterion, and gMDL based on the approximate DoF.

Figure 1 displays the results for the different kernel widths 0.01 (left), 0.1 (center) and 1 (right). The first row shows the degrees of freedom of KPLS (blue line) and approximate degrees of freedom of KPLS depending on the number of maximal components (red dashed line). As indicated by Proposition 4, the approximation becomes more accurate if  $m_{\max}$  is large. Furthermore, the approximation depends on the width of the rbf-kernel. For very small kernel widths (left), the eigenvalues of the kernel matrix decay very slowly, and more components are needed to compensate. In the second line of Figure 1, we display gMDL (blue line) and the approximate gMDL depending on the number of maximal components (red dashed line). The behavior of the approximation is qualitatively the same: It depends on the size of the kernel widths, and in general, it becomes more accurate if more components are used to compute  $\mathbf{D}$ .

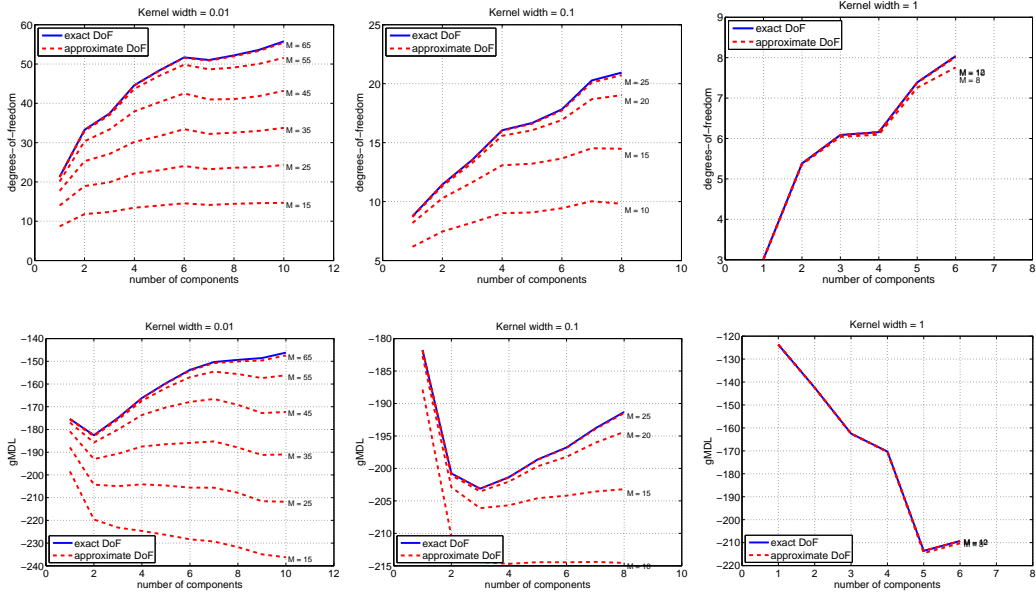


Figure 1: Quality of the approximate degrees of freedom. Results for kernel widths 0.01 (left), 0.1 (center), and 1 (right). Top row: DoF (blue line) and approximate DoF (red dashed line) for different numbers of maximal components. Bottom row: gMDL (blue line) and approximate gMDL (red dashed line) for different numbers of maximal components.

### 3.3 RUNTIME COMPARISON

As shown above, the approximation of the Degrees of Freedom of KPLS leads to reduction in runtime from cubic to quadratic. We now illustrate that this leads to a considerable speed up even for medium sized data. We used the kin regression data set from the delve repository<sup>1</sup>. This eight-dimensional synthetic 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. It consists of 8192 data points. For sub-samples of size 100, 200,  $\dots$ , 1000, we compute (a) KPLS and its Degrees of Freedom for up to  $m = 10$  components and (b) KPLS and its approximate Degrees of Freedom for up to  $m = 10$  components. In both cases, we use a Gaussian Kernel. Note that for (b), we compute  $m_{max} = 30$  components in order to obtain a close approximation, hence the number of KPLS iterations is three times higher for alternative (b). The runtime of both variants are displayed in Fig. 2. The gap between the two methods is clearly visible already for small sample sizes, and the two graphs show the expected quadratic versus cubic form. While the latter is an empirical illustration of the theoretical runtime analysis that we present above, it is important to stress that the improvement from  $\mathcal{O}(n^3)$  to  $\mathcal{O}(n^2)$  is not an asymptotic result but also leads to a significant improvement in runtime already for medium sized data.

<sup>1</sup><http://www.cs.toronto.edu/~delve/>



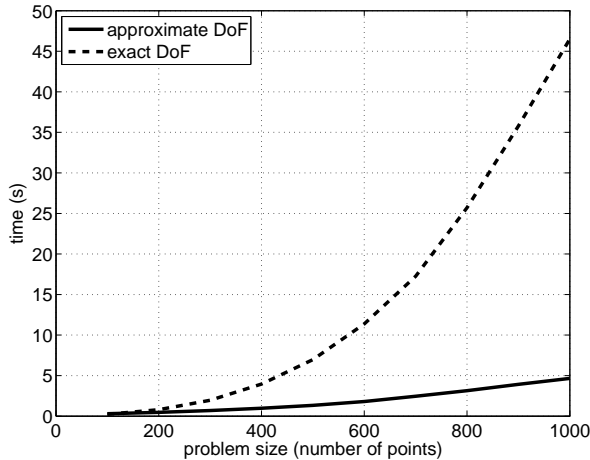


Figure 2: Comparison of runtime on the "kin" data set. Jagged line: KPLS with exact Degrees of Freedom for  $m = 10$  components. Solid line: KPLS with approximate Degrees of Freedom for  $m = 10$  components and  $m_{max} = 30$  components for the approximation of the eigenvalues of the kernel matrix. Hence, for the approximation, the effective number of components is three times higher.

## 4 CONFIDENCE INTERVALS IN QUADRATIC RUNTIME

For the derivation of (approximate) confidence intervals (9), we need to compute the quantity  $\mathbf{H}_m \mathbf{k}(\mathbf{x})$ , where  $\mathbf{H}_m \mathbf{y}$  is the first order Taylor approximation of the kernel coefficients  $\hat{\alpha}_m$ . Using the representation from proposition 1, we can directly compute this matrix-vector product, even without approximating the eigenvalues and thus compute the exact expression in quadratic runtime.

Note that Taylor expansions occur in both types of approximations, for the Degrees of Freedom as well as for the confidence intervals. However, there are essential conceptual differences. For the Degrees of Freedom, the representation in terms of derivatives (8) is in fact no approximation but due to the assumption of normally distributed errors in (7) which leads to Stein's Lemma [27]. In this case, the Degrees of Freedom are approximated using Lanczos methods and Ritz values. In contrast, for the confidence intervals, we have to use the Taylor expansion (10) to obtain an approximate distribution (9) for the KPLS parameters. The computation of the Taylor expansion  $\mathbf{H}_m$  defined in (10) is cubic in  $n$  as it involves multiplications of matrices of size  $n \times n$ . Here, we reduce the computational cost to  $\mathcal{O}(n^2)$  by cleverly exploiting the fact that the matrix-vector product  $\mathbf{H}_m \mathbf{k}(\mathbf{x})$  is a sufficient statistic.

**Proposition 5.** *We have*

$$\begin{aligned} \mathbf{H}_m^\top \mathbf{k}(\mathbf{x}) &= \sum_{j=1}^m \mathbf{K}^{j-1} \left\{ c_j \left( \mathbf{I}_n - \mathbf{K} \mathbf{T} \mathbf{N} \mathbf{R}^\top \right) \right. \\ &\quad \left. + \mathbf{K} (\mathbf{y} - \hat{\mathbf{y}}_m) \mathbf{u}_j^\top \right\} \mathbf{k}(\mathbf{x}) + \mathbf{T} \mathbf{N} \mathbf{R}^\top \mathbf{k}(\mathbf{x}). \end{aligned}$$

with  $\mathbf{R}$  denoting the matrix of normalized residuals,  $\mathbf{N}$  denoting the  $m \times m$  diagonal matrix consisting of elements  $n_{ii} = 1/\|\mathbf{K} \mathbf{r}_i\|$  and

$$\mathbf{U} = (\mathbf{u}_1, \dots, \mathbf{u}_m) = \mathbf{R} \mathbf{N} \mathbf{B}^{-\top}.$$

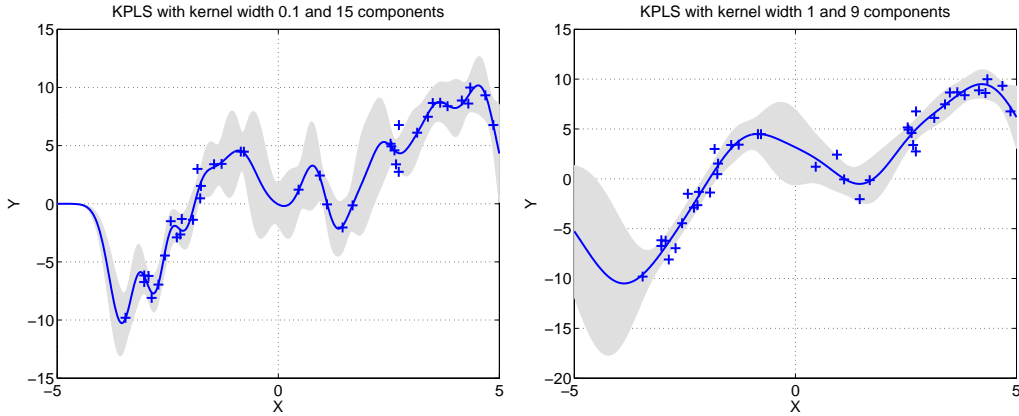


Figure 3: Confidence Intervals for KPLS. Left: KPLS with 15 components and an rbf-kernel of width 0.1 Right: KPLS with 9 components and an rbf kernel of width 1

*Proof.* As

$$\partial \hat{\mathbf{y}}_m / \partial \mathbf{y} = \mathbf{K} (\partial \hat{\boldsymbol{\alpha}}_m / \partial \mathbf{y}) ,$$

the formula can be shown by “canceling out”  $\mathbf{K}$  in the formula of the derivative of  $\hat{\mathbf{y}}_m$ , and then multiplying the formula with  $\mathbf{k}(\mathbf{x})$ .  $\square$

**Illustration** Again, we use the regression model (7), with  $f(x) = (x - 1)(x + 2)(x - 1.5)\exp(-x^2/10)$  and  $\sigma = 1$ . We draw  $n = 40$  points  $X_i$  from a mixture of two normal distribution with mean  $-2$  and  $3$  a variance of  $1$  in both cases. We fit KPLS with for two different models, (1) KPLS with 15 components and an rbf-kernel of width  $0.1$  and (2) KPLS with 9 components and an rbf kernel of width  $1$ . Figure 3 shows the KPLS fit and its confidence intervals (based on a level of  $98\%$ ) for the two models.

In areas with high data density, the prediction is quite stable with small confidence intervals. Next to such high density areas, the predictions becomes unstable, as they can depend quite sensitively on the neighboring data. Finally, when one moves far away from the data points, their influence decreases to zero. This is much more apparent in the left plot with the small kernel widths.

## 5 CONCLUSION

We proposed an implementation of the Kernel PLS method which not only computes the fit in quadratic time, but a degree-of-freedom estimate and confidence intervals based on a sensitivity analysis, which formerly required cubic runtime. The latter estimates can be used, for example, for model selection, or to measure the local stability of the learned function. The approximation schemes exploit the fact that Kernel PLS can be extended to compute Lanczos type approximations of the eigenvalues as well. Together with a novel formula for computing the derivatives of the kernel parameters  $\boldsymbol{\alpha}$ , these approximations allow us to replace costly computation of powers of the kernel matrix. In summary, one obtains a Kernel PLS algorithm which also provides relevant additional information for model selection and provide further insight into the complexity and stability of the learned function.

Our results capitalize on the close connection between the dimensionality reduction technique PLS on the one hand and Krylov methods and Lanczos approximations on the other hand. While the latter two methods are commonly used in numerical linear algebra, their benefits for data analysis have not yet been exploited sufficiently. Only recently (e.g. [17, 4, 10]) they are utilized explicitly in a machine learning framework. Recent research results on the correspondence of penalization techniques and preconditioning of linear systems [11, 12] further underpin the strong potential of these methods. We strongly believe that the interplay between numerical linear algebra and machine learning will further stimulate the field of data analysis.

**Acknowledgement** This work is funded in part by the FP7-ICT Programme of the European Community, under the PASCAL2 Network of Excellence, ICT-216886, by the BMBF grant FKZ 01-IS07007A (ReMind), by the MEXT Grant-in-Aid for Young Scientists (A), 20680007, and by the JFE 21st Century Foundation.

## References

- [1] A.-L. Boulesteix and K. Strimmer. Partial Least Squares: A Versatile Tool for the Analysis of High-Dimensional Genomic Data. *Briefings in Bioinformatics*, 8(1):32–44, 2007.
- [2] M. L. Braun, J. M. Buhmann, and K.-R. Müller. Denoising and Dimension Reduction in Feature Space. *Advances in Neural Information Processing Systems*, 19:185–192, 2007.
- [3] M.L. Braun, J.M. Buhmann, and K.R. Müller. On Relevant Dimensions in Kernel Feature Spaces. *Journal of Machine Learning Research*, 9:1875–1908, 2008.
- [4] N. de Freitas, Y. Wang, M. Mahdaviani, and D. Lang. Fast Krylov Methods for N-body Learning. *Advances in Neural Information Processing Systems*, 18:251–258, 2006.
- [5] M. C. Denham. Prediction Intervals in Partial Least Squares. *Journal of Chemometrics*, 11(1):39–52, 1997.
- [6] B. Efron. The Estimation of Prediction Error: Covariance Penalties and Cross-Validation. *Journal of the American Statistical Association*, 99(467):619–633, 2004.
- [7] M. Hansen and B. Yu. Model Selection and Minimum Description Length Principle. *Journal of the American Statistical Association*, 96:746–774, 2001.
- [8] M. Hestenes and E. Stiefel. Methods for Conjugate Gradients for Solving Linear Systems. *Journal of Research of the National Bureau of Standards*, 49:409–436, 1952.
- [9] A. Hoskuldsson. PLS Regression Methods. *Journal of Chemometrics*, 2(3):211–228, 1988.
- [10] T. Ide and K. Tsuda. Change Point Detection Using Krylov Subspace Learning. In *Proceedings of the SIAM International Conference on Data Mining*, pages 515 – 520, 2007.

- [11] A. Kondylis and J. Whittaker. Spectral Preconditioning of Krylov Spaces: Combining PLS and PC Regression. *Computational Statistics and Data Analysis*, 50(5):2588–2603, 2008.
- [12] N. Krämer, A.-L. Boulesteix, and G. Tutz. Penalized Partial Least Squares with Applications to B-Spline Transformations and Functional Data. *Chemometrics and Intelligent Laboratory Systems*, 94:60–69, 2008.
- [13] N. Krämer and M.L. Braun. Kernelizing PLS, Degrees of Freedom, and Efficient Model Selection. In *Proceedings of the 24th International Conference on Machine Learning*, pages 441 – 448, 2007.
- [14] C. Lanczos. 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, 1950.
- [15] R. Manne. Analysis of Two Partial-Least-Squares Algorithms for Multivariate Calibration. *Chemometrics and Intelligent Laboratory Systems*, 2:187–197, 1987.
- [16] H. Martens and T. Naes. *Multivariate Calibration*. Wiley, New York, 1989.
- [17] C.S. Ong, X. Mary, S. Canu, and A.J. Smola. Learning with non-positive kernels. In *Proceedings of the twenty-first international conference on Machine learning*, pages 639–646, 2004.
- [18] A. Phatak and F. de Hoog. Exploiting the Connection between PLS, Lanczos, and Conjugate Gradients: Alternative Proofs of Some Properties of PLS. *Journal of Chemometrics*, 16:361–367, 2002.
- [19] A. Phatak, P.M. Riley, and A. Penlidis. The Asymptotic Variance of the Univariate PLS Estimator. *Linear Algebra and its Applications*, 354:245–253, 2002.
- [20] S. Rännar, F. Lindgren, P. Geladi, and S. Wold. A PLS Kernel Algorithm for Data Sets with many Variables and Fewer Objects, Part I: Theory and Applications. *Journal of Chemometrics*, 8:111–125, 1994.
- [21] R. Rosipal and N. Krämer. Overview and Recent Advances in Partial Least Squares. In *Subspace, Latent Structure and Feature Selection Techniques*, Lecture Notes in Computer Science, pages 34–51. 2006.
- [22] R. Rosipal and L.J. Trejo. Kernel Partial Least Squares Regression in Reproducing Kernel Hilbert Spaces. *Journal of Machine Learning Research*, 2:97–123, 2001.
- [23] R. Rosipal, L.J. Trejo, and B. Matthews. Kernel PLS-SVC for Linear and Nonlinear Classification. In *Proceedings of the Twentieth International Conference on Machine Learning*, pages 640–647, Washington, DC, 2003.
- [24] Y. Saad. *Iterative methods for sparse linear systems*. PWS, 1st edition, 1996.
- [25] H. Saigo, N. Krämer, and K. Tsuda. Partial Least Squares Regression for Graph Mining. In *14 International Conference on Knowledge Discovery and Data Mining*, pages 578–586, 2008.

- [26] S. Serneels, P. Lemberge, and P.J. Van Espen. Calculation of PLS Prediction Intervals Using Efficient Recursive Relations for the Jacobian Matrix. *Journal of Chemometrics*, 18:76–80, 2004.
- [27] C. Stein. Estimation of the mean of a multivariate normal distribution. *Annals of Statistics*, 9(6):1135–1151, 1981.
- [28] H. Wold. Path models with Latent Variables: The NIPALS Approach. In H.M. Blalock et al., editor, *Quantitative Sociology: International Perspectives on Mathematical and Statistical Model Building*, pages 307–357. Academic Press, 1975.
- [29] S. Wold, H. Ruhe, H. Wold, and W.J. Dunn III. 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, 1984.
- [30] J. Ye. On Measuring and Correcting the Effects of Data Mining and Model Selection. *Journal of the American Statistical Association*, 93(441):120–131, 1998.