

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

Nicole Krämer<sup>(a)</sup>, Masashi Sugiyama<sup>(b)</sup>, Mikio L. Braun<sup>(a)</sup>  
 nkraemer@cs.tu-berlin.de, sugi@cs.titech.ac.jp, mikio@cs.tu-berlin.de

<sup>(a)</sup>Berlin Institute of Technology, Machine Learning Group

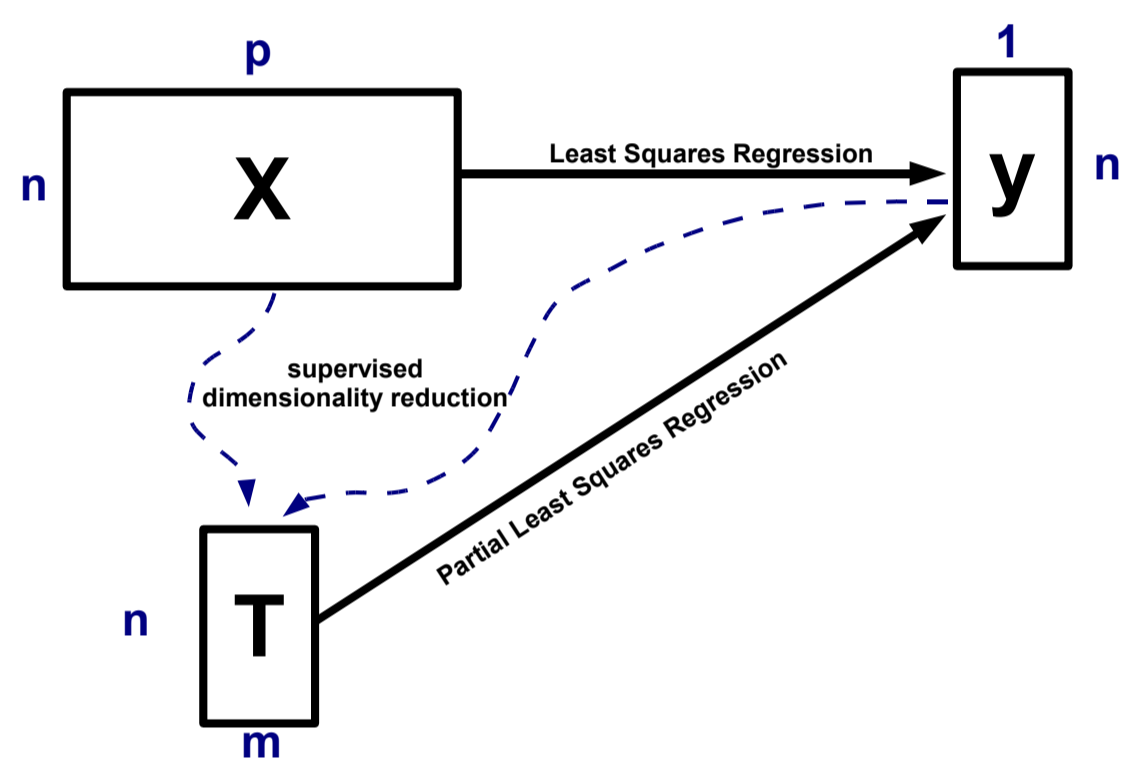
<sup>(b)</sup>Tokyo Institute of Technology, Department of Computer Science

## Summary

- The Degrees of Freedom of Kernel Partial Least Squares (KPLS) require all eigenvalues of the kernel matrix  $\mathbf{K}$ , hence the computation is cubic in the number of observations  $n$ .
- We use Kernel PLS *itself* to approximate the eigenvalues of the kernel matrix.
  - We can compute approximate Degrees of Freedom of KPS in  $\mathcal{O}(n^2)$ !
- We can also compute approximate confidence intervals for KPLS in  $\mathcal{O}(n^2)$ !

## Main Results

### Kernel Partial Least Squares (KPLS)



- Find  $m$  orthogonal latent components  $\mathbf{t}_i = \mathbf{X}\tilde{\mathbf{w}}_i$  with maximal covariance to the response  $\mathbf{y}$ .

$$\tilde{\mathbf{w}}_i = \arg \max_{\|\tilde{\mathbf{w}}\|=1} \text{cov}(\mathbf{X}\tilde{\mathbf{w}}, \mathbf{y})$$

$$\text{s.t. } \mathbf{X}\tilde{\mathbf{w}} \perp \mathbf{X}\tilde{\mathbf{w}}_l = 0 \text{ for } l < i.$$

→ The latent components  $\mathbf{T} = (\mathbf{t}_1, \dots, \mathbf{t}_m)$  depend on the output  $\mathbf{y}$  as well.

- $\mathbf{T}$  is used instead of  $\mathbf{X}$  in a least-squares fit
 
$$\hat{\mathbf{y}}_m = \mathbf{T}(\mathbf{T}^\top \mathbf{T})^{-1} \mathbf{T}^\top \mathbf{y} = \mathcal{P}_{\mathbf{T}} \mathbf{y}$$
 ( $\mathcal{P}$  = projection operator)

### Degrees of Freedom of KPLS

Unbiased estimate for the Degrees of Freedom of KPLS [2]

$$\widehat{\text{DoF}}(m) = \text{trace} \left( \frac{\partial \hat{\mathbf{y}}_m}{\partial \mathbf{y}} \right).$$

### Bad News

$$\widehat{\text{DoF}}(m) = \sum_{j=1}^m c_j \text{trace}(\mathbf{K}^j) + \mathcal{O}(n^2)$$

We need the **eigenvalues of the kernel matrix  $\mathbf{K}$**  for the computation of the Degrees of Freedom of KPLS. The computation of the degrees of freedom of KPLS is **cubic** in the number of observations.

### Lanczos Approximations

Partial Least Squares is equivalent to the Lanczos decomposition of  $\mathbf{X}$

$$\mathbf{T}\mathbf{L} = \mathbf{X}\mathbf{W}$$

with  $\mathbf{T}$  and  $\mathbf{W}$  orthogonal and  $\mathbf{L}$  upper bidiagonal. The eigenvalues of the  $m \times m$  tridiagonal matrix  $\mathbf{D} = \mathbf{L}^\top \mathbf{L}$  are good approximations of the eigenvalues of  $\mathbf{K}$  [4].

### Good News

The eigenvalues of the kernel matrix can be approximated by KPLS itself. Replace  $\mathbf{K}^j$  by  $\mathbf{D}^j$  in the formula for the degrees of freedom. The approximate Degrees of Freedom can be computed in **quadratic** runtime.

### Approximate Confidence Intervals for KPLS

First order Taylor approximation of the kernel coefficients

$$\mathbf{H}_m = (\partial \hat{\boldsymbol{\alpha}}_m / \partial \mathbf{y}).$$

leads to an approximate distribution of the predictions  $\hat{\mathbf{y}}(\mathbf{x})$ .

### More Good News

The product of  $\mathbf{H}_m$  with a vector is a sufficient statistic for the confidence intervals of KPLS. It can be computed in  $\mathcal{O}(n^2)$ .

## Details

### PLS Implementation

#### Partial Least Squares

Input:  $\mathbf{X}_1 = \mathbf{X}, \mathbf{y}, m$   
**for**  $i=1, \dots, m$  **do**  
 $\mathbf{w}_i = \mathbf{X}_i^\top \mathbf{y}$   
 $\mathbf{t}_i = \mathbf{X}_i \mathbf{w}_i$  (component)  
 $\mathbf{X}_{i+1} = \mathbf{X}_i - \mathcal{P}_{\mathbf{t}_i} \mathbf{X}_i$  (deflation)  
**end for**  
 $\mathbf{L} = \mathbf{T}^\top \mathbf{X}\mathbf{W}$  (upper diagonal  $m \times m$  matrix)  
 $\hat{\boldsymbol{\beta}}_m = \mathbf{W}\mathbf{L}^{-1} \mathbf{T}^\top \mathbf{y}$  (regression vector)

#### Kernel Partial Least Squares

Input:  $\mathbf{K}, \mathbf{y}, m, \hat{\mathbf{y}}_0 = \mathbf{t}_0 = \mathbf{0}$   
**for**  $i=1, \dots, m$  **do**  
 $\tilde{\mathbf{r}}_i = \mathbf{y} - \hat{\mathbf{y}}_{i-1}$  (residuals)  
 $\mathbf{t}_i = (\mathbf{I}_n - \mathcal{P}_{\hat{\mathbf{y}}_{i-1}}) \mathbf{K} \tilde{\mathbf{r}}_i$  (component)  
 $\hat{\mathbf{y}}_i = \hat{\mathbf{y}}_{i-1} + \mathcal{P}_{\mathbf{t}_i} \mathbf{y}$  (fitted values)  
**end for**  
 $\mathbf{L} = \mathbf{T}^\top \mathbf{K} \tilde{\mathbf{R}}$  (upper diagonal  $m \times m$  matrix)  
 $\hat{\boldsymbol{\alpha}}_m = \tilde{\mathbf{R}} \mathbf{L}^{-1} \mathbf{T}^\top \mathbf{y}$  (kernel coefficients)

The columns of  $\mathbf{T}$  span the same space as the *Krylov* sequence  $\mathbf{K}\mathbf{y}, \mathbf{K}^2\mathbf{y}, \dots, \mathbf{K}^m\mathbf{y}$ . Hence

$$\hat{\mathbf{y}}_m = \mathcal{P}(\mathbf{K}\mathbf{y}, \mathbf{K}^2\mathbf{y}, \dots, \mathbf{K}^m\mathbf{y})\mathbf{y}. \quad (1)$$

The first derivative of KPLS is computed either via an iterative algorithm [2] or via formula (1) [3]. The computation time is cubic in the number of observations (as it involves matrix-matrix multiplications). Unfortunately, this is also true for its trace:

$$\widehat{\text{DoF}}(m) = \sum_{j=1}^m c_j \text{trace}(\mathbf{K}^j) + m - \sum_{j=1}^m \left( \sum_{l=1}^m \mathbf{t}_l^\top \mathbf{K}^j \mathbf{t}_l \right) + (\mathbf{y} - \hat{\mathbf{y}}_m)^\top \sum_{j=1}^m \mathbf{K}^j \mathbf{v}_j \quad (2)$$

with  $b_{ij} = \langle \mathbf{t}_i, \mathbf{K}^j \mathbf{t}_j \rangle$ ,  $\mathbf{c} = \mathbf{B}^{-1} \mathbf{T}^\top \mathbf{y}$  and  $\mathbf{V} = (\mathbf{v}_1, \dots, \mathbf{v}_m) = \mathbf{T}\mathbf{B}^{-\top}$ .

### KPLS and Lanczos Decompositions

$$\mathbf{L} = \mathbf{T}^\top \mathbf{X}\mathbf{W} = \begin{pmatrix} * & * & 0 & 0 & \dots & 0 \\ 0 & * & * & 0 & \dots & 0 \\ & & \vdots & \vdots & & \\ 0 & \dots & 0 & 0 & * & * \\ 0 & \dots & 0 & 0 & 0 & * \end{pmatrix} \in \mathbb{R}^{m \times m} \text{ and } \mathbf{D} = \mathbf{L}^\top \mathbf{L} \in \mathbb{R}^{m \times m}$$

The eigenvalues of  $\mathbf{D}$  are good approximations of the eigenvalues of  $\mathbf{K}$  (a) for the leading eigenvalues of  $\mathbf{K}$ , (b) if  $m \gg 0$ , (c) if the eigenvalues of  $\mathbf{K}$  decay fast [4].

### Approximate Degrees of Freedom in $\mathcal{O}(n^2)$

Compute  $\mathbf{D}$  for a large number of components  $m_{\max} \geq m$  and replace  $\mathbf{K}^j$  by  $\mathbf{D}^j$  in formula (2).

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

### Approximate Confidence Intervals in $\mathcal{O}(n^2)$

$$\hat{\mathbf{y}}(\mathbf{x}) \sim \mathcal{N}(\mathbf{k}(\mathbf{x})^\top E[\hat{\boldsymbol{\alpha}}], \sigma^2 \mathbf{k}(\mathbf{x})^\top \mathbf{H}_m \mathbf{H}_m^\top \mathbf{k}(\mathbf{x}))$$

with  $\mathbf{k}(\mathbf{x}) = (k(\mathbf{x}, \mathbf{x}_1), \dots, k(\mathbf{x}, \mathbf{x}_n)) \in \mathbb{R}^n$  and  $\sigma$  the noise level.

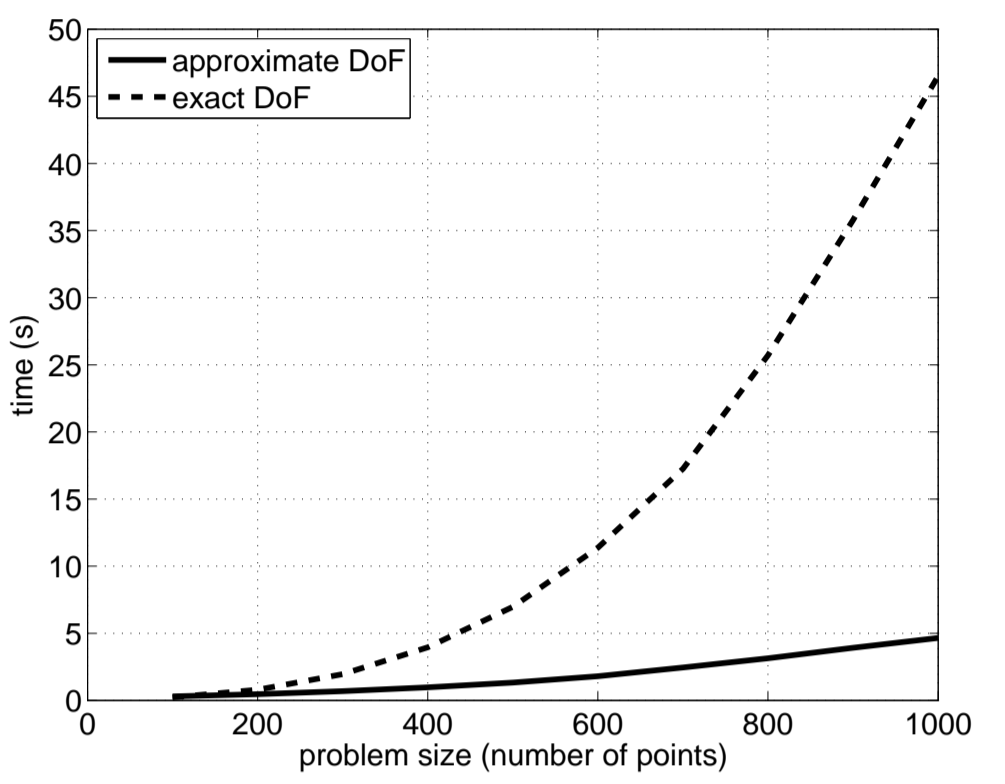
$$\mathbf{H}_m^\top \mathbf{k}(\mathbf{x}) = \sum_{j=1}^m \mathbf{K}^{j-1} \left\{ c_j (\mathbf{I}_n - \mathbf{K}\mathbf{T}\mathbf{N}\mathbf{R}^\top) + \mathbf{K}(\mathbf{y} - \hat{\mathbf{y}}_m) \mathbf{u}_j^\top \right\} \mathbf{k}(\mathbf{x}) + \mathbf{T}\mathbf{N}\mathbf{R}^\top \mathbf{k}(\mathbf{x}).$$

with  $\mathbf{r}_i = \tilde{\mathbf{r}}_i / \|\tilde{\mathbf{r}}_i\|_{\mathbf{K}}$  (residuals)  $\mathbf{N} = \text{diag}(1/\|\mathbf{K}\tilde{\mathbf{r}}_i\|)$  and  $\mathbf{U} = (\mathbf{u}_1, \dots, \mathbf{u}_m) = \mathbf{R}\mathbf{N}\mathbf{B}^{-\top}$ .

## Experiments

### Runtime Comparison

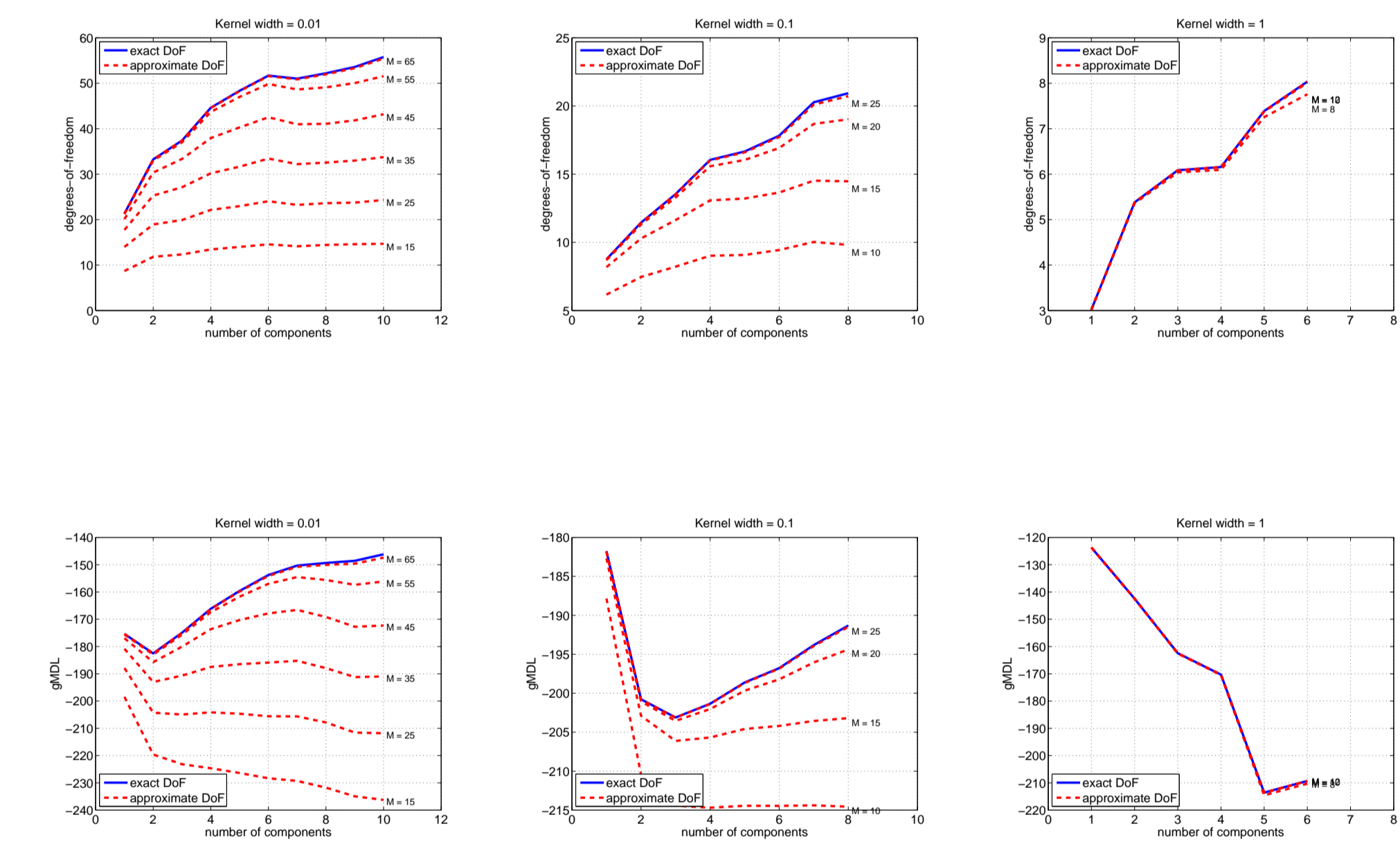
Comparison of runtime on the "kin" data set [1] ( $p = 8$  dimensions and  $n = 8192$  observations). Jagged line: KPLS with exact Degrees of Freedom for  $m = 10$  components. Solid line: KPLS with approximate Degrees of Freedom for  $m = 10$  components and  $m_{\max} = 30$  components for the approximation of the eigenvalues of the kernel matrix. Hence, for the approximation, the effective number of components is three times higher.



### Quality of the Approximation

Computation of KPLS with rbf-kernels on simulated data

$$f(\mathbf{x}) = \text{sinc}(\mathbf{x}) + \varepsilon, \varepsilon \sim \mathcal{N}(0, 0.1^2) \quad \mathbf{x}_i \sim \mathcal{U}[-\pi, \pi], i = 1, \dots, 100.$$



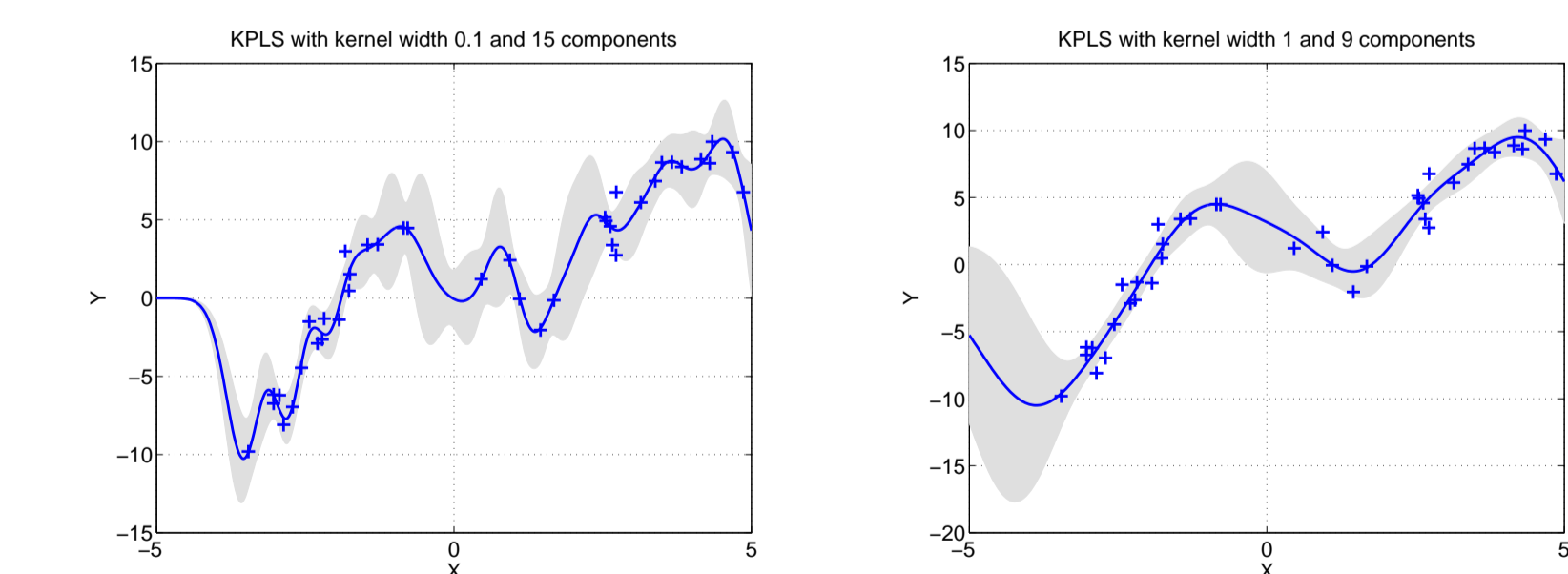
Results for kernel widths 0.01 (left), 0.1 (center), and 1 (right). Top row: DoF (blue line) and approximate DoF (red dashed line) for different numbers of maximal components. Bottom row: generalized Minimum Description Length (gMDL) (blue line) and approximate gMDL (red dashed line) for different numbers of maximal components.

### Approximate Confidence Intervals for KPLS

Computation of KPLS with rbf-kernels on simulated data

$$f(x) = (x-1)(x+2)(x-1.5) \exp(-x^2/10) + \varepsilon, \varepsilon \sim \mathcal{N}(0, 1)$$

$n = 40$  observations  $x_i$  are drawn from a mixture of  $\mathcal{N}(-2, 1)$  and  $\mathcal{N}(3, 1)$ .



Confidence intervals for two different kernel widths. Left: KPLS with 15 components and an rbf-kernel of width 0.1 and Right: KPLS with 9 components and an rbf kernel of width 1.

## References

- [1] Delve repository, <http://www.cs.toronto.edu/~delve/>.
- [2] N. Krämer and M.L. Braun, *Kernelizing PLS, Degrees of Freedom, and Efficient Model Selection*, Proceedings of the 24th International Conference on Machine Learning, 2007, pp. 441–448.
- [3] A. Phatak, P.M. Rille, and A. Penlidis, *The Asymptotic Variance of the Univariate PLS Estimator*, Linear Algebra and its Applications **354** (2002), 245–253.
- [4] Y. Saad, *Iterative methods for sparse linear systems*, 1st ed., PWS, 1996.