



Kernel PLS-SVC for Linear and Nonlinear Classification

Roman Rosipal^{1,3}, Leonard J. Trejo¹, Bryan Matthews²

¹*NASA Ames Research Center,
Computational Sciences Division, Moffett Field, CA*

²*QSS Group Inc., NASA Ames Research Center,
Computational Sciences Division, Moffett Field, CA*

³*Department of Theoretical Methods
Slovak Academy of Sciences, Bratislava, Slovak Republic*

Outline

1. Introduction to PLS
2. PLS and Kernel PLS Discrimination
(relation to CCA and Fisher's LDA)
3. Experimental Results

Partial Least Squares

- PLS - a class of techniques for modeling relations between blocks of observed variables by means of latent variables
- Herman Wold ('66,'81) - NIPALS - to linearize models nonlinear in the parameters
- Svante Wold et. al ('83) - extended PLS for the overdetermined regression problems
- Chemometrics - strong latent variable structure

Partial Least Squares

- data sets:

$$\mathbf{X} \ (n_{objects} \times N_{variables})$$

$$\mathbf{Y} \ (n_{objects} \times M_{responses})$$

– zero-mean

- bilinear decomposition:

$$\mathbf{X} = \mathbf{TP}^T + \mathbf{E}$$

$$\mathbf{Y} = \mathbf{UQ}^T + \mathbf{F}$$

where:

\mathbf{T}, \mathbf{U} matrix of score variables (LV, components)

\mathbf{P}, \mathbf{Q} matrix of loadings

\mathbf{E}, \mathbf{F} matrix of residuals (errors)

- PLS - bilinear decomposition of \mathbf{X} and \mathbf{Y} maximizing covariance between score vectors $\mathbf{t} = \mathbf{X}\mathbf{w}$ and $\mathbf{u} = \mathbf{Y}\mathbf{c}$

$$\begin{aligned} \max_{|\mathbf{r}|=|\mathbf{s}|=1} [\text{cov}(\mathbf{X}\mathbf{r}, \mathbf{Y}\mathbf{s})]^2 &= [\text{cov}(\mathbf{X}\mathbf{w}, \mathbf{Y}\mathbf{c})]^2 \\ &= \text{var}(\mathbf{X}\mathbf{w}) [\text{corr}(\mathbf{X}\mathbf{w}, \mathbf{Y}\mathbf{c})]^2 \text{var}(\mathbf{Y}\mathbf{c}) \\ &= [\text{cov}(\mathbf{t}, \mathbf{u})]^2 \end{aligned}$$

- | | |
|---|--|
| $\mathbf{X}^T \mathbf{Y} \mathbf{Y}^T \mathbf{X} \mathbf{w} = \lambda \mathbf{w}$
$\mathbf{t} = \mathbf{X} \mathbf{w}$ | $\mathbf{X} \mathbf{X}^T \mathbf{Y} \mathbf{Y}^T \mathbf{t} = \lambda \mathbf{t}$
$\mathbf{u} = \mathbf{Y} \mathbf{Y}^T \mathbf{t}$ |
|---|--|

- $\mathbf{p} = \mathbf{X}^T \mathbf{t} / (\mathbf{t}^T \mathbf{t})$; $\mathbf{q} = \mathbf{Y}^T \mathbf{u} / (\mathbf{u}^T \mathbf{u})$

- deflation schemes define different forms of PLS

Forms of Partial Least Squares

- PLS1, PLS2: rank-one approximation on \mathbf{X}, \mathbf{Y} with \mathbf{t}
 $\mathbf{X} \rightarrow \mathbf{X} - \mathbf{t}\mathbf{p}^T$; $\mathbf{Y} \rightarrow \mathbf{Y} - \mathbf{t}\mathbf{c}^T$
mutually orthogonal components \mathbf{t}_i , $i = 1, \dots, p$
- PLS-SB: SVD on $\mathbf{Y}^T \mathbf{X}$
mutually orthogonal weight vectors $\mathbf{w}_i, \mathbf{c}_i$
generally not orthogonal \mathbf{t}_i and \mathbf{u}_i
- 1st $SV_{i+1} \geq 2\text{nd } SV_i \rightarrow$ we select one component at a time

Partial Least Squares Discrimination

$$\mathbf{Y} = \begin{pmatrix} \mathbf{1}_{n_1} & \mathbf{0}_{n_1} & \cdots & \mathbf{0}_{n_1} \\ \mathbf{0}_{n_2} & \mathbf{1}_{n_2} & \cdots & \mathbf{0}_{n_2} \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{0}_{n_g} & \mathbf{0}_{n_g} & \cdots & \mathbf{1}_{n_{g-1}} \\ \mathbf{0}_{n_g} & \mathbf{0}_{n_g} & \cdots & \mathbf{0}_{n_g} \end{pmatrix}$$

- orthonormalized PLS

$$\tilde{\mathbf{Y}} = \mathbf{Y}(\mathbf{Y}^T \mathbf{Y})^{-1/2}$$

$$\tilde{\mathbf{Y}}^T \tilde{\mathbf{Y}} = \mathbf{I}$$

Orthonormalized PLS vs. CCA, Fisher's LDA

[Barker & Rayens 2003]

- orthonormalized PLS

$$\max_{|\mathbf{r}|=|\mathbf{s}|=1} [\text{cov}(\mathbf{X}\mathbf{r}, \tilde{\mathbf{Y}}\mathbf{s})]^2 = \text{var}(\mathbf{X}\mathbf{w}) [\text{corr}(\mathbf{X}\mathbf{w}, \tilde{\mathbf{Y}}\mathbf{c})]^2$$

$$\mathbf{X}^T \tilde{\mathbf{Y}} \tilde{\mathbf{Y}}^T \mathbf{X} \mathbf{w} = \lambda \mathbf{w}$$

$$\mathbf{X}^T \mathbf{Y} (\mathbf{Y}^T \mathbf{Y})^{-1} \mathbf{Y}^T \mathbf{X} \mathbf{w} = \lambda \mathbf{w}$$

$$\mathbf{H} \mathbf{w} = \lambda \mathbf{w}$$

- CCA, Fisher's LDA

$$\max_{|\mathbf{r}|=|\mathbf{s}|=1} [\text{corr}(\mathbf{X}\mathbf{r}, \mathbf{Y}\mathbf{s})]^2 = [\text{corr}(\mathbf{X}\mathbf{a}, \mathbf{Y}\mathbf{b})]^2$$

$$(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y} (\mathbf{Y}^T \mathbf{Y})^{-1} \mathbf{Y}^T \mathbf{X} \mathbf{a} = \lambda \mathbf{a}$$

$$\mathbf{E}^{-1} \mathbf{H} \mathbf{a} = \frac{\lambda}{1-\lambda} \mathbf{a}$$

Kernel PLS Discrimination

- linear PLS discrimination in a feature space \mathcal{F}
- nonlinear kernel-based PLS:

$$\begin{aligned} \mathbf{X}\mathbf{X}^T\mathbf{Y}\mathbf{Y}^T\mathbf{t} &= \lambda\mathbf{t} \\ \mathbf{u} &= \mathbf{Y}\mathbf{Y}^T\mathbf{t} \end{aligned}$$

\Rightarrow

$$\begin{aligned} \mathbf{K}\mathbf{Y}\mathbf{Y}^T\mathbf{t} &= \lambda\mathbf{t} \\ \mathbf{u} &= \mathbf{Y}\mathbf{Y}^T\mathbf{t} \end{aligned}$$

- nonlinear kernel-based orthonormalized PLS:

$$\begin{aligned} \mathbf{K}\mathbf{Y}(\mathbf{Y}^T\mathbf{Y})^{-1}\mathbf{Y}^T\mathbf{t} &= \mathbf{K}\tilde{\mathbf{Y}}\tilde{\mathbf{Y}}^T\mathbf{t} = \lambda\mathbf{t} \\ \tilde{\mathbf{Y}} &= \mathbf{Y}(\mathbf{Y}^T\mathbf{Y})^{-1/2} \end{aligned}$$

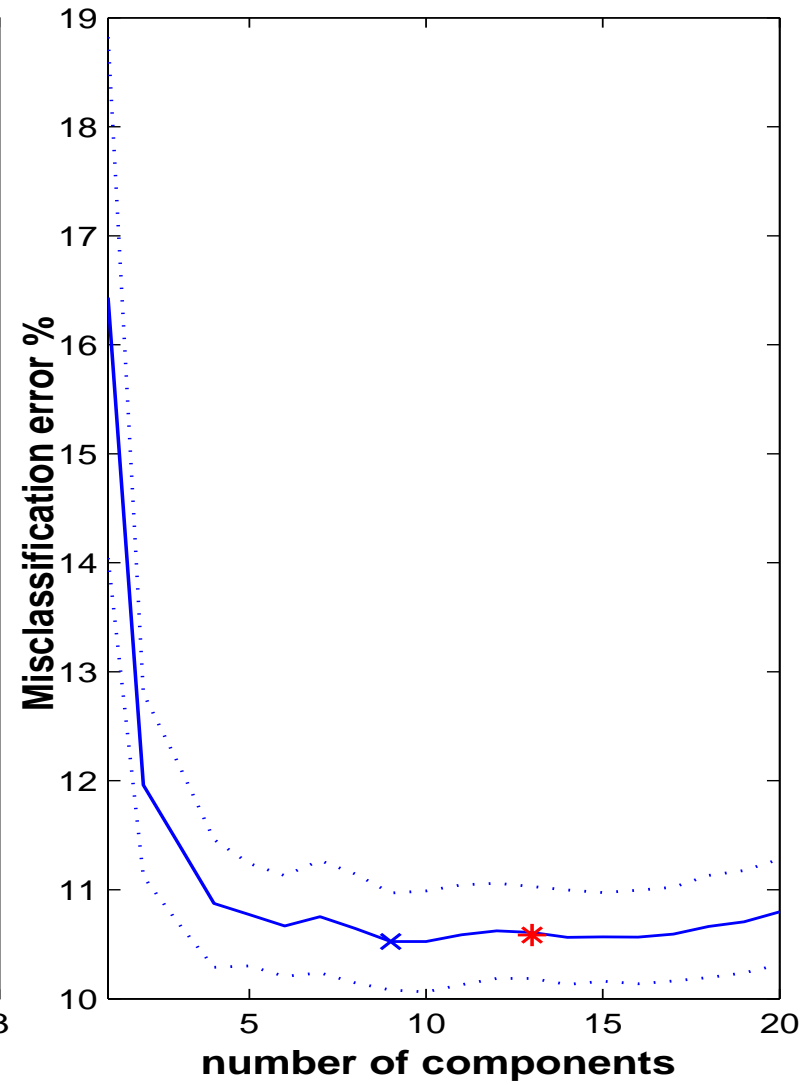
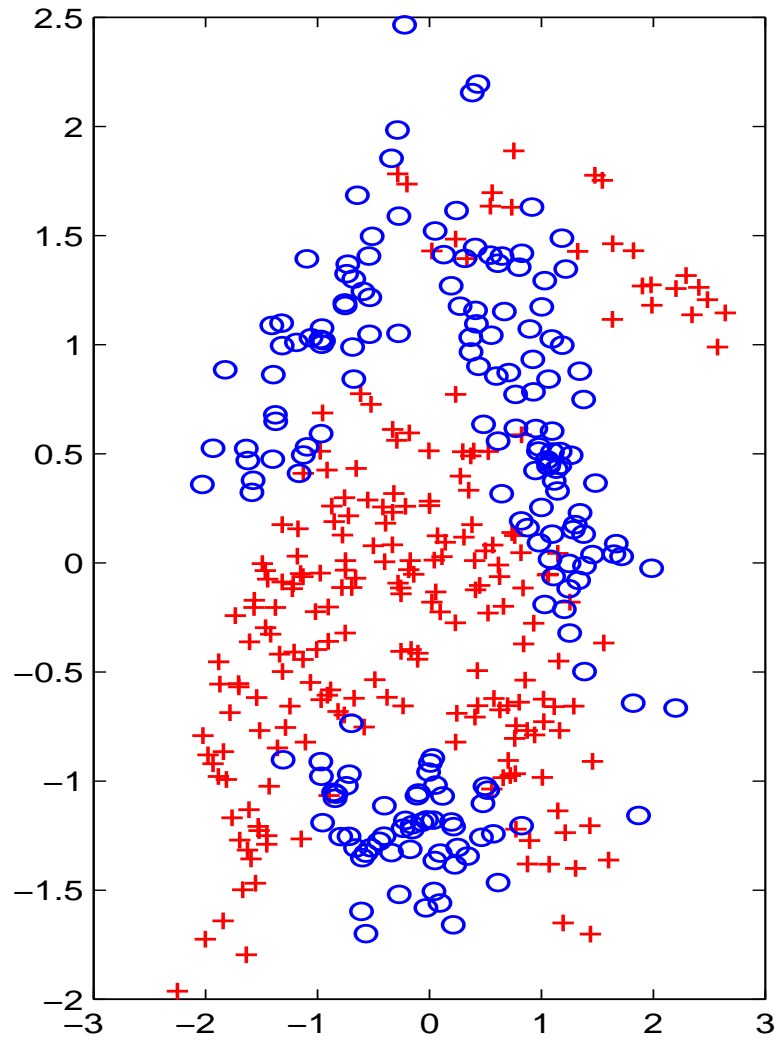
Kernel PLS-SVC Classification

- orthonormalized kernel PLS + SVC (KPLS-SVC)
- orthonormalized kernel PLS can be combined with other existing classifiers (e.g. LDA, logistic regression)

Experiments:

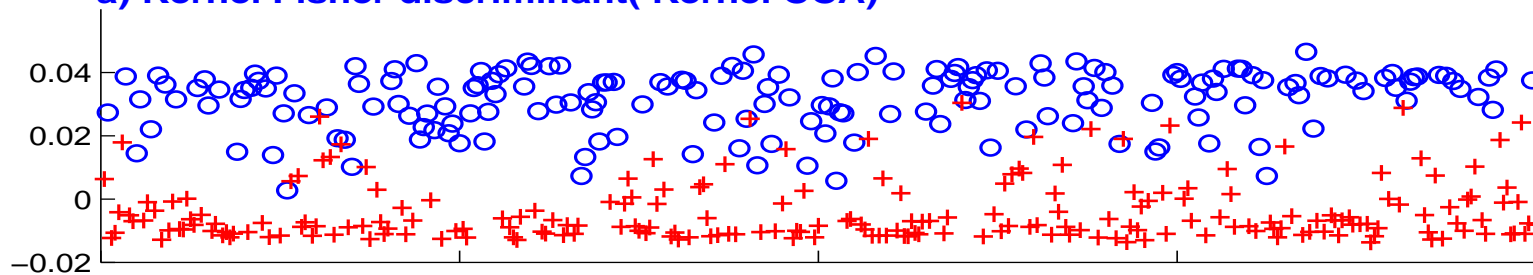
- 13 benchmark data sets of two-class classification problem
<http://www.first.gmd.de/~raetsch>
- vowel sounds data set - multi-class problem (11 classes)
- classification of finger movement periods from non-movement periods based on electroencephalograms (EEG)

Banana data set

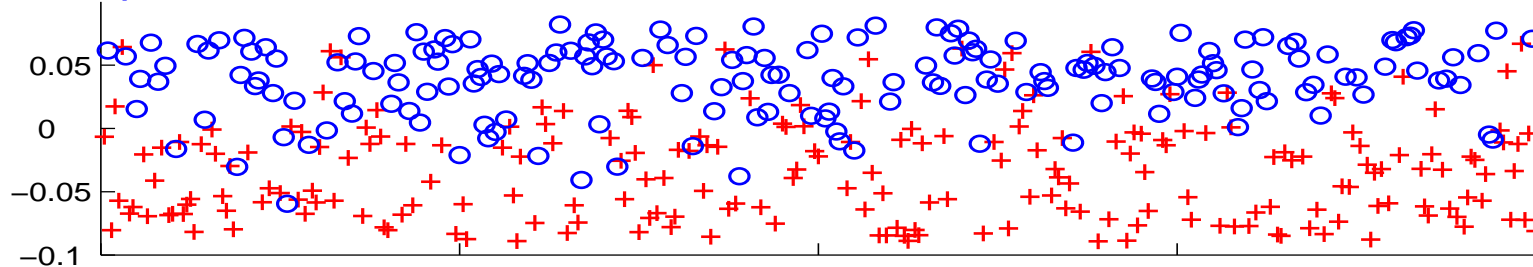


Data projection onto direction given by:

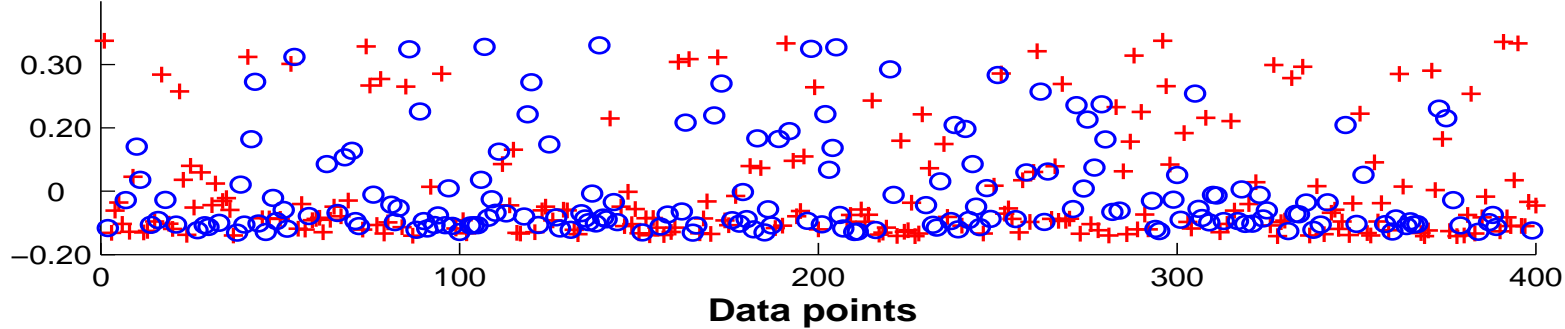
a) Kernel Fisher discriminant(Kernel CCA)

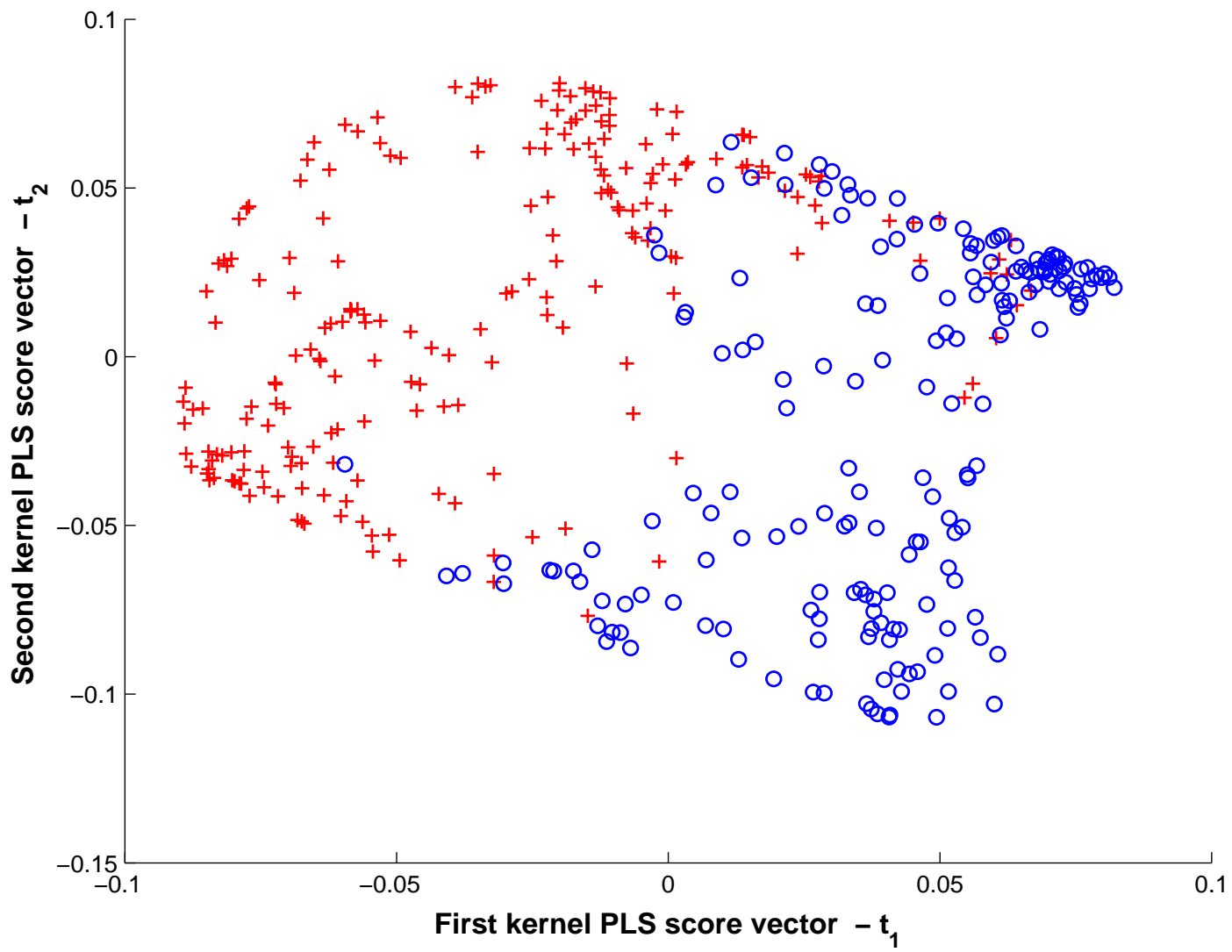


b) First kernel PLS score vector



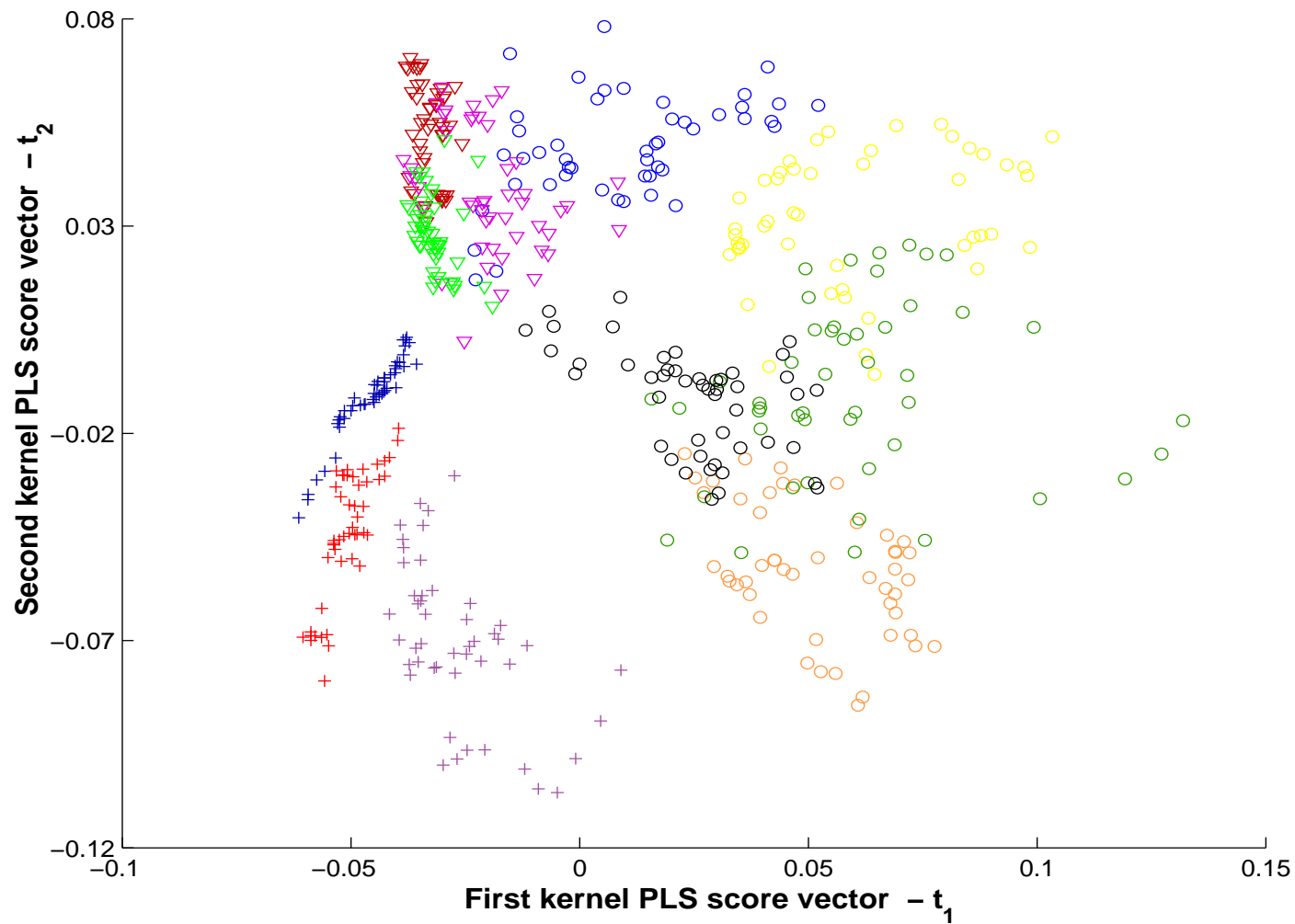
c) First kernel PCA principal component





Data Set	KFD	C-SVC	KPLS-SVC
Banana	10.8±0.5	11.5±0.5	10.5±0.4
B.Cancer	25.8±4.6	26.0±4.7	25.1±4.5*
Diabetes	23.2±1.6	23.5±1.7	23.0±1.7
German	23.7±2.2	23.6±2.1	23.5±1.6
Heart	16.1±3.4	16.0±3.3	16.5±3.6
Image	4.76±0.58	2.96±0.60	3.03±0.61
Ringnorm	1.49±0.12	1.66±0.12	1.43±0.10
F.Solar	33.2±1.7	32.4±1.8	32.4±1.8
Splice	10.5±0.6	10.9±0.7	10.9±0.8
Thyroid	4.20±2.07	4.80±2.19	4.39±2.10
Titanic	23.2±2.06	22.4±1.0	22.4±1.1*
Twonorm	2.61±0.15	2.96±0.23	2.34±0.11
Waveform	9.86±0.44	9.88±0.43	9.58±0.36

Vowel sounds data set: 11 classes, 10 predictors

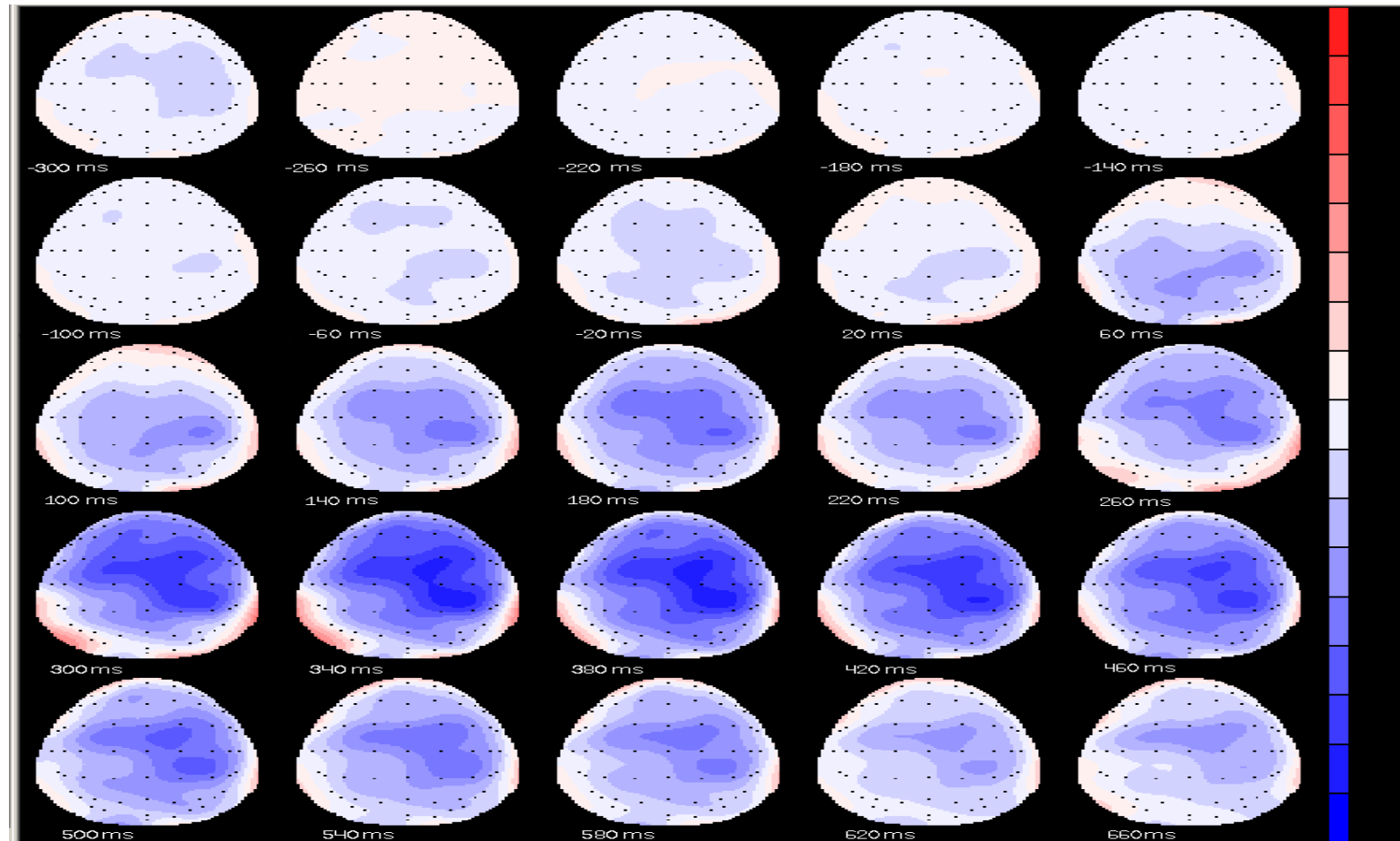


Method	Training Error	Testing Error
LDA	0.32	0.56
SVC (linear) - 1vs1	0.19	0.51
KPLS-SVC (linear) - 1vs1	0.16	0.47
FDA/MARS (df=2)	0.02	0.42
FDA/MARS (df=6,red. dim.)	0.13	0.39
SVC (gauss) - 1vs1	0.01	0.37
KPLS-SVC (gauss) - 1vs1	0.01	0.35
SVC (gauss, $w \leq 5$) - 1vs1	0.002	0.29
KPLS-SVC (gauss, $w \leq 5$) - 1vs1	0.002	0.33

Finger movement periods vs. non-movement periods

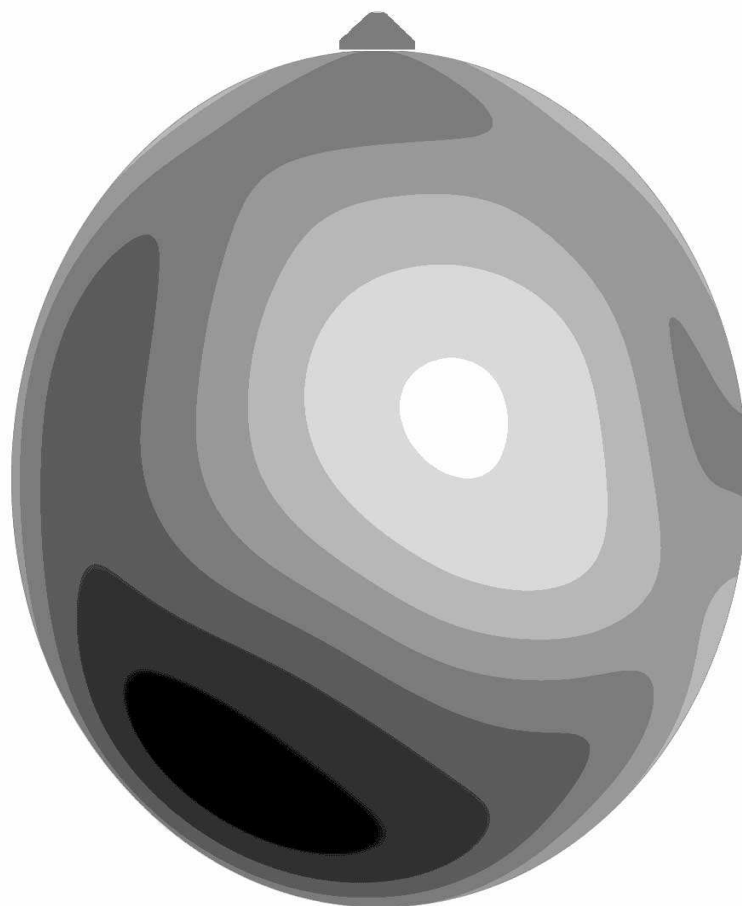


PLS-derived Spatio-temporal Filter - 01/09/2003

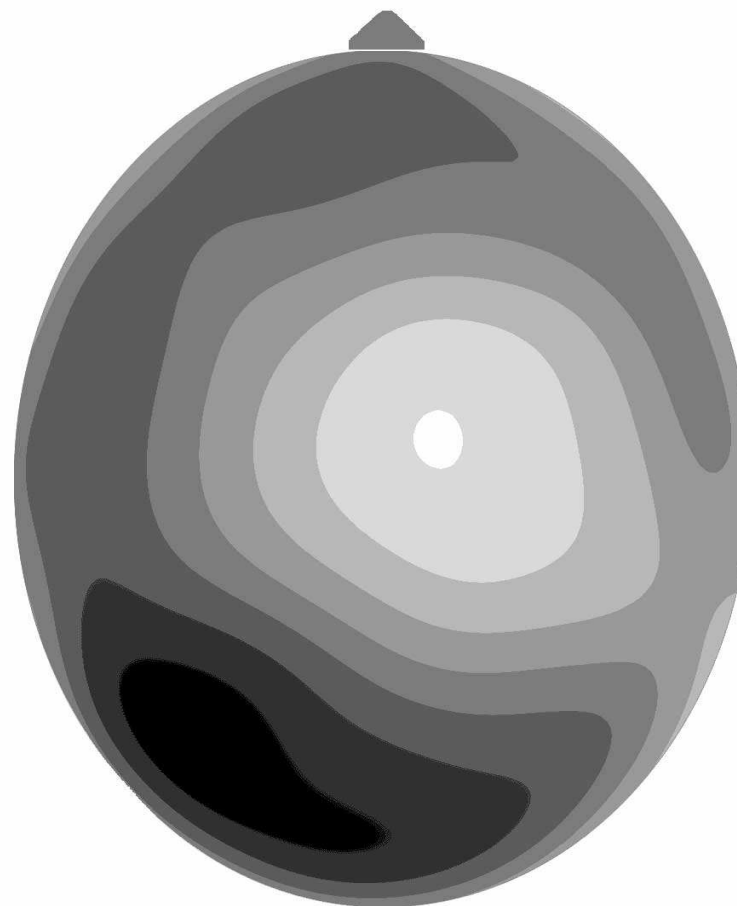


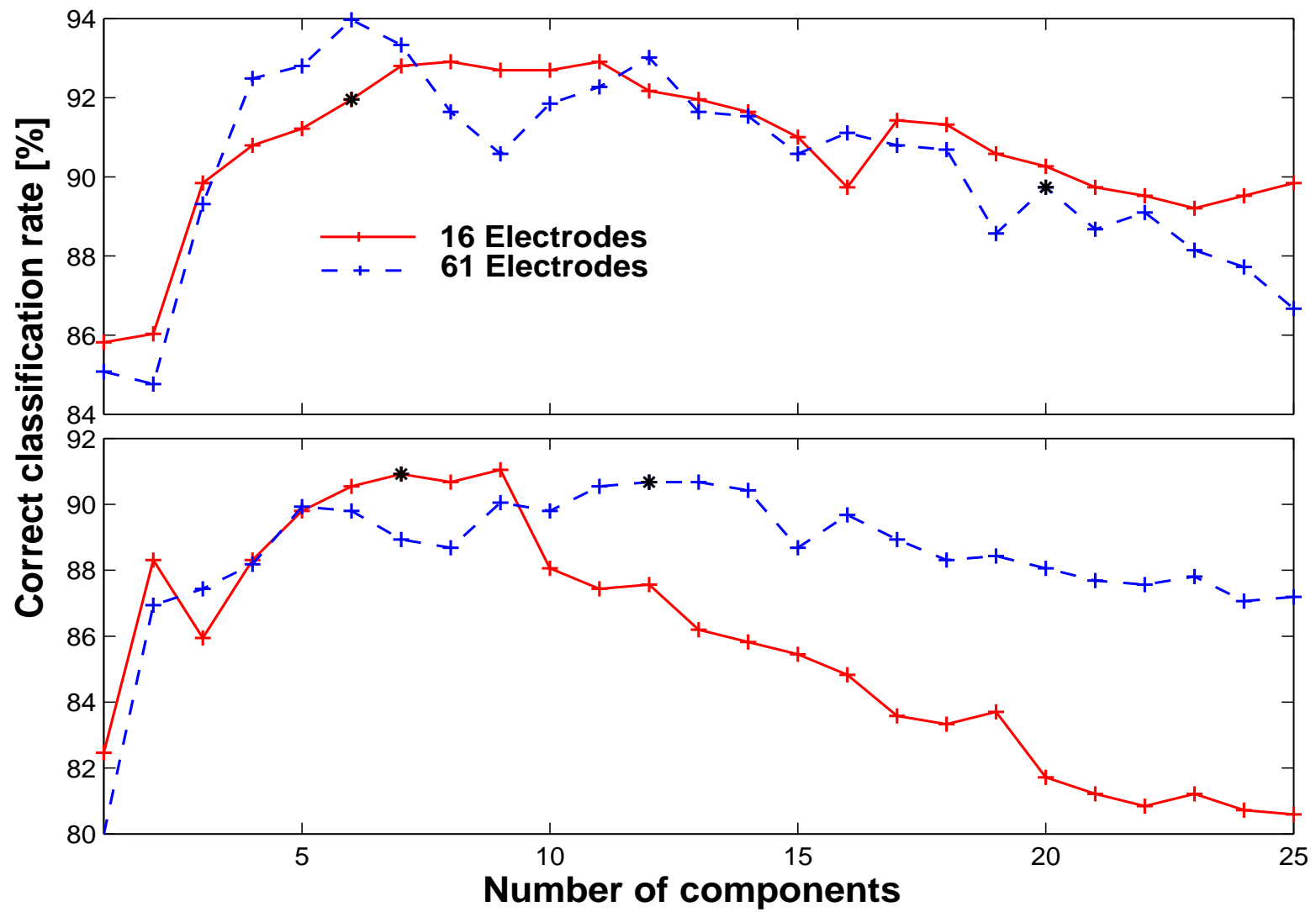
PLS-derived Spatio-temporal Filter
(370ms after button press)

11/14/2002



01/09/2003





Conclusions

- PLS discrimination - useful method for dimensionality reduction, visualization
- PLS discrimination preferred over PCA
- nonlinear kernel-based version of PLS discrimination provided
- KPLS-SVC achieved good results on the used data sets

?