

**Kernel Versions of Some Orthogonal Transformations**

Allan A. Nielsen  
Technical University of Denmark  
Applied Mathematics and Computer Science  
people.compute.dtu.dk/alan  
alan@dtu.dk

Symposium i anvendt Statistik, 26.-28. januar 2015

**Intro**

- Kernel versions of PCA and MAF/MNF analysis
- Parsimonious representation of multi- or hypervariate data
- Interesting projections in feature space with different measures of ‘interestingness’: max variance, autocorrelation, signal-to-noise ratio

**Data matrix**

- $n$  by  $p$  data matrix or design matrix  $X = \begin{bmatrix} x_1^T \\ x_2^T \\ \vdots \\ x_n^T \end{bmatrix}$
- Each row consists of a vector of measurements  $x_i^T$  from  $p$  variables for a particular observation;  $X$  is often column centered

**PCA**

- R-mode or primal analysis: decompose  $p$  by  $p$  variance-covariance matrix  $S = X^T X / (n - 1) = 1 / (n - 1) \sum_{i=1}^n x_i x_i^T$
- $\frac{1}{n-1} X^T X u_i = \lambda_i u_i$
- Projections or scores are  $x_i^T u_i$
- Variance of scores maximized

$$u_i^T S u_i = \lambda_i u_i^T u_i = \lambda_i$$

**PCA**

- Q-mode or dual analysis: decompose  $n$  by  $n$  (Gram) matrix  $XX^T / (n - 1)$ , multiply from left with  $X$ ,  $v_i \propto Xu_i$

$$\frac{1}{n-1} XX^T (Xu_i) = \lambda_i (Xu_i)$$

$$\frac{1}{n-1} XX^T v_i = \lambda_i v_i$$


**PCA**

- Multiply from left with  $X^T$

$$\frac{1}{n-1} X^T X (X^T v_i) = \lambda_i (X^T v_i)$$

- Solutions related by (Eckart-Young, 1936)

$$u_i = X^T v_i / \sqrt{(n-1)\lambda_i}$$

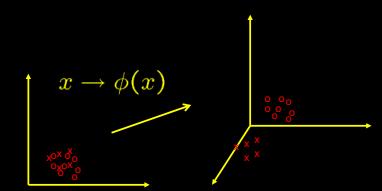
$$v_i = Xu_i / \sqrt{(n-1)\lambda_i}$$

**Gram matrix**

- Consists of inner products only

$$XX^T = \begin{bmatrix} x_1^T x_1 & x_1^T x_2 & \cdots & x_1^T x_n \\ x_2^T x_1 & x_2^T x_2 & \cdots & x_2^T x_n \\ \vdots & \vdots & \ddots & \vdots \\ x_n^T x_1 & x_n^T x_2 & \cdots & x_n^T x_n \end{bmatrix}$$

**Nonlinear mapping**



**Nonlinear mapping**

- Example:  $x = [z_1 \ z_2]^T$  mapped into  $\phi(x) = [z_1 \ z_2 \ z_1^2 \ z_2^2 \ z_1 z_2]^T$
- Maps original 2-dimensional vector into a 5-dimensional feature space so that for example a linear decision rule becomes general enough to differentiate between all linear and quadratic forms including ellipsoids

## Nonlinear mapping

- Data matrix  $X \rightarrow \Phi = \begin{bmatrix} \phi(x_1)^T \\ \phi(x_2)^T \\ \vdots \\ \phi(x_n)^T \end{bmatrix}$
- Each row consists of a  $q$ -vector of nonlinear mappings  $\phi(x_i)^T$ ,  $q \geq p$

## Nonlinear mapping

- Maps Gram matrix into  $\Phi\Phi^T = \begin{bmatrix} \phi(x_1)^T\phi(x_1) & \phi(x_1)^T\phi(x_2) & \cdots & \phi(x_1)^T\phi(x_n) \\ \phi(x_2)^T\phi(x_1) & \phi(x_2)^T\phi(x_2) & \cdots & \phi(x_2)^T\phi(x_n) \\ \vdots & \vdots & \ddots & \vdots \\ \phi(x_n)^T\phi(x_1) & \phi(x_n)^T\phi(x_2) & \cdots & \phi(x_n)^T\phi(x_n) \end{bmatrix}$
- Consists of inner products of mappings only, still  $n$  by  $n$

## PCA

- Q-mode or dual analysis: decompose  $n$  by  $n$  (Gram) matrix  $\Phi\Phi^T/(n-1)$

$$\frac{1}{n-1}\Phi\Phi^Tv_i = \lambda_i v_i$$

- Solutions related by

$$u_i = \Phi^Tv_i/\sqrt{(n-1)\lambda_i}$$

$$v_i = \Phi u_i/\sqrt{(n-1)\lambda_i}$$

## Kernel trick

- We need not know the actual mappings (let  $\lambda_i$  subsume  $n-1$ )

$$\Phi\Phi^Tv_i = \lambda_i v_i$$

$$Kv_i = \lambda_i v_i$$

- (Mercer) kernel matrix with elements  $K_{ij} = \kappa(x_i, x_j)$

## Kernel solution

- Solutions to  $\Phi\Phi^T = \lambda_i v_i$  related by  $u_i = \Phi^Tv_i/\sqrt{\lambda_i}$   
 $v_i = \Phi u_i/\sqrt{\lambda_i}$
- Scores are  $\phi(x)^T u_i$

$$\phi(x)^T u_i = [\kappa(x, x_1) \ \kappa(x, x_2) \ \cdots \ \kappa(x, x_n)] v_i / \sqrt{\lambda_i}$$

- i.e., expressed by kernel elements (memory based, non-para.)

## Kernel solution

- Several basic properties may all be expressed in terms of the kernel function without using the mapping  $\phi(x)$  explicitly, including
  - the norm in kernel feature space
  - the distance between observations in feature space
  - the norm of the mean in feature space
  - **centering to zero mean in feature space**
  - standardisation to unit variance in feature space

## Kernel functions

- Popular kernels
  - Stationary, depend on vector difference, invariant under translation  
 $\kappa(x_i, x_j) = \kappa(x_i - x_j)$
  - Homogeneous, RBF  
 $\kappa(x_i, x_j) = \kappa(\|x_i - x_j\|)$

## Kernel functions

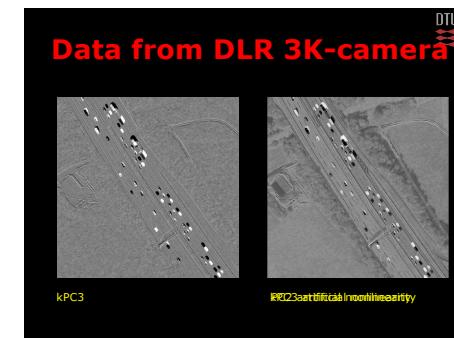
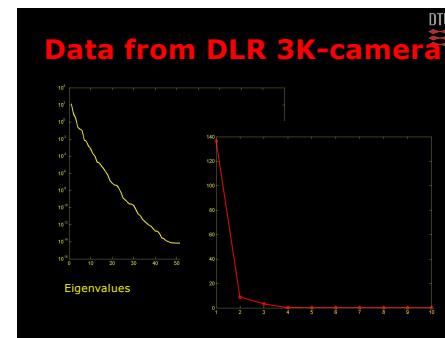
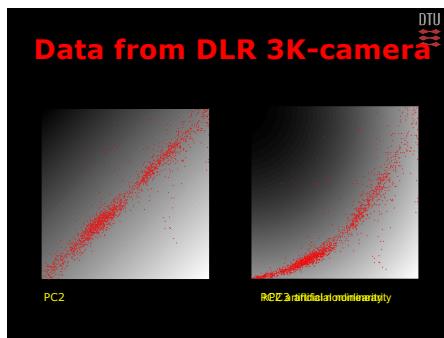
- Often used RBFs,  $h = \|x_i - x_j\|$ 
  - Multiquadric  
 $\kappa(h) = (h^2 + h_0^2)^{1/2}$
  - Inverse multiquadric  
 $\kappa(h) = (h^2 + h_0^2)^{-1/2}$
  - Gaussian  
 $\kappa(h) = \exp(-\frac{1}{2}(h/h_0)^2)$

## Data from DLR 3K-camera




t<sub>1</sub> as R, t<sub>2</sub> as GB

Spatial sub-sample



**MAF**

- R-mode or primal analysis: max autocorrelation of linear combinations  
 $a^T x(r)$

$$\begin{aligned} \rho &= 1 - \frac{1}{2} \frac{a^T S_{\Delta} a}{a^T S a} \\ &= 1 - \frac{1}{2} \frac{a^T X_{\Delta}^T X_{\Delta} a}{a^T X^T X a} \end{aligned}$$

**MAF**

- Q-mode or dual analysis: set  $a \propto X^T b$  and kernelize

$$\begin{aligned} \rho &= 1 - \frac{1}{2} \frac{b^T \Phi \Phi_{\Delta}^T \Phi_{\Delta} \Phi^T b}{b^T \Phi \Phi^T \Phi \Phi^T b} \\ &= 1 - \frac{1}{2} \frac{b^T K_{\Delta} K_{\Delta}^T b}{b^T K^2 b} \end{aligned}$$

**MAF kernel solution**

- Solutions related by  
 $a_i = \Phi^T b_i$   
 $b_i = \Phi a_i$
- Scores are  $\phi(x)^T a_i$   
 $\phi(x)^T a_i = [\kappa(x, x_1) \ \kappa(x, x_2) \ \cdots \ \kappa(x, x_n)] b_i$
- i.e., expressed by kernel elements (memory based, non-para.)

**MNF**

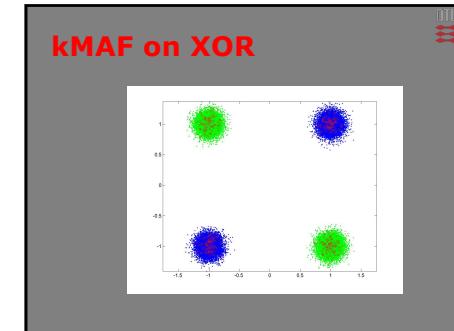
- R-mode or primal analysis: min NF (or max SNR) of linear combinations  
 $a^T x(r)$

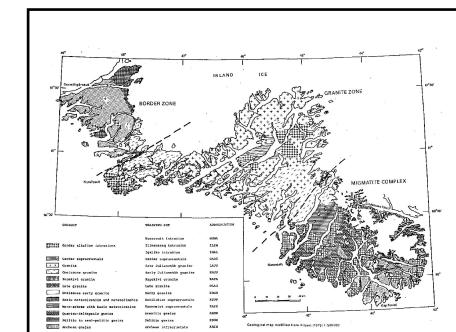
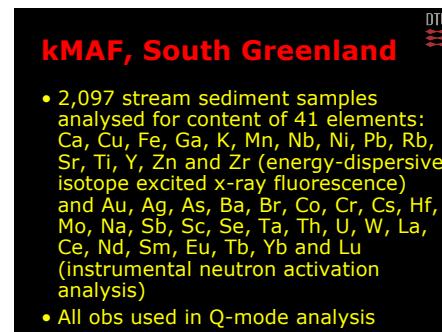
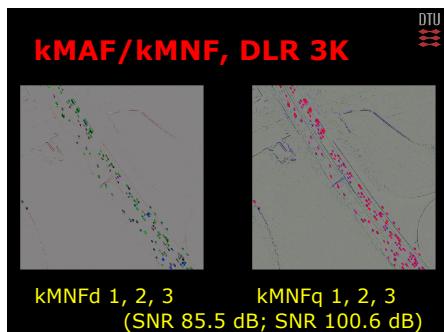
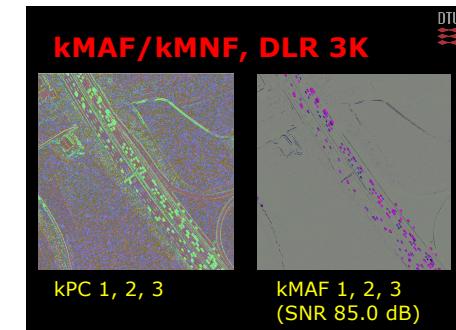
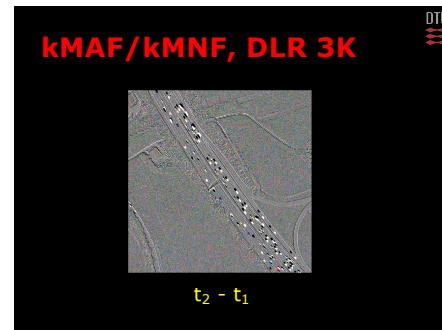
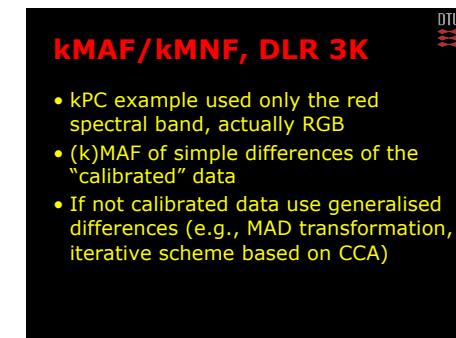
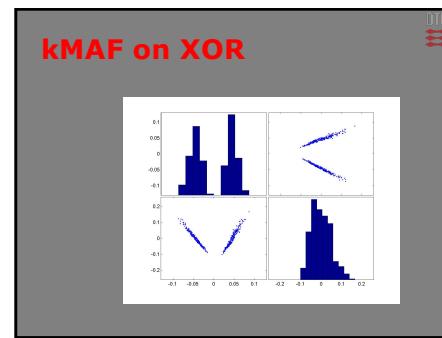
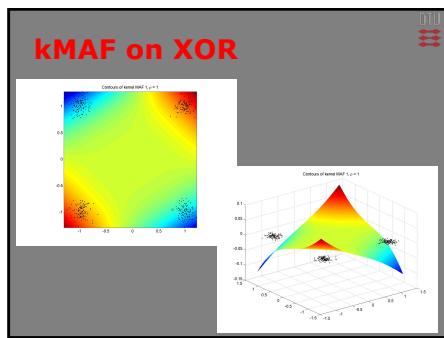
$$\begin{aligned} \frac{1}{\text{NF}} &= \frac{a^T S a}{a^T S_N a} \\ &= \frac{a^T X^T X a}{a^T X_N^T X_N a} \end{aligned}$$

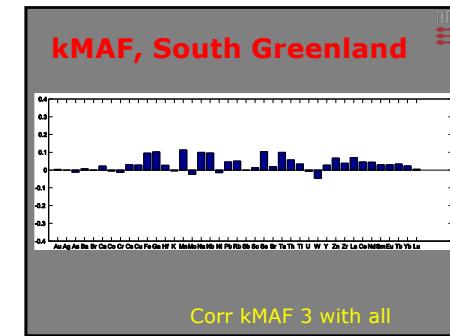
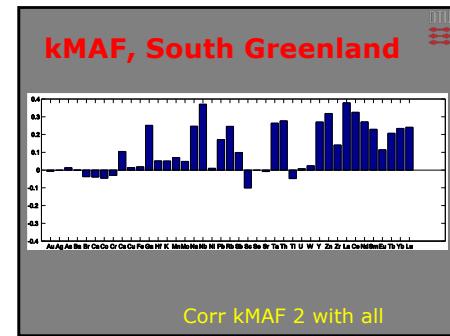
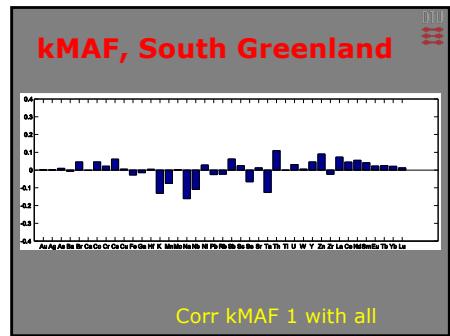
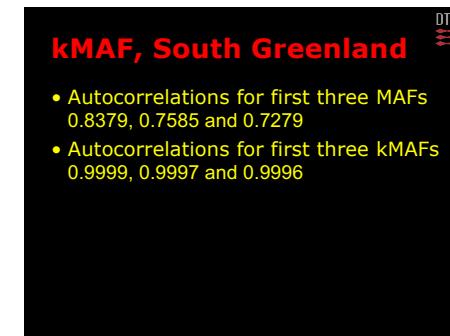
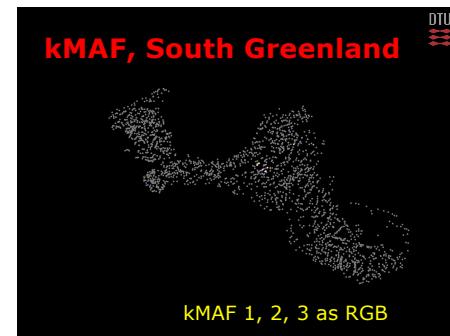
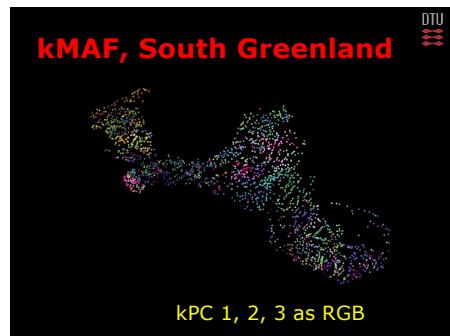
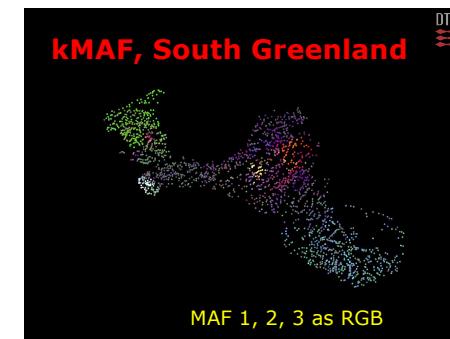
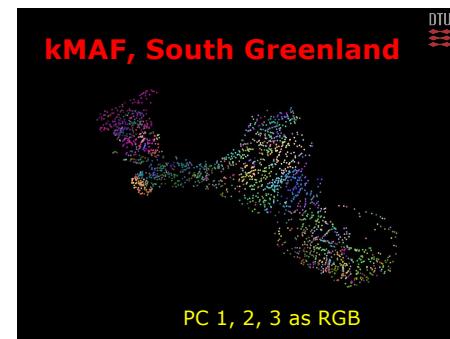
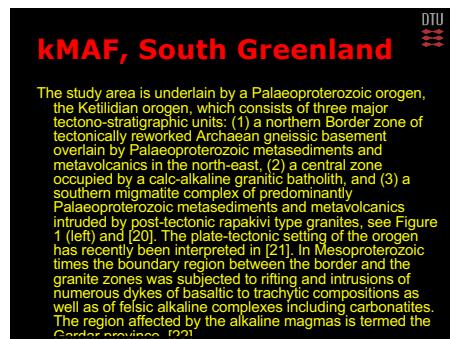
**MNF**

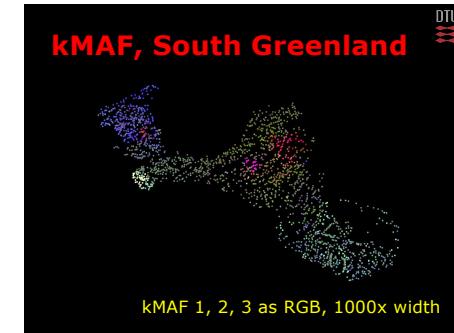
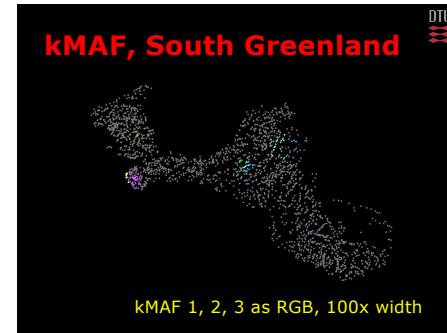
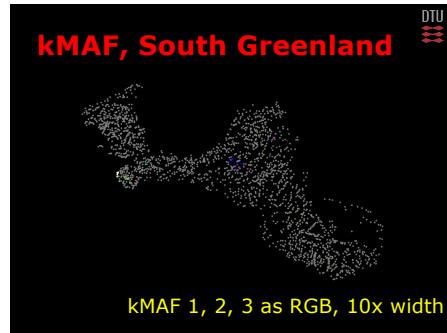
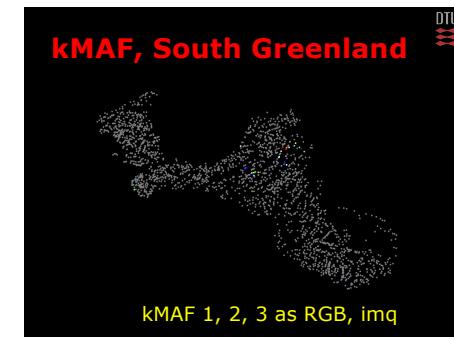
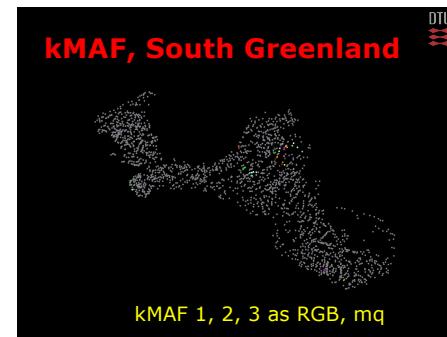
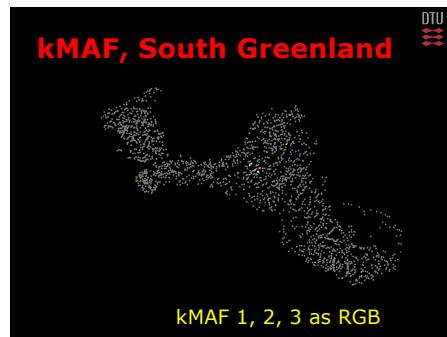
- Q-mode or dual analysis: set  $a \propto X^T b$  and kernelize

$$\begin{aligned} \frac{1}{\text{NF}} &= \frac{b^T \Phi \Phi^T \Phi \Phi^T a}{b^T \Phi \Phi_N^T \Phi_N \Phi^T b} \\ &= \frac{b^T K^2 b}{b^T K_N K_N^T b} \end{aligned}$$









**kMAF, South Greenland**

- Partial reconstruction of simplified geological map: nice, makes one trust the technique
- Potentially more important: where does the analysis not agree with conventional wisdom?

## Future work

- Kernel function
- Scale parameter, kernel width
- Dependence on particular training samples
- Parallel implementation (test data part), or run on graphics board GPU
- Pre-image problem
- ...

## Conclusions

- Kernel orthogonalisation is based on Q-mode analysis
- Data are replaced by nonlinear mappings
- Inner products of nonlinear mappings in turn replaced by kernel function, 'kernel trick'
- kPCA unlike linear PCA successfully finds change observations in a DLR 3K case where nonlinearities are introduced artificially

## Conclusions

- kMAF/kMNF unlike kPCA successfully finds change observations in DLR 3K colour image data
- (k)MAF unlike (k)PCA successfully finds relevant geological structures
- kMAF varies with increasing kernel width find geologically relevant regions of increasing scale; focus on extreme observations – smooth background esp. for small kernel widths

## References

- H. Hotelling, "Analysis of a complex of statistical variables into principal components," *Journal of Educational Psychology* **24**, 417-441 and 498-520, 1933.
- P. Switzer and A. A. Green, "Min/max autocorrelation factor analysis", Technical Report 6, Department of Statistics, Stanford University, 1984.
- A. A. Green, M. Berman, P. Switzer, and M. D. Craig, "A transformation for ordering multispectral data in terms of image quality with implications for noise removal", *IEEE Transactions on Geoscience and Remote Sensing* **26**(1), 65-74 (1988).
- B. Schölkopf, A. Smola, and K.-R. Müller, "Nonlinear component analysis as a kernel eigenvalue problem," *Neural Computation* **10**(5), 1299-1319 (1998).
- J. Shawe-Taylor and N. Cristianini, *Kernel Methods for Pattern Analysis*, Cambridge University Press (2004).
- C. M. Bishop, *Pattern Recognition and Machine Learning*, Springer (2006).

## References

- H. Hotelling, "Analysis of a complex of statistical variables into principal components," *Journal of Educational Psychology* **24**, 417-441 and 498-520, 1933.
- P. Switzer and A. A. Green, "Min/max autocorrelation factor analysis", Technical Report 6, Department of Statistics, Stanford University, 1984.
- A. A. Green, M. Berman, P. Switzer, and M. D. Craig, "A transformation for ordering multispectral data in terms of image quality with implications for noise removal", *IEEE Transactions on Geoscience and Remote Sensing* **26**(1), 65-74 (1988).
- A. A. Nielsen, K. Conradsen, J. L. Pedersen and A. Steenfelt, "Spatial Factor Analysis of Stream Sediment Geochemistry Data from South Greenland," in *Vera Pawlowsky-Glahn (Ed.) Proceedings of the Third Annual Conference of the International Association for Mathematical Geology, IAGG'97*, Barcelona, Spain, 22-27 September 1997, pp. 95-96.
- B. Schölkopf, A. Smola, and K.-R. Müller, "Nonlinear component analysis as a kernel eigenvalue problem," *Neural Computation* **10**(5), 1299-1319 (1998).
- A. A. Nielsen, K. Conradsen, J. L. Pedersen and A. Steenfelt, "Maximum Autocorrelation Factorial Kriging," in *W.J. Kleingeld and D.G. Krige (editors)*, *Proceedings of the 6th International Geostatistics Congress, Geostats 2000*, pp. 538-547, Cape Town, South Africa, 10-14 April 2000.
- J. Shawe-Taylor and N. Cristianini, *Kernel Methods for Pattern Analysis*, Cambridge University Press, 2004.
- L. Bottou, *Machine Learning and Pattern Recognition*, Springer, 2004.
- A. A. Nielsen and M. J. Carr, "Kernel principal component and maximum autocorrelation factor analysis for change detection," *SPIE Europe Remote Sensing Conference*, Berlin, Germany, 31 August-3 September 2009.
- A. A. Nielsen, "Kernel maximum autocorrelation factor and minimum noise fraction transformations," *IEEE Transactions on Image Processing* **20**(3), 612-624, 2011.