

why not change the world?<sup>®</sup>

## Projection-based Chemometrics and Deep Reconstruction

## **Dr. Uwe Kruger**

Department of Biomedical Engineering Jonsson Engineering Center Rensselaer Polytechnic Institute

## **Presentation Outline**

- Motivation for kernel-based methods (kernel density estimation)
- Principal Component Analysis (PCA) and Kernel principal component analysis (KPCA)
- Partial Least Squares (PLS) and Kernel partial least squares (KPLS)
- Some ideas on how to integrate nonlinear projectionbased methods for network pruning and detecting/diagnosing anomalies.



- Let's examine a very simple approach to motivate Cover's theorem and the idea behind reproducing kernels:
- How can we estimate the cumulative distribution function of a random variable X using a set of n observations drawn from the distribution of X?
- Let's try the following naïve estimator:

$$\hat{F}(x) = \frac{\# x_i \le x}{n} = \frac{S(x)}{n}$$



**Dr. Uwe Kruger Slide 3** Projection-Based Data Chemometrics and Deep Reconstruction Troy, November 19., 2017



- OK, the *n* observations, if assumed to be drawn independently, can be used to formulate a total of *n* Bernoulli trials (like flipping a coin)
  - two outcomes, the value can be larger or smaller than x;
  - the probability to be smaller then x (success) is equal to the cumulative probability distribution function for x, i.e. F(x); and
  - for the *i*th draw (drawing the *i*th value of the random variable *X*), the probability that  $x_i$  is smaller than or equal to x is F(x) for  $1 \le i \le n$ .
- Under these assumptions, S(x) has a binomial distribution with n degrees of freedom and the probability of success is F(x):

$$S(x) \sim B\{n, p = F(x)\} \qquad E\{S(x)\} = np = nF(x)$$
  
$$f(x) = \binom{n}{x} p^{x} (1-p)^{n-x} \qquad V\{S(x)\} = np(1-p) = nF(x)(1-F(x))$$

Dr. Uwe Kruger Slide 4 Projection-Based Data Chemometrics and Deep Reconstruction Troy, November 19., 2017



• OK, this implies that the naïve estimator is unbiased:

$$E\left\{\hat{F}(x)\right\} = \frac{E\left\{S(x)\right\}}{n} = \frac{nF(x)}{n} = F(x)$$

$$V\left\{\hat{F}(x)\right\} = \frac{V\left\{S(x)\right\}}{n^2} = \frac{nF(x)(1-F(x))}{n^2} = \frac{F(x)(1-F(x))}{n}$$

$$\lim_{n \to \infty} V\left\{\hat{F}(x)\right\} \to 0$$

$$\lim_{n \to \infty} \hat{F}(x) = \lim_{n \to \infty} F(x)$$

- This follows from simple asymptotics!
- We can develop this one step further by utilizing the fact that the Binomial distribution can be approximated by a normal distribution with a reasonable degree of accuracy, meaning a large enough sample size: *np* > 5 and *n*(1 - *p*) > 5!



• Let's define a new random variable first:

$$Z(x) = \frac{S(x) - nF(x)}{\sqrt{nF(x)(1 - F(x))}} \sim N\{0,1\}$$
  

$$Z(x) = \frac{(\# x_i \le x) - nF(x)}{\sqrt{nF(x)(1 - F(x))}}$$
  

$$-1.96 \le \frac{(\# x_i \le x) - nF(x)}{\sqrt{nF(x)(1 - F(x))}} \le 1.96$$
  

$$nF(x) - 1.96\sqrt{nF(x)(1 - F(x))} \le (\# x_i \le x) \le nF(x) + 1.96\sqrt{nF(x)(1 - F(x))}$$

- The above confidence interval is computed for a significance of α=0.05!
- OK, let's move on and convert this into an integral equation, one second...



$$nF(x) - 1.96\sqrt{nF(x)(1 - F(x))} \le \int_{-\infty}^{x} \sum_{i=1}^{n} \delta(\xi - x_i) d\xi \le nF(x) + 1.96\sqrt{nF(x)(1 - F(x))}$$

$$\int_{-\infty}^{x} \delta(\xi - x_i) d\xi = \begin{cases} 1 \text{ if } x_i \le x \\ 0 \text{ if } x_i > x \end{cases}$$

$$F(x) - 1.96\sqrt{\frac{F(x)(1 - F(x))}{n}} \le \int_{-\infty}^{x} \frac{1}{n} \sum_{i=1}^{n} \delta(\xi - x_i) d\xi \le F(x) + 1.96\sqrt{\frac{F(x)(1 - F(x))}{n}}$$

$$F(x) - 1.96\sqrt{\frac{F(x)(1 - F(x))}{n}} \le \int_{-\infty}^{x} \frac{1}{n} \sum_{i=1}^{n} \frac{K(\xi - x_i)}{\text{slightlyless "spiky"}} d\xi \le F(x) + 1.96\sqrt{\frac{F(x)(1 - F(x))}{n}}$$

$$f(x) - 1.96\frac{d\left(\sqrt{\frac{F(x)(1 - F(x))}{n}}\right)}{dx} \le \frac{1}{n} \sum_{i=1}^{n} K(x - x_i) \le f(x) + 1.96\frac{d\left(\sqrt{\frac{F(x)(1 - F(x))}{n}}\right)}{dx}$$

Dr. Uwe Kruger Slide 7 Projection-Based Data Chemometrics and Deep Reconstruction Troy, November 19., 2017



• So what have we got?

$$f(x)-1.96\frac{d\left(\sqrt{\frac{F(x)(1-F(x))}{n}}\right)}{dx} \le \frac{1}{n}\sum_{i=1}^{n}K(x-x_i) \le f(x)+1.96\frac{d\left(\sqrt{\frac{F(x)(1-F(x))}{n}}\right)}{dx}$$
$$\lim_{n\to\infty}f(x)\pm 1.96\frac{d\left(\sqrt{\frac{F(x)(1-F(x))}{n}}\right)}{dx} = \lim_{n\to\infty}f(x)\pm \frac{1}{\sqrt{n}}\frac{d\left(\sqrt{F(x)(1-F(x))}\right)}{dx} \to f(x)$$
$$\lim_{n\to\infty}\frac{1}{n}\sum_{i=1}^{n}K(x-x_i) \to f(x)$$

$$\lim_{n \to \infty} \frac{1}{n} \sum_{i=1}^{n} K(x - x_i) \to f(x)$$

 All we said about the slightly less spiky Dirac delta function is that its integral must be equal to one, so how about defining it as follows:

$$K(x-x_i) = \frac{1}{\sqrt{2\pi\sigma}} e^{-\frac{1}{2}\left(\frac{x-x_i}{\sigma}\right)} \implies \lim_{\sigma \to 0} K(x-x_i) \to \delta(x-x_i)$$

Dr. Uwe Kruger Slide 8 Projection-Based Data Chemometrics and Deep Reconstruction Troy, November 19., 2017



### **Kernel Density Estimation**

• The function  $K(x - x_i)$  is referred to as a kernel function and the derivative shows that, asymptotically, the estimate:

$$\frac{1}{n}\sum_{i=1}^{n}K(x-x_{i})$$

converges to the true probability density function for any value of *x*. The above estimator is defined as a kernel density estimator.

- Along the same lines, we can also develop an approach to develop nonlinear counterpart of data-driven chemometric modeling techniques, such as principal component analysis (PCA) and partial least squares (PLS).
- Essentially, an artificial neural network can be seen as a kernel-based nonlinear modeling technique, *i.e.* the neurons are, effectively, small kernels.
- Let's start with PCA first, after some more discussions on kernels.



#### **Kernel Density Estimation**

• Theoretically, kernel functions other than the Gaussian kernel:

$$K(x-x_i) = \frac{1}{\sqrt{2\pi\sigma}} e^{-\frac{1}{2}\left(\frac{x-x_i}{\sigma}\right)^2}$$

can be considered if their area is equal to 1 and include the Epanechnikov, the triangular and the uniform kernel among others.

- Theoretically, the derivative showed that the shape of the kernel function does not influence the estimate in an asymptotic sense.
- Practically, however, the shape of the kernel function does influence the accuracy of the estimate. This yields the following general form of the kernel density estimator:

$$\frac{1}{nh}\sum_{i=1}^{n}K\left(\frac{x-x_{i}}{h}\right), \quad K\left(\frac{x-x_{i}}{h}\right) = \frac{1}{\sqrt{2\pi}}e^{-\frac{1}{2}\left(\frac{x-x_{i}}{h}\right)^{2}},$$

 $h \rightarrow$  bandwidth



- Kernel PCA is a generic nonlinear extension to linear PCA (Kruger *et al.*, 2008).
- Let's look at some basics before we go into the kernel stuff.

$$\mathbf{z} = \mathbf{A}\mathbf{s} \quad \dim \{\mathbf{z}\} > \dim \{\mathbf{s}\} \quad E\{\mathbf{z}\} = \mathbf{A}E\{\mathbf{s}\} = \mathbf{0}$$
$$\mathbf{Z} = \begin{bmatrix} \mathbf{z}_1^T \\ \mathbf{z}_2^T \\ \vdots \\ \mathbf{z}_n^T \end{bmatrix} = \begin{bmatrix} \mathbf{s}_1^T \\ \mathbf{s}_2^T \\ \vdots \\ \mathbf{s}_n^T \end{bmatrix} \mathbf{A}^T = \mathbf{U}\mathbf{L}\mathbf{P}^T \implies \text{singular value decomposition}$$

• Next, let's define the following two matrices:

 $\Sigma_{z} = \frac{1}{n} \mathbf{Z}^{T} \mathbf{Z} = \mathbf{P} \Big[ \frac{1}{n} \mathbf{L}^{2} \Big] \mathbf{P}^{T} \rightarrow \text{data covariance matrix and its eigendecom position}$  $\Phi_{z} \Big( \mathbf{Z}, \mathbf{Z} \Big) = \mathbf{Z} \mathbf{Z}^{T} = \mathbf{U} \Big[ \mathbf{L}^{2} \Big] \mathbf{U}^{T} \rightarrow \text{Gram matrix and its eigendecom position}$ 



• Let's see how we can determine the unknown source variables (up to a similarity transformation) – which are the principal components:

$$\begin{bmatrix} \mathbf{s}_1 \\ \mathbf{s}_2^T \\ \vdots \\ \mathbf{s}_n^T \end{bmatrix} \mathbf{A}^T = \mathbf{U} \mathbf{L} \mathbf{P}^T \Longrightarrow \mathbf{S} \propto \mathbf{U} \mathbf{L} = \mathbf{T}, \mathbf{A} \propto \mathbf{P}$$

 $c^{T}$ 

 $\mathbf{t} = \mathbf{P}^T \mathbf{z} = \mathbf{L}^{-1} \mathbf{U}^T \mathbf{Z} \mathbf{z} = \mathbf{L}^{-1} \mathbf{U}^T \mathbf{\Phi}_z(\mathbf{Z}, \mathbf{z})$  given that  $\mathbf{Z} = \mathbf{U} \mathbf{L} \mathbf{P}^T \Rightarrow \mathbf{P}^T = \mathbf{L}^{-1} \mathbf{U}^T \mathbf{Z}$ 

• Let's make the relationship between the source variables and the measured variables nonlinear, *i.e.*:

$$\mathbf{z} = \mathbf{\Theta}(\mathbf{s}), \quad \mathbf{f} = \mathbf{\psi}(\mathbf{z}), \quad \mathbf{t} = \mathbf{P}^{\mathrm{T}}\mathbf{f}, \quad \mathbf{F} = \begin{vmatrix} \mathbf{\psi}^{T}(\mathbf{z}_{1}) \\ \mathbf{\psi}^{T}(\mathbf{z}_{2}) \\ \vdots \\ \mathbf{\psi}^{T}(\mathbf{z}_{n}) \end{vmatrix} \mathbf{W}$$

which we assume to be bijective!

Dr. Uwe Kruger Slide 12 Projection-Based Data Chemometrics and Deep Reconstruction Troy, November 19., 2017



• Let's define the Gram matrix

$$\Phi_{z}(\mathbf{Z}, \mathbf{Z}) = \underbrace{\left[\mathbf{I} - \frac{1}{n}\mathbf{1}\mathbf{1}^{T}\right]}_{\text{incorporating mean}} \underbrace{\mathbf{FF}^{T}}_{\text{defined as the kernel matrix}} \underbrace{\left[\mathbf{I} - \frac{1}{n}\mathbf{1}\mathbf{1}^{T}\right]^{T}}_{\text{incorporating mean}}$$

$$\mathbf{FF}^{T} = \begin{bmatrix} \boldsymbol{\psi}^{T}(\mathbf{z}_{1})\boldsymbol{\psi}(\mathbf{z}_{1}) & \boldsymbol{\psi}^{T}(\mathbf{z}_{1})\boldsymbol{\psi}(\mathbf{z}_{2}) & \cdots & \boldsymbol{\psi}^{T}(\mathbf{z}_{1})\boldsymbol{\psi}(\mathbf{z}_{n}) \\ \boldsymbol{\psi}^{T}(\mathbf{z}_{2})\boldsymbol{\psi}(\mathbf{z}_{1}) & \boldsymbol{\psi}^{T}(\mathbf{z}_{2})\boldsymbol{\psi}(\mathbf{z}_{2}) & \cdots & \boldsymbol{\psi}^{T}(\mathbf{z}_{2})\boldsymbol{\psi}(\mathbf{z}_{n}) \\ \vdots & \vdots & \ddots & \vdots \\ \boldsymbol{\psi}^{T}(\mathbf{z}_{n})\boldsymbol{\psi}(\mathbf{z}_{1}) & \boldsymbol{\psi}^{T}(\mathbf{z}_{n})\boldsymbol{\psi}(\mathbf{z}_{2}) & \cdots & \boldsymbol{\psi}^{T}(\mathbf{z}_{n})\boldsymbol{\psi}(\mathbf{z}_{n}) \\ \end{bmatrix} = \mathbf{K}(\mathbf{Z}, \mathbf{Z})$$

$$\mathbf{K}(\mathbf{Z}, \mathbf{Z}) = \begin{bmatrix} k_{11} & k_{12} & \cdots & k_{1n} \\ k_{21} & k_{22} & \cdots & k_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ k_{n1} & k_{n2} & \cdots & k_{nn} \end{bmatrix}$$

Dr. Uwe Kruger
Slide 13 Projection-Based Data Chemometrics and Deep Reconstruction
Troy, November 19., 2017



• Let's repeat the "trick" we did when estimating the probability density function using the kernel density estimator using kernels:

$$k_{ij} = \mathbf{\psi}^{T}(\mathbf{z}_{i})\mathbf{\psi}(\mathbf{z}_{j}) = \frac{1}{\sqrt{2\pi\sigma}} e^{-\frac{1}{2}\left(\frac{\|\mathbf{z}_{i}-\mathbf{z}_{j}\|}{\sigma}\right)^{2}} \cdots e^{-\frac{1}{2}\left(\frac{\|\mathbf{z}_{1}-\mathbf{z}_{2}\|}{\sigma}\right)^{2}} \cdots e^{-\frac{1}{2}\left(\frac{\|\mathbf{z}_{1}-\mathbf{z}_{n}\|}{\sigma}\right)^{2}} \\ \mathbf{K}(\mathbf{Z},\mathbf{Z}) = \frac{1}{\sqrt{2\pi\sigma}} \begin{bmatrix} 1 & e^{-\frac{1}{2}\left(\frac{\|\mathbf{z}_{2}-\mathbf{z}_{1}\|}{\sigma}\right)^{2}} & \cdots & e^{-\frac{1}{2}\left(\frac{\|\mathbf{z}_{2}-\mathbf{z}_{n}\|}{\sigma}\right)^{2}} \\ e^{-\frac{1}{2}\left(\frac{\|\mathbf{z}_{n}-\mathbf{z}_{1}\|}{\sigma}\right)^{2}} & 1 & \cdots & e^{-\frac{1}{2}\left(\frac{\|\mathbf{z}_{2}-\mathbf{z}_{n}\|}{\sigma}\right)^{2}} \\ \vdots & \vdots & \ddots & \vdots \\ e^{-\frac{1}{2}\left(\frac{\|\mathbf{z}_{n}-\mathbf{z}_{1}\|}{\sigma}\right)^{2}} & e^{-\frac{1}{2}\left(\frac{\|\mathbf{z}_{n}-\mathbf{z}_{2}\|}{\sigma}\right)^{2}} & \cdots & 1 \end{bmatrix}$$

Dr. Uwe Kruger Slide 14 Projection-Based Data Chemometrics and Deep Reconstruction Troy, November 19., 2017



Let's finalize the definition of the Gram matrix:

 $\Phi_{z}(\mathbf{Z},\mathbf{Z}) = \mathbf{K}(\mathbf{Z},\mathbf{Z}) - \frac{1}{n} \mathbf{K}(\mathbf{Z},\mathbf{Z}) \mathbf{1} \mathbf{1}^{T} - \frac{1}{n} \mathbf{1} \mathbf{1}^{T} \mathbf{K}(\mathbf{Z},\mathbf{Z}) + \frac{1}{n^{2}} \mathbf{1} \mathbf{1}^{T} \mathbf{K}(\mathbf{Z},\mathbf{Z}) \mathbf{1} \mathbf{1}^{T}$ 

Next, we carry out the eigendecomposition of the Gram matrix:

$$\boldsymbol{\Phi}_{z}(\mathbf{Z},\mathbf{Z}) = \mathbf{U}\mathbf{L}\mathbf{U}^{T}$$

In a similar fashion to PCA, we can now determine the principal components:  $\mathbf{t} = \underbrace{\mathbf{L}^{-1}\mathbf{U}^{T}\left[\mathbf{I} - \frac{1}{n}\mathbf{1}\mathbf{1}^{T}\right]}_{\mathbf{T}}\mathbf{F}\left(\psi(\mathbf{z}) - \underbrace{\overline{\psi}}_{\frac{1}{n}\mathbf{F}^{T}\mathbf{1}}\right) = \underbrace{\mathbf{A}^{T}\left(\mathbf{k}(\mathbf{Z}, \mathbf{z}) - \frac{1}{n}\mathbf{K}(\mathbf{Z}, \mathbf{Z})\mathbf{1}\right)}_{\text{like a neural network, this is a weighted sum of basis functions}}$ 

$$\mathbf{t} = \mathbf{A}^T \left( \mathbf{k} (\mathbf{Z}, \mathbf{z}) - \overline{\mathbf{k}} \right)$$

Dr. Uwe Kruger Projection-Based Data Chemometrics and Deep Reconstruction Slide 15 Troy, November 19., 2017



- Asymptotically, n→∞, the shape of the basic kernel function is not important.
- Theoretically, and this follows from the properties of reproducing kernels, any function can be constructed in the feature space that maps the nonlinear surface in the data space to become a plane (subspace) in the feature space.
- The projection in the feature space then yields linear principal components in the feature space that are related to the source variables in the original variable space – connected through the following mappings:

$$\mathbf{z} = \mathbf{\Theta}(\mathbf{s}), \quad \mathbf{f} = \mathbf{\psi}(\mathbf{z}), \quad \mathbf{t} = \mathbf{A}^T (\mathbf{f} - \bar{\mathbf{f}}), \quad \mathbf{z} = \widetilde{\mathbf{\Theta}}(\mathbf{t})$$

Dr. Uwe Kruger Slide 16 Projection-Based Data Chemometrics and Deep Reconstruction Troy, November 19., 2017



 Let's examine the geometric framework that underpins the partial least squares concept: orthogonally projecting the data points onto directions for the predictor space:

 $\cos(\alpha) = \frac{x^T w}{\|x\| \|w\|}$ 

with ||w|| = 1, we get

 $\alpha \quad t \qquad \cos(\alpha) \|\mathbf{x}\| = \mathbf{x}^T \mathbf{w} = t$ and the response space:  $\cos(\beta) \|\mathbf{y}\| = \mathbf{y}^T \mathbf{v} = u$ if  $\|\mathbf{v}\| = 1$ 

Slide 17

Dr. Uwe Kruger Projection-Based Data Chemometrics and Deep Reconstruction Troy, November 19., 2017



- Let's examine the random vectors these can be of a considerable dimension – X and Y describing the predictor and response sets that are related as follows:
  - Y = BX + E E being a random vector describing uncertainty
- We could use ordinary least squares to determine the parameter matrix **B**:  $B = S_{YX}S_{XX}^{-1}$
- The problem is that if X has a very large dimension, the inverse of the covariance matrix  $S_{XX}$  may not exist or is badly conditioned!
- Here is where PLS comes in! Using the projections we discussed before:  $T = \mathbf{X}^T \mathbf{w}$  and  $U = \mathbf{Y}^T \mathbf{v}$
- Now, we select the random variables *T* and *U* such that they maximize their covariance!



- This yields the following objective function:  $I = E\{TU\} - \lambda_1(\boldsymbol{w}^T\boldsymbol{w} - 1) - \lambda_2(\boldsymbol{v}^T\boldsymbol{v} - 1)$  $I = \boldsymbol{w}^T E\{\boldsymbol{X}\boldsymbol{Y}^T\}\boldsymbol{v} - \lambda_1(\boldsymbol{w}^T\boldsymbol{w} - 1) - \lambda_2(\boldsymbol{v}^T\boldsymbol{v} - 1)$  $I = \boldsymbol{w}^T \boldsymbol{S}_{XY} \boldsymbol{v} - \lambda_1 (\boldsymbol{w}^T \boldsymbol{w} - 1) - \lambda_2 (\boldsymbol{v}^T \boldsymbol{v} - 1)$  $\frac{\partial J}{\partial w} = S_{XY} v - 2\lambda_1 w = 0$  $\frac{\partial J}{\partial w} = \boldsymbol{S}_{YX} \boldsymbol{w} - 2\lambda_2 \boldsymbol{v} = \boldsymbol{0}$  $S_{XY}S_{YX}w = 4\lambda_1\lambda_2w$  $S_{YX}S_{XY}v = 4\lambda_1\lambda_2v$
- So, we now have the direction vectors **w** and **v** in both spaces!
- That also means that we have the random variables T and U!
- Whilst X predicts Y, PLS utilizes T instead of X to predict Y!



• With  $T = \mathbf{X}^T \mathbf{w}$ , we get:

F = X - Tp and E = Y - Tq – these being the residual vectors for **X** and **Y**, respectively.

 The parameter vectors *p* and *q* can be obtained by solving two least squares regression problems – minimizing the length of the residual vectors:

$$p = \frac{E\{XT\}}{E\{T^2\}}$$
 and  $q = \frac{E\{YT\}}{E\{T^2\}}$ 

- After that, the PLS algorithm can be repeated using the residual vectors *F* and *E* instead of the original random vectors *X* and *Y*.
- This gives rise to the following iterative algorithm, which is referred to as the standard PLS algorithm and detailed on the next slide. This algorithm was first published by Herman Wold in 1966.



 Let X and Y be two data matrices that store a total of n data points drawn from the random vectors X and Y, respectively:

$$\mathbf{X} = \begin{bmatrix} x_{11} & x_{12} & \cdots & x_{1N} \\ x_{21} & x_{22} & \cdots & x_{2N} \\ \vdots & \vdots & \ddots & \vdots \\ x_{n1} & x_{n2} & \cdots & x_{nN} \end{bmatrix} \quad \mathbf{Y} = \begin{bmatrix} y_{11} & y_{12} & \cdots & y_{1N} \\ y_{21} & y_{22} & \cdots & y_{2N} \\ \vdots & \vdots & \ddots & \vdots \\ y_{n1} & y_{n2} & \cdots & y_{nN} \end{bmatrix}$$

• The first step is to normalize the data matrices, *i.e.* the observations in each column are mean centered and scaled to have a unit variance: Sample mean vectors  $: \overline{\mathbf{x}} = \frac{1}{n} \mathbf{X}^T \mathbf{1}$   $: \overline{\mathbf{y}} = \frac{1}{n} \mathbf{Y}^T \mathbf{1}$ Sample variance vectors  $: \mathbf{\sigma}_X^2 = \frac{1}{n-1} \operatorname{diag} \left( [\mathbf{X} - \mathbf{1} \overline{\mathbf{x}}^T]^T [\mathbf{X} - \mathbf{1} \overline{\mathbf{x}}^T] \right)$  $: \mathbf{\sigma}_Y^2 = \frac{1}{n-1} \operatorname{diag} \left( [\mathbf{Y} - \mathbf{1} \overline{\mathbf{x}}^T]^T [\mathbf{Y} - \mathbf{1} \overline{\mathbf{x}}^T] \right)$ 

Normalizin g both matrices :  $\mathbf{X}_0 = \left[\mathbf{X} - \mathbf{1}\overline{\mathbf{x}}^T\right] \left[\operatorname{diag}\left(\mathbf{\sigma}_X^2\right)\right]^{\frac{1}{2}} \quad \mathbf{Y}_0 = \left[\mathbf{Y} - \mathbf{1}\overline{\mathbf{y}}^T\right] \left[\operatorname{diag}\left(\mathbf{\sigma}_Y^2\right)\right]^{\frac{1}{2}}$ 

Dr. Uwe Kruger Slide 21 Projection-Based Data Chemometrics and Deep Reconstruction Troy, November 19., 2017



Next, define the covariance and cross-covariance matrices: Sample covariance matrix  $: \mathbf{S}_{XX}^{(0)} = \frac{1}{n-1} \mathbf{X}_{0}^{T} \mathbf{X}_{0}, \quad \mathbf{M}_{XX}^{(0)} = (n-1) \mathbf{S}_{XX}^{(0)}$ Sample cross-covariance matrix:  $\mathbf{S}_{xy}^{(0)} = \frac{1}{n-1} \mathbf{X}_{0}^{T} \mathbf{Y}_{0}, \quad \mathbf{M}_{xy}^{(0)} = (n-1) \mathbf{S}_{xy}^{(0)}$ and determining the regression vectors  $\mathbf{w}_i = \mathbf{w}_0$ ; Setup the PLS iteration for i = 1: m $\mathbf{p}_{i} = \mathbf{M}_{XX}^{(i-1)} \mathbf{w}_{i} / (\mathbf{w}_{i}^{T} \mathbf{M}_{XX}^{(i-1)} \mathbf{w}_{i});$  $\varepsilon$  =100:  $\mathbf{q}_{i} = \mathbf{M}_{YX}^{(i-1)} \mathbf{w} / \left( \mathbf{w}_{i}^{T} \mathbf{M}_{YY}^{(i-1)} \mathbf{w}_{i} \right);$  $\mathbf{w}_{0} = \mathbf{M}_{XX}^{(i-1)}(:,1) / \operatorname{norm} \left( \mathbf{M}_{XY}^{(i-1)}(:,1) \right);$  $\mathbf{X}_i = \mathbf{X}_{i-1} - \mathbf{X}_{i-1} \mathbf{w}_i \mathbf{p}_i^T;$ while  $\varepsilon < 1e-10$  $\mathbf{Y}_{i} = \mathbf{Y}_{i-1} - \mathbf{X}^{(i-1)} \mathbf{w}_{i} \mathbf{q}_{i}^{T};$  $\mathbf{v} = \mathbf{M}_{yy}^{(i-1)} \mathbf{w}_{0};$  $\mathbf{v} = \mathbf{v}/\operatorname{norm}(\mathbf{v});$  $\mathbf{M}_{XX}^{(i)} = \mathbf{X}_i^T \mathbf{X}_i;$  $\mathbf{w} = \mathbf{M}_{\mathbf{v}\mathbf{v}}^{(i-1)}\mathbf{v};$  $\mathbf{M}_{\mathbf{Y}\mathbf{Y}}^{(i)} = \mathbf{X}_{i}^{T}\mathbf{Y}_{i};$  $\mathbf{w}_{u} = \mathbf{w}/\text{norm}(\mathbf{w});$  $P(:, i) = p_i;$  $\varepsilon = \operatorname{norm}(\mathbf{w}_{1} - \mathbf{w}_{0});$  $\mathbf{Q}(:,i) = \mathbf{q}_i;$  $\mathbf{W}(:, i) = \mathbf{w}_i;$  $\mathbf{W}_0 = \mathbf{W}_u;$ end end

Dr. Uwe Kruger Slide 22 Projection-Based Data Chemometrics and Deep Reconstruction Troy, November 19., 2017



- The regression matrix can now be estimated as follows:  $\mathbf{B} = \mathbf{Q} [\mathbf{W}^T \mathbf{P}]^{-1} \mathbf{W}^T$
- To establish a nonlinear extension of the standard PLS algorithm, let's look at the standard algorithm again:

$$\lambda \mathbf{v} = \mathbf{Y}_0^T \underbrace{\mathbf{X}_0 \mathbf{X}_0^T}_{\text{This is a Gram matrix}} \mathbf{Y}_0 \mathbf{v}$$

 We can "kernelize" the above Gram matrix by using the following nonlinear transformation involving the random vectors X, Y and E:

$$\boldsymbol{G} = \boldsymbol{\psi}(\boldsymbol{X}) \quad \boldsymbol{Y} = \mathbf{B}\,\boldsymbol{G} + \boldsymbol{E}$$

• Based on the data matrix **X**, we get:

$$\mathbf{G} = \boldsymbol{\Psi} (\mathbf{X}_0) \quad \mathbf{Y}_0 = \mathbf{G}_0 \mathbf{B}^T + \mathbf{E}$$

Dr. Uwe Kruger Slide 23 Projection-Based Data Chemometrics and Deep Reconstruction Troy, November 19., 2017



Using the Gram matrix based on \u03c6(X\_0), we can compute the projection vector v as follows:

$$\lambda \mathbf{v} = \mathbf{Y}_0^T \quad \left[ \mathbf{I} - \frac{1}{n} \mathbf{1} \mathbf{1}^T \right] \mathbf{G} \mathbf{G}^T \left[ \mathbf{I} - \frac{1}{n} \mathbf{1} \mathbf{1}^T \right] \quad \mathbf{Y}_0 \mathbf{v}$$

This is  $\Phi_X(\mathbf{X}_0, \mathbf{X}_0)$ , the Gram matrix for  $\psi(\mathbf{X}_0)$ 

- Once we have **v**, it is easy to compute the vector **u**:  $\mathbf{u} = \mathbf{Y}_0 \mathbf{v}$
- The next step is to compute the vector **t**. For linear PLS, we can derive the following relationship:

$$\mathbf{t} = \mathbf{X}_0 \mathbf{w}; \quad \mathbf{w} \propto \mathbf{M}_{XY} \mathbf{v} = \mathbf{X}_0^T \underbrace{\mathbf{Y}_0 \mathbf{v}}_{\mathbf{u}}$$

$$\mathbf{t} \propto \mathbf{X}_0 \mathbf{X}_0^T \mathbf{u}$$



- Instead of the linear Gram matrix  $\mathbf{X}_0 \mathbf{X}_0^T$ , we can also use the nonlinear Gram matrix  $\mathbf{\Phi}_X(\mathbf{X}_0, \mathbf{X}_0)$ , which gives rise to:  $\mathbf{t} \propto \mathbf{\Phi}(\mathbf{X}_0, \mathbf{X}_0)\mathbf{u}$
- To address the scaling problem, as we cannot compute the projection vector w, we can scale the vector t to unit length:  $\mathbf{t} = \mathbf{t}/\text{norm}(\mathbf{t})$
- Now, we can deflate the Gram matrix,  $\mathbf{G}_0 \mathbf{G}_0^T = \mathbf{\Phi}_X(\mathbf{X}_0, \mathbf{X}_0)$ :

$$\mathbf{G}_{i} = \mathbf{G}_{i-1} - \mathbf{t}_{i-1}\mathbf{p}_{i-1}^{T} = \mathbf{G}_{i-1} - \mathbf{t}_{i-1}\frac{\mathbf{t}_{i-1}^{T}\mathbf{G}_{i-1}}{\underbrace{\mathbf{t}_{i-1}^{T}\mathbf{t}_{i-1}}_{1}} = \left[\mathbf{I} - \mathbf{t}_{i-1}\mathbf{t}_{i-1}^{T}\right]\mathbf{G}_{i-1}$$
$$\mathbf{G}_{i}\mathbf{G}_{i}^{T} = \left[\mathbf{I} - \mathbf{t}_{i-1}\mathbf{t}_{i-1}^{T}\right]\mathbf{G}_{i-1}\mathbf{G}_{i-1}^{T}\left[\mathbf{I} - \mathbf{t}_{i-1}\mathbf{t}_{i-1}^{T}\right]$$

Slide 25

Dr. Uwe Kruger Projection-Based Data Chemometrics and Deep Reconstruction Troy, November 19., 2017



and the response matrix  $Y_0$ :

$$\mathbf{Y}_{i} = \mathbf{Y}_{i-1} - \mathbf{t}_{i-1}\mathbf{q}_{i-1}^{T} = \mathbf{Y}_{i-1} - \mathbf{t}_{i-1}\frac{\mathbf{t}_{i-1}^{T}\mathbf{Y}_{i-1}}{\underbrace{\mathbf{t}_{i-1}^{T}\mathbf{t}_{i-1}}} = \left[\mathbf{I} - \mathbf{t}_{i-1}\mathbf{t}_{i-1}^{T}\right]\mathbf{Y}_{i-1}$$

- The last step is to compute the regression matrix. Again, let's look again at the linear PLS algorithm first:  $\mathbf{B}^{T} = \mathbf{W} \begin{bmatrix} \mathbf{P}^{T} \mathbf{W} \end{bmatrix}^{-1} \mathbf{Q}^{T} \quad \mathbf{Q} = \mathbf{Y}_{0}^{T} \mathbf{T} \begin{bmatrix} \mathbf{T}^{T} \mathbf{T} \end{bmatrix}^{-1} \quad \mathbf{P} = \mathbf{X}_{0}^{T} \mathbf{T} \begin{bmatrix} \mathbf{T}^{T} \mathbf{T} \end{bmatrix}^{-1} \quad \mathbf{W} \propto \mathbf{X}_{0}^{T} \mathbf{U}$  $\mathbf{B}^{T} = \mathbf{X}_{0}^{T} \mathbf{U} \begin{bmatrix} \mathbf{T}^{T} \mathbf{T} \end{bmatrix}^{-1} \mathbf{T}^{T} \mathbf{X}_{0} \mathbf{X}_{0}^{T} \mathbf{U} \end{bmatrix}^{-1} \begin{bmatrix} \mathbf{T}^{T} \mathbf{T} \end{bmatrix}^{-1} \mathbf{T}^{T} \mathbf{Y}_{0}$  $\mathbf{B}^{T} = \mathbf{X}_{0}^{T} \mathbf{U} \begin{bmatrix} \mathbf{T}^{T} \mathbf{T} \end{bmatrix}^{-1} \mathbf{T}^{T} \mathbf{X}_{0} \mathbf{X}_{0}^{T} \mathbf{U} \end{bmatrix}^{-1} \mathbf{T}^{T} \mathbf{Y}_{0}$
- Using this regression matrix for predicting a new observation yields:  $\hat{\mathbf{y}}_{0}^{T} = \mathbf{x}_{0}^{T}\mathbf{B}^{T} = \mathbf{x}_{0}^{T}\mathbf{X}_{0}^{T}\mathbf{U}[\mathbf{T}^{T}\mathbf{X}_{0}\mathbf{X}_{0}^{T}\mathbf{U}]^{-1}\mathbf{T}^{T}\mathbf{Y}_{0}$

Slide 26

Dr. Uwe Kruger Projection-Based Data Chemometrics and Deep Reconstruction Troy, November 19., 2017



- Finally, replacing the linear Gram matrix and vector by there nonlinear counterparts gives rise to:  $\hat{\mathbf{y}} = \mathbf{B}\overline{\mathbf{\psi}}(\mathbf{x}_0) = \mathbf{Y}^T \mathbf{T} \left[ \mathbf{U}^T \overline{\mathbf{\psi}}(\mathbf{X}_0) \overline{\mathbf{\psi}}^T (\mathbf{X}_0) \mathbf{T} \right]^{-1} \mathbf{U}^T \overline{\mathbf{\psi}}(\mathbf{X}_0) \overline{\mathbf{\psi}}(\mathbf{x}_0)$
- Let's put this all together and define the KPLS algorithm.
- Besides the construction of the nonlinear transformation, and with it its Gram matrix, the rest of the algorithm is related to the linear PLS algorithm.
- Compared to artificial neural networks, which have many network weights, the "only" parameter that needs to be specified is the kernel parameter. The remaining parameter are obtained by a linear regression problem an solved using the robust PLS algorithm!



# How could chemometric techniques assist in dealing with very large network architechtures?

- Disclaimer: KPLS is not a substitute to deep learning architectures!
- KPLS does run into problems if the number of data points increase, say beyond 10,000 (remember the size of the Gram matrix is equal to the number of data points squared)
- A KPLS model has the potential to outperforms competitive artificial neural network models when the number of variables x or y are larger and/or the number of data points is small.
- To see how KPCA and KPLS could be useful tools, let's examine the structure of large network topologies on the next slide in more detail



# How could chemometric techniques assist in dealing with very large network architechtures?

• Starting from a "small" (trained) network:



- how do we know that this neuron is important or could be discarded if it contributes negligibly to the accuracy of the network prediction – e.g. for specific tasks (set of lung images)?
- 2. secondly, how can we statistically examine significant differences in the individual layers/layer combinations if we have two sets of images (one set that is labeled normal, whilst the other set is labeled as containing anomalies)?



#### Abnormality Detection (Basis Multivariate Approach)

• Hotelling's T<sup>2</sup> Statistic

•Q Statistic



Slide 30

Dr. Uwe Kruger Projection-Based Data Chemometrics and Deep Reconstruction Troy, November 19., 2017



- Analysis of data from a Volkswagen 1.9L TDI diesel engine.
- Various fault conditions were recorded and diagnosed.



**Slide 31 Projection-Based Data Chemometrics and Deep Reconstruction** Troy, November 19., 2017



#### Variables analysed

No	Engine Variable	Unit	Note
1	Fuel Flow	kg/h	
2	Air Flow	kg/h	
3	Inlet Manifold Pressure	Bar	
4	Inlet Manifold Temperature	°C	output
5	Turbine Inlet Pressure	Bar	
6	Turbine Inlet Temperature	°C	

Principal Component	Variance Captured (%)	Variance Total (%)
1	79.5998	79.5998
2	16.4492	96.0490
3	2.4169	98.4659
4	1.0745	99.5404
5	0.4010	99.9414
6	0.0586	100.000

Modelling results

RPM	1500	2500	3500	4500
	30%	49%	57%	62%
	40%	59%	64%	65%
Pedal Position	54%	74%	74%	76%
	62%	78%	80%	83%
	100%	100%	100%	100%

Number of Bottleneck Nodes	Variance Captured (%)	Note
1	97.8160	Important variation
2	99.4212	
3	99.8336	
4	99.8725	Negligible
5	99.9401	
6	99.9414	

#### Slide 32

Dr. Uwe Kruger Projection-Based Data Chemometrics and Deep Reconstruction Troy, November 19., 2017



#### Air leak of 2mm in the manifold plenum chamber



**Slide 33 Projection-Based Data Chemometrics and Deep Reconstruction** Troy, November 19., 2017



- An incipient hole in the air intake system could be successfully detected.
- However, a detailed diagnosis as to which recorded engine variable is affected by this event could not be obtained.
- Traditional techniques fail to detect or diagnose this event.
- Model-based fault detection and diagnosis is expensive, whilst datadriven techniques are a viable alternative that are cost-effective.





#### Air leak of 6mm in the manifold plenum chamber







Dr. Uwe Kruger Projection-Based Data Chemometrics and Deep Reconstruction Troy, November 19., 2017



#### Application to Chemistry (RAMAN Spectroscopy) Variable Selection



## Application to Chemistry (RAMAN Spectroscopy)



Slide 38

Dr. Uwe Kruger Projection-Based Data Chemometrics and Deep Reconstruction Troy, November 19., 2017

