

# Machine Intelligence II

## Kernel PCA

### Tutorial

Neural Information Processing Group (Prof. Dr. Klaus Obermayer)

26.05.2016

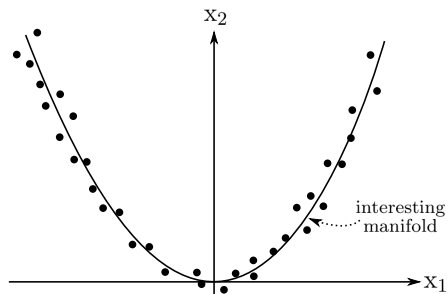
# Projection Methods

- Principal Component Analysis (PCA)
- Online-PCA  $\Rightarrow$  Hebbian Learning
- Nonlinear Structure  $\Rightarrow$  Kernel PCA
- Source Separation
  - Model based  $\Rightarrow$  Independent Component Analysis (ICA)
  - Cost Function Based  $\Rightarrow$  Projection Pursuit

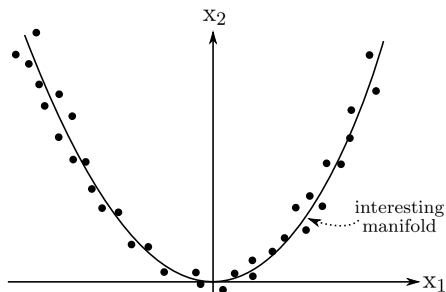
# Kernel Principal Component Analysis

**Goal:** finding nonlinear regularities

# Kernel Principal Component Analysis: Motivation



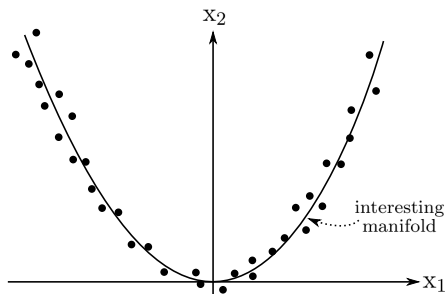
# Kernel Principal Component Analysis: Motivation



in the original space:

- standard PCA: two directions with high variance
- "interesting" feature is a non-linear combination of elementary features

# Kernel Principal Component Analysis: Motivation

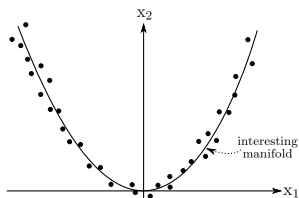


in the original space:

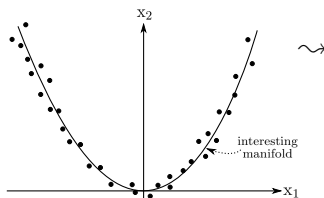
- standard PCA: two directions with high variance
- "interesting" feature is a non-linear combination of elementary features

⇒ problem for standard PCA

# Kernel Principal Component Analysis: Intuition

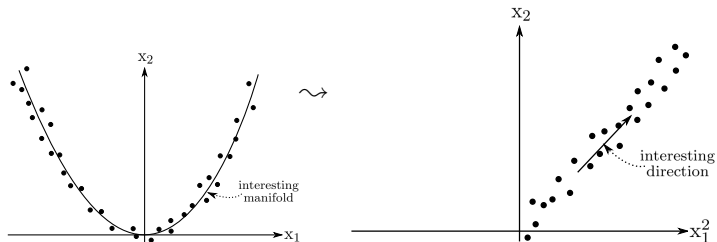


# Kernel Principal Component Analysis: Intuition

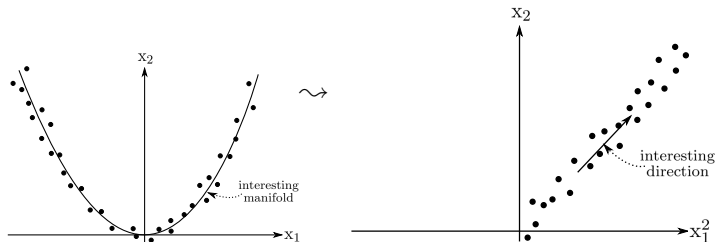




# Kernel Principal Component Analysis: Intuition

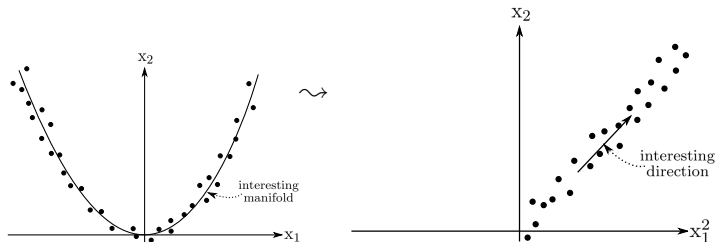


# Kernel Principal Component Analysis: Intuition



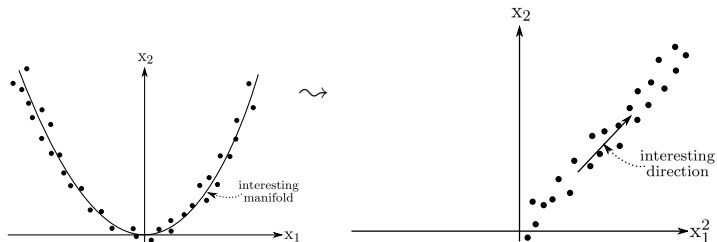
original (data) space       $\leadsto$       transformed (feature) space:

# Kernel Principal Component Analysis: Intuition



original (data) space  $\leadsto$  transformed (feature) space:  
 $\leadsto$  "interesting" features correspond to directions of high variance

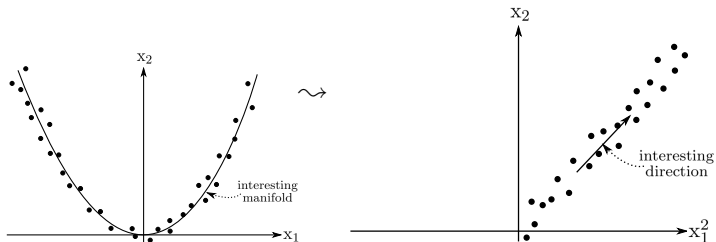
# Kernel Principal Component Analysis: Intuition



original (data) space  $\leadsto$  transformed (feature) space:  
 $\leadsto$  "interesting" features correspond to directions of high variance

## Agenda

# Kernel Principal Component Analysis: Intuition

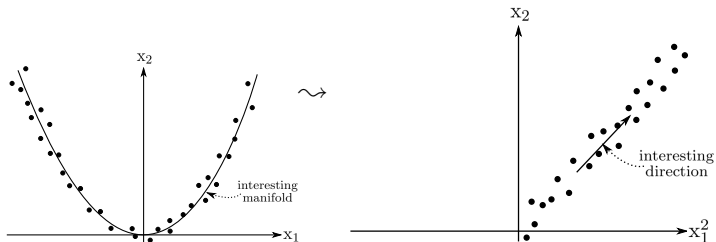


original (data) space  $\leadsto$  transformed (feature) space:  
 $\leadsto$  "interesting" features correspond to directions of high variance

## Agenda

- 1 non-linear preprocessing:

# Kernel Principal Component Analysis: Intuition

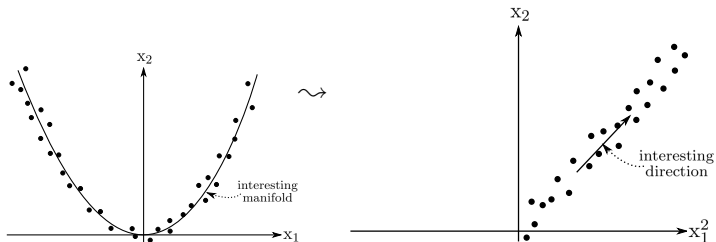


original (data) space  $\sim$  transformed (feature) space:  
 $\sim$  "interesting" features correspond to directions of high variance

## Agenda

- 1 non-linear preprocessing:  
 transformation into an "appropriate" feature space  $\underline{\phi} : \underline{x} \rightarrow \underline{\phi}(\underline{x})$

# Kernel Principal Component Analysis: Intuition



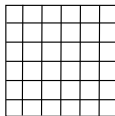
original (data) space  $\sim$  transformed (feature) space:  
 $\sim$  "interesting" features correspond to directions of high variance

## Agenda

- ① non-linear preprocessing:  
 transformation into an "appropriate" feature space  $\underline{\phi} : \underline{x} \rightarrow \underline{\phi}(\underline{x})$
- ② apply standard linear methods

# Why we need the Kernel Trick

Problem: Some feature spaces may be extremely high-dimensional

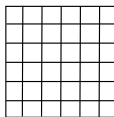


pixel image



# Why we need the Kernel Trick

Problem: Some feature spaces may be extremely high-dimensional

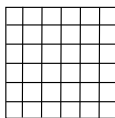


pixel image

- interesting structure is hidden in correlations between pixel values

# Why we need the Kernel Trick

Problem: Some feature spaces may be extremely high-dimensional



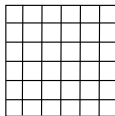
pixel image

- interesting structure is hidden in correlations between pixel values
- suitable feature space: space spanned by all  $d^{\text{th}}$ -order monomials

e.g.  $d = 2$  :  $x_1^2, x_1x_2, x_2^2, x_1x_3, x_2x_3, x_3^2, \dots$

# Why we need the Kernel Trick

Problem: Some feature spaces may be extremely high-dimensional



pixel image

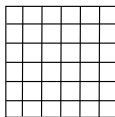
- interesting structure is hidden in correlations between pixel values
- suitable feature space: space spanned by all  $d^{\text{th}}$ -order monomials

e.g.  $d = 2$  :  $x_1^2, x_1x_2, x_2^2, x_1x_3, x_2x_3, x_3^2, \dots$

- dimensionality  $O(N^d)$  prohibits “direct” application of this idea

# Why we need the Kernel Trick

Problem: Some feature spaces may be extremely high-dimensional



pixel image

- interesting structure is hidden in correlations between pixel values
- suitable feature space: space spanned by all  $d^{\text{th}}$ -order monomials

e.g.  $d = 2$  :  $x_1^2, x_1x_2, x_2^2, x_1x_3, x_2x_3, x_3^2, \dots$

- dimensionality  $O(N^d)$  prohibits “direct” application of this idea

~> **kernel trick** allows to avoid this problem

## Insight: PCA can be reformulated to use Scalar Products

PCA involves solution of the eigenvalue problem of the covariance matrix

$$\underline{\mathbf{C}}\underline{\mathbf{e}}_k = \lambda_k \underline{\mathbf{e}}_k$$

# Insight: PCA can be reformulated to use Scalar Products

PCA involves solution of the eigenvalue problem of the covariance matrix

$$\underline{\mathbf{C}}\underline{\mathbf{e}}_k = \lambda_k \underline{\mathbf{e}}_k \quad \leadsto \quad \underline{\mathbf{K}}\underline{\mathbf{a}}_k = p\lambda_k \underline{\mathbf{a}}_k$$

## Insight: PCA can be reformulated to use Scalar Products

PCA involves solution of the eigenvalue problem of the covariance matrix

$$\underline{\mathbf{C}}\underline{\mathbf{e}}_k = \lambda_k \underline{\mathbf{e}}_k \quad \leadsto \quad \underline{\mathbf{K}}\underline{\mathbf{a}}_k = p\lambda_k \underline{\mathbf{a}}_k$$

where  $\underline{\mathbf{a}}_k^{(\beta)}$  describe an expansion of the eigenvectors in terms of data points

# Insight: PCA can be reformulated to use Scalar Products

PCA involves solution of the eigenvalue problem of the covariance matrix

$$\underline{\mathbf{C}}\underline{\mathbf{e}}_k = \lambda_k \underline{\mathbf{e}}_k \quad \leadsto \quad \underline{\mathbf{K}}\underline{\mathbf{a}}_k = p\lambda_k \underline{\mathbf{a}}_k$$

where  $a_k^{(\beta)}$  describe an expansion of the eigenvectors in terms of data points

$$\underline{\mathbf{e}}_k = \sum_{\beta=1}^p a_k^{(\beta)} \underline{\mathbf{x}}^{(\beta)}$$



## Insight: PCA can be reformulated to use Scalar Products

PCA involves solution of the eigenvalue problem of the covariance matrix

$$\underline{\mathbf{C}}\underline{\mathbf{e}}_k = \lambda_k \underline{\mathbf{e}}_k \quad \leadsto \quad \underline{\mathbf{K}}\underline{\mathbf{a}}_k = p\lambda_k \underline{\mathbf{a}}_k$$

where  $a_k^{(\beta)}$  describe an expansion of the eigenvectors in terms of data points

$$\underline{\mathbf{e}}_k = \sum_{\beta=1}^p a_k^{(\beta)} \underline{\mathbf{x}}^{(\beta)}$$

always possible: PCs "live" in the subspace spanned by the (centered) data

# Insight: PCA can be reformulated to use Scalar Products

PCA involves solution of the eigenvalue problem of the covariance matrix

$$\underline{\mathbf{C}}\underline{\mathbf{e}}_k = \lambda_k \underline{\mathbf{e}}_k \quad \leadsto \quad \underline{\mathbf{K}}\underline{\mathbf{a}}_k = p\lambda_k \underline{\mathbf{a}}_k$$

where  $a_k^{(\beta)}$  describe an expansion of the eigenvectors in terms of data points

$$\underline{\mathbf{e}}_k = \sum_{\beta=1}^p a_k^{(\beta)} \underline{\mathbf{x}}^{(\beta)}$$

always possible: PCs "live" in the subspace spanned by the (centered) data

$\Rightarrow$  We can write EV-problem s.t. data-points occur only in scalar products!

## Insight: PCA can be reformulated to use Scalar Products

PCA involves solution of the eigenvalue problem of the covariance matrix

$$\underline{\mathbf{C}}\underline{\mathbf{e}}_k = \lambda_k \underline{\mathbf{e}}_k \quad \leadsto \quad \underline{\mathbf{K}}\underline{\mathbf{a}}_k = p\lambda_k \underline{\mathbf{a}}_k$$

where  $a_k^{(\beta)}$  describe an expansion of the eigenvectors in terms of data points

$$\underline{\mathbf{e}}_k = \sum_{\beta=1}^p a_k^{(\beta)} \underline{\mathbf{x}}^{(\beta)}$$

always possible: PCs "live" in the subspace spanned by the (centered) data

$\Rightarrow$  We can write EV-problem s.t. data-points occur only in scalar products!

**note:** Eigenvectors  $\underline{\mathbf{e}}_k \in \mathbb{R}^N$  but coefficients  $\underline{\mathbf{a}}_k \in \mathbb{R}^p$

# Why Scalar Products?

- If PCA doesn't work nonlinear pre-processing can simplify the problem

# Why Scalar Products?

- If PCA doesn't work nonlinear pre-processing can simplify the problem

$$\underline{\phi} : \underline{\mathbf{x}} \xrightarrow{\text{non-linear transformation}} \underline{\phi(\mathbf{x})}$$

# Why Scalar Products?

- If PCA doesn't work nonlinear pre-processing can simplify the problem

$$\underline{\phi} : \underline{\mathbf{x}} \xrightarrow{\text{non-linear transformation}} \underline{\phi}(\underline{\mathbf{x}})$$

## Mercer's theorem

Every **positive definite** kernel  $k$  corresponds to a scalar product in some metric feature space (*cf. MI I*).

# Why Scalar Products?

- If PCA doesn't work nonlinear pre-processing can simplify the problem

$$\underline{\phi} : \underline{\mathbf{x}} \xrightarrow{\text{non-linear transformation}} \underline{\phi}(\underline{\mathbf{x}})$$

## Mercer's theorem

Every **positive definite** kernel  $k$  corresponds to a scalar product in some metric feature space (*cf. MI I*).

## Kernel trick

$\Rightarrow$  avoid direct transformation  $\underline{\phi}$

# Why Scalar Products?

- If PCA doesn't work nonlinear pre-processing can simplify the problem

$$\underline{\phi} : \underline{\mathbf{x}} \xrightarrow{\text{non-linear transformation}} \underline{\phi}(\underline{\mathbf{x}})$$

## Mercer's theorem

Every **positive definite** kernel  $k$  corresponds to a scalar product in some metric feature space (*cf. MI I*).

## Kernel trick

- $\Rightarrow$  avoid direct transformation  $\underline{\phi}$
- $\Rightarrow$  replace all scalar products by "kernel functions"

$$\underline{\phi}_{(\underline{\mathbf{x}})}^T \underline{\phi}_{(\underline{\mathbf{x}'})} \longleftrightarrow k(\underline{\mathbf{x}}, \underline{\mathbf{x}'})$$



# Why Scalar Products?

- If PCA doesn't work nonlinear pre-processing can simplify the problem

$$\underline{\phi} : \underline{\mathbf{x}} \xrightarrow{\text{non-linear transformation}} \underline{\phi}(\underline{\mathbf{x}})$$

## Mercer's theorem

Every **positive definite** kernel  $k$  corresponds to a scalar product in some metric feature space (*cf. MI I*).

## Kernel trick

- $\Rightarrow$  avoid direct transformation  $\underline{\phi}$
- $\Rightarrow$  replace all scalar products by "kernel functions"

$$\underline{\phi}^T(\underline{\mathbf{x}})\underline{\phi}(\underline{\mathbf{x}}') \longleftrightarrow k(\underline{\mathbf{x}}, \underline{\mathbf{x}}')$$

---

if an algorithm can be formulated solely in terms of scalar products, a non-linear version can be derived "without actually projecting" into the (highdimensional) feature space

# Kernels and Applications

## Popular Kernel Functions

$$k(\underline{\mathbf{x}}, \underline{\mathbf{x}}') = (\underline{\mathbf{x}}^T \underline{\mathbf{x}}' + 1)^d$$

polynomial kernel of degree  $d$   
*image processing (pixel correlation)*

$$k(\underline{\mathbf{x}}, \underline{\mathbf{x}}') = \exp \left\{ -\frac{(\underline{\mathbf{x}} - \underline{\mathbf{x}}')^2}{2\sigma^2} \right\}$$

RBF-kernel with range  $\sigma$   
*infinite dimensional feature space*

$$k(\underline{\mathbf{x}}, \underline{\mathbf{x}}') = \tanh \{ K \underline{\mathbf{x}}^T \underline{\mathbf{x}}' + \theta \}$$

neural network kernel with parameters  $K$  and  $\theta$   
*not necessarily positive definite*

## Other Kernelizable Methods

- Support Vector Machines (MI I)
- Fisher Discriminant Analysis, Canonical Correlation Analysis
- K-Means Clustering, Self-Organizing Maps (both MI II)

- 1 calculate the un-normalized kernel matrix  $\tilde{\mathbf{K}} \in \mathbb{R}^{p,p}$

$$\tilde{K}_{\alpha\beta} = k(\underline{\mathbf{x}}^{(\alpha)}, \underline{\mathbf{x}}^{(\beta)})$$

- ① calculate the un-normalized kernel matrix  $\tilde{\mathbf{K}} \in \mathbb{R}^{p,p}$

$$\tilde{K}_{\alpha\beta} = k(\underline{\mathbf{x}}^{(\alpha)}, \underline{\mathbf{x}}^{(\beta)})$$

- ② center the kernel matrix ( $\implies$  centered data in feature space)

$$K_{\alpha\beta} = \tilde{K}_{\alpha\beta} - \frac{1}{p} \sum_{\delta=1}^p \tilde{K}_{\alpha\delta} - \frac{1}{p} \sum_{\gamma=1}^p \tilde{K}_{\gamma\beta} + \frac{1}{p^2} \sum_{\gamma,\delta=1}^p \tilde{K}_{\gamma\delta}$$

- ① calculate the un-normalized kernel matrix  $\tilde{\mathbf{K}} \in \mathbb{R}^{p,p}$

$$\tilde{K}_{\alpha\beta} = k(\underline{\mathbf{x}}^{(\alpha)}, \underline{\mathbf{x}}^{(\beta)})$$

- ② center the kernel matrix ( $\implies$  centered data in feature space)

$$K_{\alpha\beta} = \tilde{K}_{\alpha\beta} - \frac{1}{p} \sum_{\delta=1}^p \tilde{K}_{\alpha\delta} - \frac{1}{p} \sum_{\gamma=1}^p \tilde{K}_{\gamma\beta} + \frac{1}{p^2} \sum_{\gamma,\delta=1}^p \tilde{K}_{\gamma\delta}$$

- ③ solve the eigenvalue problem  $\boxed{\frac{1}{p} \mathbf{K} \tilde{\mathbf{a}}_k = \lambda_k \tilde{\mathbf{a}}_k}$

- ① calculate the un-normalized kernel matrix  $\tilde{\mathbf{K}} \in \mathbb{R}^{p,p}$

$$\tilde{K}_{\alpha\beta} = k(\underline{\mathbf{x}}^{(\alpha)}, \underline{\mathbf{x}}^{(\beta)})$$

- ② center the kernel matrix ( $\implies$  centered data in feature space)

$$K_{\alpha\beta} = \tilde{K}_{\alpha\beta} - \frac{1}{p} \sum_{\delta=1}^p \tilde{K}_{\alpha\delta} - \frac{1}{p} \sum_{\gamma=1}^p \tilde{K}_{\gamma\beta} + \frac{1}{p^2} \sum_{\gamma,\delta=1}^p \tilde{K}_{\gamma\delta}$$

- ③ solve the eigenvalue problem  $\boxed{\frac{1}{p} \mathbf{K} \tilde{\mathbf{a}}_k = \lambda_k \tilde{\mathbf{a}}_k}$

- ④ normalize eigenvectors to unit length ( $\implies \underline{\mathbf{e}}_k^2 = 1$  in feature space)

$$\underline{\mathbf{a}}_k = \frac{1}{\sqrt{p\lambda_k} \|\tilde{\mathbf{a}}_k\|} \tilde{\mathbf{a}}_k$$

- ① calculate the un-normalized kernel matrix  $\tilde{\mathbf{K}} \in \mathbb{R}^{p,p}$

$$\tilde{K}_{\alpha\beta} = k(\underline{\mathbf{x}}^{(\alpha)}, \underline{\mathbf{x}}^{(\beta)})$$

- ② center the kernel matrix ( $\implies$  centered data in feature space)

$$K_{\alpha\beta} = \tilde{K}_{\alpha\beta} - \frac{1}{p} \sum_{\delta=1}^p \tilde{K}_{\alpha\delta} - \frac{1}{p} \sum_{\gamma=1}^p \tilde{K}_{\gamma\beta} + \frac{1}{p^2} \sum_{\gamma,\delta=1}^p \tilde{K}_{\gamma\delta}$$

- ③ solve the eigenvalue problem  $\boxed{\frac{1}{p} \mathbf{K} \tilde{\mathbf{a}}_k = \lambda_k \tilde{\mathbf{a}}_k}$

- ④ normalize eigenvectors to unit length ( $\implies \underline{\mathbf{e}}_k^2 = 1$  in feature space)

$$\underline{\mathbf{a}}_k = \frac{1}{\sqrt{p\lambda_k} \|\tilde{\mathbf{a}}_k\|} \tilde{\mathbf{a}}_k$$

- ⑤ calculate projections of data points  $\underline{\mathbf{x}}^{(\alpha)}$  onto eigenvectors  $\underline{\mathbf{e}}_k$

$$u_k(\underline{\mathbf{x}}^{(\alpha)}) = \sum_{\beta=1}^p a_k^{(\beta)} K_{\beta\alpha} \quad \leftarrow \text{use centered kernel matrix and normalized eigenvector!}$$

# Projections of new points onto PCs in feature space

For data points  $\underline{\mathbf{x}}^{(\alpha)}$  we have onto the  $k$ -th PC:

$$u_k \left( \underline{\mathbf{x}}^{(\alpha)} \right) = \sum_{\beta=1}^p a_k^{(\beta)} K_{\beta\alpha} \quad \leftarrow \text{use centered kernel matrix and normalized eigenvector!}$$



# Projections of new points onto PCs in feature space

For data points  $\underline{\mathbf{x}}^{(\alpha)}$  we have onto the  $k$ -th PC:

$$u_k \left( \underline{\mathbf{x}}^{(\alpha)} \right) = \sum_{\beta=1}^p a_k^{(\beta)} K_{\beta\alpha} \quad \leftarrow \text{use centered kernel matrix and normalized eigenvector!}$$

More generally, for arbitrary  $\underline{\mathbf{x}}$  the projection is computed as:

$$\begin{aligned} u_k(\underline{\mathbf{x}}) &= \sum_{\beta=1}^p a_k^{(\beta)} \underline{\phi}_{(\underline{\mathbf{x}}^{(\beta)})}^T \underline{\phi}_{(\underline{\mathbf{x}})} \quad \leftarrow \text{centered feature vectors} \\ &= \sum_{\beta=1}^p a_k^{(\beta)} \left( k(\underline{\mathbf{x}}^{(\beta)}, \underline{\mathbf{x}}) - \frac{1}{p} \sum_{\delta=1}^p \tilde{K}_{\beta\delta} - \frac{1}{p} \sum_{\gamma=1}^p k(\underline{\mathbf{x}}^{(\gamma)}, \underline{\mathbf{x}}) + \frac{1}{p^2} \sum_{\gamma, \delta=1}^p \tilde{K}_{\gamma\delta} \right) \end{aligned}$$

# KPCA: Additional Thoughts

- $\text{kernel-PCA} \hat{=} \text{PCA in feature space (!!!)}$

# KPCA: Additional Thoughts

- kernel-PCA  $\hat{=}$  PCA in feature space (!!!)
  - ↪ projections onto PCs are uncorrelated
  - ↪  $\lambda_l$ : variance of the data along principal component  $k$  (in feature space)
  - ↪ “minimal” mean-squared approximation of data (in feature space)

# KPCA: Additional Thoughts

- kernel-PCA  $\hat{=}$  PCA in feature space (!!!)
  - ~> projections onto PCs are uncorrelated
  - ~>  $\lambda_k$ : variance of the data along principal component  $k$  (in feature space)
  - ~> “minimal” mean-squared approximation of data (in feature space)
- #PCs can exceed #dimensions in the original space

# KPCA: Additional Thoughts

- kernel-PCA  $\hat{=}$  PCA in feature space (!!!)
  - ↪ projections onto PCs are uncorrelated
  - ↪  $\lambda_k$ : variance of the data along principal component  $k$  (in feature space)
  - ↪ “minimal” mean-squared approximation of data (in feature space)
- #PCs can exceed #dimensions in the original space
- kernel matrices may be very large
  - ↪ only first few eigenvectors (largest eigenvalues) are of interest
  - ↪ use specialized (iterative) routines (e.g. ARPACK via eigs)