

# Machine Intelligence 2

## 1.3 Kernel PCA

Prof. Dr. Klaus Obermayer

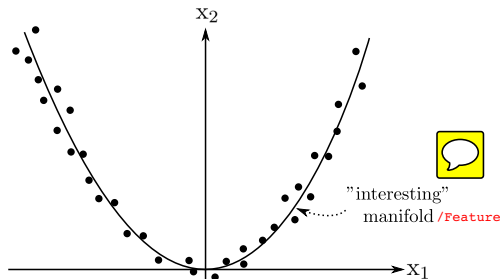
Fachgebiet Neuronale Informationsverarbeitung (NI)

SS 2017

# Kernel Principal Component Analysis

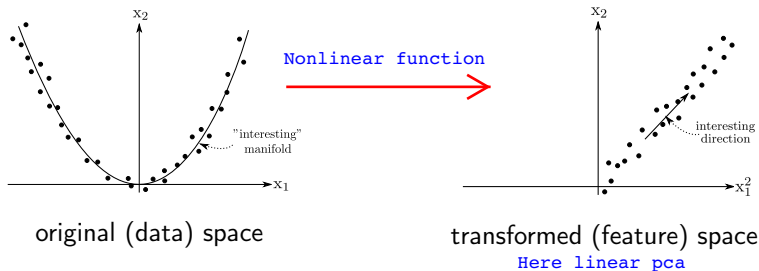
**Goal:** extraction of nonlinear features

# Kernel Principal Component Analysis: motivation



- standard PCA: two directions with high variance
- but: only one "interesting" manifold (nonlinear combination of the elementary features)

# Kernel Principal Component Analysis: intuition

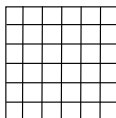


## Agenda

- 1 data preprocessing:  
nonlinear transformation into an "appropriate" feature space  
$$\underline{\phi} : \underline{\mathbf{x}} \mapsto \underline{\phi}(\underline{\mathbf{x}})$$
- 2 application of standard (linear) PCA

# Projections & Kernels

relevant feature spaces may be extremely high-dimensional



pixel image

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

example:  $d = 2$



$$\underline{\phi}(\underline{\mathbf{x}}) = (1, x_1, x_2, \dots, x_N, x_1^2, x_1x_2, x_2^2, x_1x_3, x_2x_3, x_3^2, \dots, x_N^2)^T$$

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

→ application of the **kernel trick** to avoid this problem (cf. *MI I*)

# PCA & scalar products

eigenvalue problem of PCA:

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

expansion of the eigenvectors:



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



PC in feature Space

PCs always lie in the subspace spanned by the (centered) data.

eigenvectors  $\underline{\mathbf{e}}_k \in \mathbb{R}^N$ , coefficients  $\underline{\mathbf{a}}_k \in \mathbb{R}^p$ : potential problem:  $p \gg N$

# PCA & scalar products

eigenvalue problem:

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

ansatz:

$$\underline{\mathbf{e}}_k = \sum_{\beta=1}^p a_k^{(\beta)} \underline{\mathbf{x}}^{(\beta)} \quad \underline{\mathbf{C}} = \underbrace{\frac{1}{p} \sum_{\alpha=1}^p \underline{\mathbf{x}}^{(\alpha)} \left( \underline{\mathbf{x}}^{(\alpha)} \right)^T}_{\text{centered data}}$$

$$\frac{1}{p} \sum_{\alpha, \beta=1}^p a_k^{(\beta)} \overbrace{\left[ \left( \underline{\mathbf{x}}^{(\alpha)} \right)^T \underline{\mathbf{x}}^{(\beta)} \right]}^{\text{scalar product}} \underline{\mathbf{x}}^{(\alpha)} = \lambda_k \sum_{\beta=1}^p a_k^{(\beta)} \underline{\mathbf{x}}^{(\beta)}$$

Multiply from left with  $\left(\underline{\mathbf{x}}^{(\gamma)}\right)^T$ ,  $\gamma = 1, \dots, p$ :

$$\frac{1}{p} \sum_{\alpha, \beta=1}^p a_k^{(\beta)} \overbrace{\left[ \left(\underline{\mathbf{x}}^{(\alpha)}\right)^T \underline{\mathbf{x}}^{(\beta)} \right]}^{\text{scalar product}} \overbrace{\left[ \left(\underline{\mathbf{x}}^{(\gamma)}\right)^T \underline{\mathbf{x}}^{(\alpha)} \right]}^{\text{scalar product}} = \lambda_k \sum_{\beta=1}^p a_k^{(\beta)} \overbrace{\left[ \left(\underline{\mathbf{x}}^{(\gamma)}\right)^T \underline{\mathbf{x}}^{(\beta)} \right]}^{\text{scalar product}}$$

$$\left(\underline{\mathbf{x}}^{(\alpha)}\right)^T \underline{\mathbf{x}}^{(\beta)} =: K_{\alpha\beta}$$

in matrix notation:



$$\underline{\mathbf{K}}^2 \underline{\mathbf{a}}_k = p \lambda_k \underline{\mathbf{K}} \underline{\mathbf{a}}_k$$

$\underline{\mathbf{K}}$ :  $p \times p$  matrix of scalar products between data,  $K_{\alpha\beta} = \left(\underline{\mathbf{x}}^{(\alpha)}\right)^T \underline{\mathbf{x}}^{(\beta)}$

$\lambda_k$ : variance along principal component  $\underline{\mathbf{e}}_k$

$\underline{\mathbf{a}}_k$ : Principal Component, represented in the ~~basis~~  $\left\{ \underline{\mathbf{x}}^{(\alpha)} \right\}, \alpha = 1, \dots, p$

$\underline{\mathbf{K}}$  is symmetric and positive semidefinite  $\Rightarrow$  transformed eigenvalue problem:

$$\underline{\mathbf{K}} \underline{\mathbf{a}}_k = p \lambda_k \underline{\mathbf{a}}_k$$

$\underline{\mathbf{a}}_k$  ist nicht PC!





## Remark

$\mathbf{K}$  is positive semidefinite.

For an arbitrary vector  $\mathbf{y}$ :  
(beliebig)

$$\begin{aligned}\underline{\mathbf{y}}^T \underline{\mathbf{K}} \underline{\mathbf{y}} &= \sum_{\alpha, \beta=1}^p y^{(\alpha)} \left( \underline{\mathbf{x}}^{(\alpha)} \right)^T \underline{\mathbf{x}}^{(\beta)} y^{(\beta)} \\ &= \left( \sum_{\alpha=1}^p y^{(\alpha)} \underline{\mathbf{x}}^{(\alpha)} \right)^2 \\ &\geq 0\end{aligned}$$

- ↪ only non-negative eigenvalues (but potentially 0)
- ↪ solutions of the transformed eigenvalue problem before differ only by components corresponding to PCs with zero eigenvalues

## Normalization

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



$$\underline{\mathbf{e}}_k = \sum_{\alpha,\beta=1}^p a_k^{(\alpha)} \left( \underline{\mathbf{x}}^{(\alpha)} \right)^T \underline{\mathbf{x}}^{(\beta)} a_k^{(\beta)}$$

$$= \underline{\mathbf{a}}_k^T \underline{\mathbf{K}} \underline{\mathbf{a}}_k = p \lambda_k \underline{\mathbf{a}}_k^2 \stackrel{!}{=} 1$$

$$\underline{\mathbf{a}}_k^{\text{norm.}} = \frac{1}{\sqrt{p \lambda_k}} \underline{\mathbf{a}}_k$$



$p$ : number of Data points

Comes from The Definition of The covariance matrix

# Projecting onto PCs

feature extraction:




$$u_k(\underline{\mathbf{x}}) = \underline{\mathbf{e}}_k^T \cdot \underline{\mathbf{x}}$$

$$= \sum_{\beta=1}^p a_k^{(\beta)} \underbrace{\left[ \left( \underline{\mathbf{x}}^{(\beta)} \right)^T \cdot \underline{\mathbf{x}} \right]}_{\text{scalar product}}$$

Centered Data does Not imply centered feature Space?

# The kernel trick

Allows work in high dimensional space

$$\underline{\phi} : \underline{\mathbf{x}} \xrightarrow{\text{nonlinear transformation}} \underline{\phi}_{(\underline{\mathbf{x}})}$$


## Kernel trick

⇒ formulate PCA in feature space (replace  $\underline{\mathbf{x}}^{(\alpha)}$  by  $\underline{\phi}_{(\underline{\mathbf{x}}^{(\alpha)})}$ )

⇒ replace all scalar products by "kernel functions"

$$\overset{\text{Not easy}}{\underline{\phi}_{(\underline{\mathbf{x}})}^T \underline{\phi}_{(\underline{\mathbf{x}}')}} \longleftrightarrow \overset{\text{Easy}}{k(\underline{\mathbf{x}}, \underline{\mathbf{x}}')}$$

## Mercer's theorem

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

---

If a linear method can be formulated solely in terms of scalar products, a nonlinear version can be derived without an explicit projection into the (high-dimensional) feature space!

# Mercer's theorem

## Statement of the theorem

Every **positive semidefinite** kernel  $k$  corresponds to a scalar product in some metric feature space.

Consider

- $D \subset \mathbb{R}^N$  compact *subset of data space*
- $k : D \times D \rightarrow \mathbb{R}$  is a continuous and symmetric function ("kernel")
- $T_k$  is the corresponding integral operator



$$T_k : L_2(D) \rightarrow L_2(D),$$
$$(T_k f)_{(\underline{x})} := \int_D k(\underline{x}, \underline{x}') f(\underline{x}') d\underline{x}'$$

with eigenvalues  $\lambda_j$  and normalized eigenfunctions  $\psi_j \in L_2(D)$

# Mercer's theorem: Condition and its consequences

## Essential part

If  $T_k$  is positive semidefinite, i.e.

$$\langle T_k f, f \rangle = \int_{D \times D} k(\underline{\mathbf{x}}, \underline{\mathbf{x}}') f(\underline{\mathbf{x}}) f(\underline{\mathbf{x}}') d\underline{\mathbf{x}} d\underline{\mathbf{x}}' \geq 0 \quad \forall f \in L_2(D)$$

then  $k(\underline{\mathbf{x}}, \underline{\mathbf{x}}') = \sum_{j=1}^M \lambda_j \psi_j(\underline{\mathbf{x}}) \psi_j(\underline{\mathbf{x}}')$  with  $\lambda_j \geq 0$

$k$  corresponds to a scalar product in a  $M$ -dim. space:

$$\underline{\phi} : \underline{\mathbf{x}} \mapsto \left( \sqrt{\lambda_1} \psi_1(\underline{\mathbf{x}}), \sqrt{\lambda_2} \psi_2(\underline{\mathbf{x}}), \dots, \sqrt{\lambda_M} \psi_M(\underline{\mathbf{x}}), \right)^T$$

$$\implies k(\underline{\mathbf{x}}, \underline{\mathbf{x}}') = \underline{\phi}(\underline{\mathbf{x}})^T \underline{\phi}(\underline{\mathbf{x}}') \quad (\text{with } M \leq \infty)$$



# Common 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



$$k(\underline{\mathbf{x}}, \underline{\mathbf{x}}') = \tanh \left\{ \underbrace{K \underline{\mathbf{x}}^T \underline{\mathbf{x}}' + \theta}_{\text{In Original Space}} \right\}$$

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

# Centering the kernel matrix



$$\frac{1}{p} \sum_{\alpha=1}^p \underline{\mathbf{x}}^{(\alpha)} \stackrel{!}{=} \underline{\mathbf{0}} \quad \nrightarrow \quad \frac{1}{p} \sum_{\alpha=1}^p \underline{\phi}(\underline{\mathbf{x}}^{(\alpha)}) = \underline{\mathbf{0}}$$

"centered" feature vectors:

$$\underbrace{\underline{\phi}(\underline{\mathbf{x}}^{(\alpha)})}_{\text{"centered" feature vectors}} = \tilde{\underline{\phi}}(\underline{\mathbf{x}}^{(\alpha)}) - \frac{1}{p} \sum_{\gamma=1}^p \underbrace{\tilde{\underline{\phi}}(\underline{\mathbf{x}}^{(\gamma)})}_{\text{uncentered feature vectors}}$$



# Centering the kernel matrix

$$\begin{aligned}
 K_{\alpha\beta} &= \phi(\mathbf{x}^{(\alpha)}) \cdot \phi(\mathbf{x}^{(\beta)}) \\
 &= \left( \tilde{\phi}(\mathbf{x}^{(\alpha)}) - \frac{1}{p} \sum_{\gamma=1}^p \tilde{\phi}(\mathbf{x}^{(\gamma)}) \right) \left( \tilde{\phi}(\mathbf{x}^{(\beta)}) - \frac{1}{p} \sum_{\delta=1}^p \tilde{\phi}(\mathbf{x}^{(\delta)}) \right) \\
 &= \tilde{\phi}(\mathbf{x}^{(\alpha)}) \tilde{\phi}(\mathbf{x}^{(\beta)}) - \frac{1}{p} \sum_{\delta=1}^p \tilde{\phi}(\mathbf{x}^{(\alpha)}) \tilde{\phi}(\mathbf{x}^{(\delta)}) \\
 &\quad - \frac{1}{p} \sum_{\gamma=1}^p \tilde{\phi}(\mathbf{x}^{(\gamma)}) \tilde{\phi}(\mathbf{x}^{(\beta)}) + \frac{1}{p^2} \sum_{\gamma,\delta=1}^p \tilde{\phi}(\mathbf{x}^{(\gamma)}) \tilde{\phi}(\mathbf{x}^{(\delta)})
 \end{aligned}$$
  

$$\begin{aligned}
 \text{[Yellow box with speech bubble]} &= \underbrace{\tilde{K}_{\alpha\beta}}_{=k(\mathbf{x}^{(\alpha)}, \mathbf{x}^{(\beta)})} - \underbrace{\frac{1}{p} \sum_{\delta=1}^p \tilde{K}_{\alpha\delta}}_{\text{row avg.}} - \underbrace{\frac{1}{p} \sum_{\gamma=1}^p \tilde{K}_{\gamma\beta}}_{\text{col. avg.}} + \underbrace{\frac{1}{p^2} \sum_{\gamma,\delta=1}^p \tilde{K}_{\gamma\delta}}_{\text{matrix avg.}}
 \end{aligned}$$

# Centering & projections

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

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

More generally, for new/arbitrary  $\underline{\mathbf{x}}^{(\alpha)}$  the projection is computed as:

"just a longer forumula" --> Falls Data point "außerhalb" liegt

$$u_k \left( \underline{\phi}_{(\underline{\mathbf{x}}^{(\alpha)})} \right) = \sum_{\beta=1}^p a_k^{(\beta)} \underline{\phi}_{(\underline{\mathbf{x}}^{(\beta)})}^T \underline{\phi}_{(\underline{\mathbf{x}}^{(\alpha)})} \quad \leftarrow \text{centered feature vectors}$$

$$= \sum_{\beta=1}^p a_k^{(\beta)} \left( \left[ \underline{\phi}_{(\underline{\mathbf{x}}^{(\beta)})} - \frac{1}{p} \sum_{\gamma=1}^p \underline{\phi}_{(\underline{\mathbf{x}}^{(\gamma)})} \right]^T \left[ \underline{\phi}_{(\underline{\mathbf{x}}^{(\alpha)})} - \frac{1}{p} \sum_{\delta=1}^p \underline{\phi}_{(\underline{\mathbf{x}}^{(\delta)})} \right] \right)$$


$$= \sum_{\beta=1}^p a_k^{(\beta)} \left( k(\underline{\mathbf{x}}^{(\beta)}, \underline{\mathbf{x}}^{(\alpha)}) - \frac{1}{p} \sum_{\delta=1}^p \tilde{K}_{\beta\delta} - \frac{1}{p} \sum_{\gamma=1}^p k(\underline{\mathbf{x}}^{(\gamma)}, \underline{\mathbf{x}}^{(\alpha)}) + \frac{1}{p^2} \sum_{\gamma,\delta=1}^p \tilde{K}_{\gamma\delta} \right)$$


Backprojection

- ① calculation of the **un-normalized** kernel matrix 

$$\tilde{K}_{\alpha\beta} = k(\underline{\mathbf{x}}^{(\alpha)}, \underline{\mathbf{x}}^{(\beta)}), \quad \alpha, \beta = 1, \dots, p$$

- ② centering **of the Kernel Matrix**


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

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

- ④ **normalization** of **eigenvectors** to unit length

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

- ⑤ calculation of projections

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

# Comments

- kernel-PCA  $\hat{=}$  PCA in feature space : Same properties
  - $\leadsto$  projections onto PCs are uncorrelated
  - $\leadsto \lambda_k$ : variance of the data along PC  $k$  (in feature space)
- #PCs typically exceed #dimensions in the original space
- expansion of PCs into data points is not sparse
  - $\leadsto$  calculating projections can be computationally expensive
  - $\leadsto$  use expansions with less data points  $q < p$



$$\underline{\mathbf{e}}_k = \sum_{\beta=1}^p a_k^{(\beta)} \underline{\phi}(\underline{\mathbf{x}}^{(\beta)}) \quad \leadsto \quad \hat{\underline{\mathbf{e}}}_k = \sum_{\gamma=1}^q \hat{a}_k^{(\gamma)} \underline{\phi}(\overbrace{\underline{\mathbf{z}}^{(\gamma)}}^{\underline{\mathbf{x}}^{(i_\gamma)}})$$

$$(\underline{\mathbf{e}}_k - \hat{\underline{\mathbf{e}}}_k)^2 =: \varphi \rightarrow \min_{\hat{a}_k^{(\gamma)}, \underline{\mathbf{z}}^{(\gamma)}}$$

$$\varphi = 1 + \sum_{\gamma, \delta=1}^q \hat{a}_k^{(\gamma)} \hat{a}_k^{(\delta)} k(\underline{\mathbf{z}}^{(\gamma)}, \underline{\mathbf{z}}^{(\delta)}) - 2 \sum_{\beta=1}^p \sum_{\gamma=1}^q a_k^{(\beta)} \hat{a}_k^{(\gamma)} k(\underline{\mathbf{x}}^{(\beta)}, \underline{\mathbf{z}}^{(\gamma)})$$

# Comments

$p \times p$

- $p \gg N$ : kernel matrices may be very large
  - ↪ only eigenvectors with largest eigenvalues are of interest
  - ↪ use specialized (iterative) routines (e.g. ARPACK via `eigs`)

⇒ analysis is performed in feature space (not data space)
- kernel PCA can be used for feature extraction & dimensionality reduction
  - e.g. solve classification problems in feature space
- optimal kernel parameters ( $\sigma$ ,  $d$ , etc.) depend on data & task
  - selection via cross-validation possible for classification tasks
  - no general measure-of-goodness of the PC projections available
- custom kernels can be used (any positive definite kernel matrix)

## Application: feature extraction

0 1 2 3 4 5 6 7 8 9

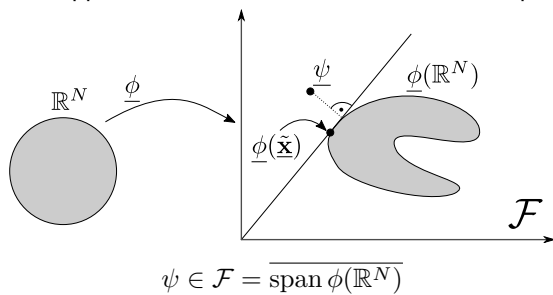
		test error (%) for different polynomial kernels						
# of component		1	2	3	4	5	6	7
32		9.6	8.8	8.1	8.5	9.1	9.3	10.8
64		8.8	7.3	6.8	6.7	6.7	7.2	7.5
128		8.6	5.8	5.9	5.9	5.8	6.0	6.8
256		8.7	5.5	5.3	5.2	5.2	5.4	5.4
512		n.a.	4.9	4.6	4.4	5.1	4.6	4.9
1024		n.a.	4.9	4.3	4.4	4.6	4.8	4.6
2048		n.a.	4.9	4.2	4.1	4.0	4.3	4.4

- Test error rates on the USPS handwritten digit database
- linear SVMs trained on nonlinear Principal Components
- nonlinear PCs extracted by PCA with a polynomial kernel (degrees 1 through 7)
- dimensionality of the space is 256 (16x16 pixel images)

*Source: Schölkopf, 2002*

# Reconstruction

- reconstruction in data space non-trivial
  - data space is mapped to a low-dim. manifold in feature space

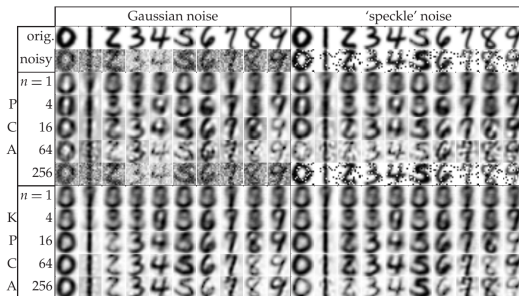


- problem: in general there is no "pre-image"  $\tilde{\underline{\mathbf{x}}}$  s.t.  $\underline{\psi} = \underline{\phi}(\tilde{\underline{\mathbf{x}}})$
- solution: calculation of approximate "pre-images":

$$\tilde{\underline{\mathbf{x}}} = \underset{\underline{\mathbf{x}}}{\operatorname{argmin}} \left\| \underline{\phi}(\underline{\mathbf{x}}) - \underline{\psi} \right\|^2$$

- algorithms: Schölkopf & Smola, ch. 18 (e.g. impl. in scikit-learn)

# Application: denoising

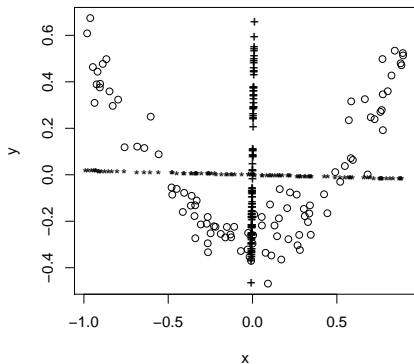


- Denoising of USPS data
- *First row:* original data (digits); *Second row:* noise added to original digits (Gaussian and "speckle")
- *Following five rows:* reconstruction of the noisy digits achieved with linear PCA using  $n = 1, 4, 16, 64, 256$  components
- *Last five rows:* reconstruction of the noisy digits achieved with Kernel PCA using the same number of components
- dimensionality of the space is 256 (16x16 pixel images)

Source: Schölkopf, 2002



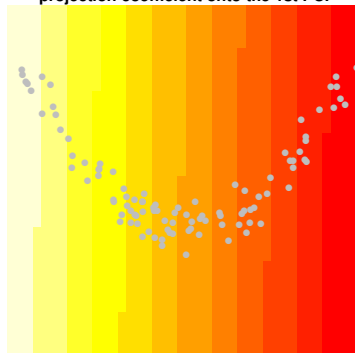
# Parabola example: PCA



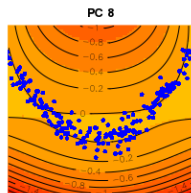
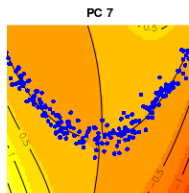
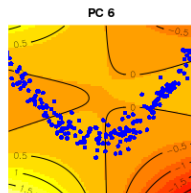
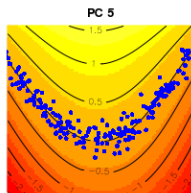
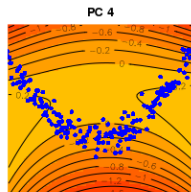
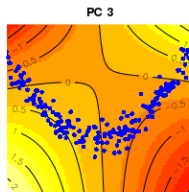
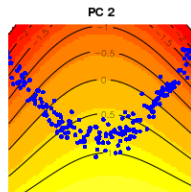
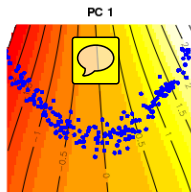
Negative to positive First pc



Color indicates intervals of the projection coefficient onto the 1st PC.

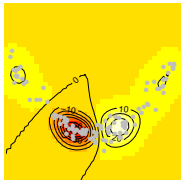


# Parabola example: kPCA with polynomial kernel

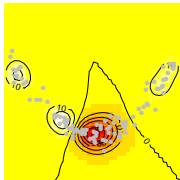


# Parabola example: kPCA with RBF-kernel

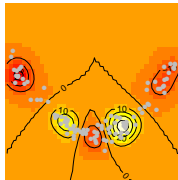
PC 1



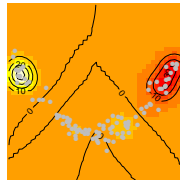
PC 2



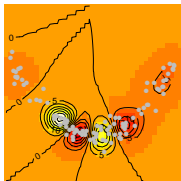
PC 3



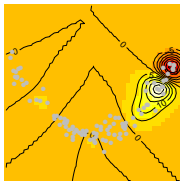
PC 4



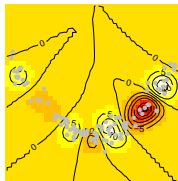
PC 5



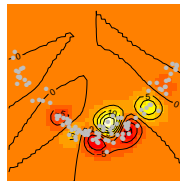
PC 6



PC 7



PC 8



# Parabola example: dimension reduction

