

Kernel PCA and Relevant Dimension Estimation

Lecture 11 – Machine Learning 2 (SS 2024), TU Berlin

Pat Chormai
(2024.07.04)

Recap: Kernel Methods

Linear model on nonlinear feature space (through $\phi : \mathcal{X} \rightarrow \mathcal{F}$)

$$f(\mathbf{x}) = \langle \boldsymbol{\beta}, \phi(\mathbf{x}) \rangle_{\mathcal{F}},$$

Choose the model class from the span of the data (in the feature space \mathcal{F})

$$\boldsymbol{\beta} = \sum_{i=1}^N \alpha_i \phi(\mathbf{x}_i).$$

Then,

$$f(\mathbf{x}) = \sum_i \alpha_i \langle \phi(\mathbf{x}_i), \phi(\mathbf{x}) \rangle_{\mathcal{F}}$$

Recap: Kernel Trick

Challenge: finding $\phi(\cdot)$ and computing $\phi(\mathbf{x})$ might be difficult.

Positive Definite Kernels (p.d.)¹

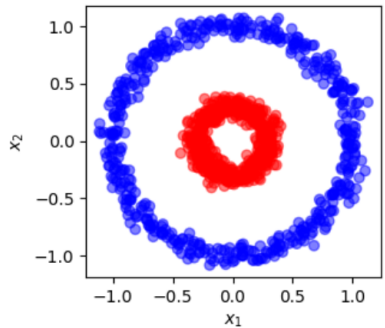
If $k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ is a p.d. kernel function, then

$$k(\mathbf{x}, \mathbf{x}') = \langle \phi(\mathbf{x}), \phi(\mathbf{x}') \rangle$$

Consequently, if our algorithm only relies on the dot product between data points, then

- no need to compute $\phi(\mathbf{x})$ explicitly

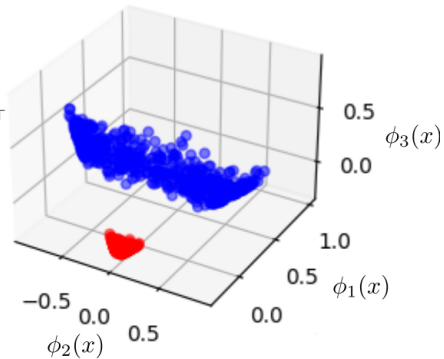
¹ Associated kernel matrix K (with $K_{ij} = k(\mathbf{x}_i, \mathbf{x}_j)$) is positive semi-definite.



Input Space

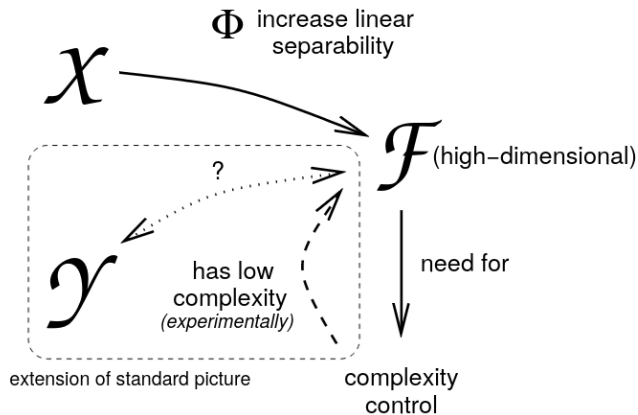
$$\phi(\mathbf{x}) = (x_1^2, \sqrt{2}x_1x_2, x_2^2)^\top$$

$$k(\mathbf{x}, \mathbf{x}') = (\mathbf{x}^\top \mathbf{x}')^2$$



Feature Space

How does the transformation relate to 'label' information?



[Braun et al., 2008]

Outline

- ▶ Kernel PCA (Schölkopf et al. 1998)
- ▶ Relevant Dimension Estimation (Braun et al. 2008)

Recap: Principal Component Analysis

Let $\mathcal{D} = \{\mathbf{x}_1, \dots, \mathbf{x}_N \in \mathbb{R}^d\}$ and assume that $\sum_i \mathbf{x}_i = 0$

Formulation

$$\min_{\mathbf{u} \in \mathbb{R}^d} \sum_{i=1}^N \|\mathbf{x}_i - \mathbf{u}\mathbf{u}^\top \mathbf{x}_i\|_2^2 \quad \text{subject to} \quad \mathbf{u}^\top \mathbf{u} = 1$$

Solution: 1st eigenvector of $\Sigma = \frac{1}{N} \sum_i \mathbf{x}_i \mathbf{x}_i^\top$

Kernel PCA

Let $\phi : \mathcal{X} \rightarrow \mathcal{F}$ and assume that $\sum_i \phi(\mathbf{x}_i) = 0$

Formulation

Perform PCA on

$$\hat{\Sigma} = \frac{1}{N} \sum_i \phi(\mathbf{x}_i) \phi(\mathbf{x}_i)^\top$$

Kernel PCA: Solution¹

Let $K \in \mathbb{R}^{N \times N}$ be the kernel matrix associated to a p.d. kernel function

$$k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$$

Properties of Kernel Matrix K

- ▶ $K = K^\top \rightsquigarrow$ orthogonal eigenvectors (Spectral Theorem)
- ▶ K is positive semi-definite \rightsquigarrow all eigenvalues $\lambda_n(K) \geq 0$

From Eigenvalue Decomposition,

$$K = U \Lambda U^\top,$$

where the columns of $U \in \mathbb{R}^{N \times N}$ are its eigenvectors and $\Lambda_{\tau\tau} = \lambda_\tau(K)$.

¹ We only state results and refer to the original paper (Schölkopf et al. 1998) or (Bishop 2007, Chapter 12.3) for the derivation.

Kernel PCA: Solution (cont.)

Let $\mathbf{u}_\tau \in \mathbb{R}^N$ be the τ th eigenvector of K and define $\boldsymbol{\alpha}^{(\tau)} = \frac{1}{\sqrt{\lambda_\tau(K)}} \mathbf{u}_\tau$

τ th Eigenvector of Kernel PCA $\hat{\Sigma}$

$$\mathbf{v}_\tau = \sum_i \alpha_i^{(\tau)} \phi(\mathbf{x}_i)$$

Computing Principal Components

Consider a datapoint \mathbf{x} and the τ th eigenvector of $\hat{\Sigma}$

τ th Principal Component of Datapoint \mathbf{x}

$$\begin{aligned}\langle \mathbf{v}_\tau, \phi(\mathbf{x}) \rangle &= \left\langle \sum_i \alpha_i^{(\tau)} \phi(\mathbf{x}_i), \phi(\mathbf{x}) \right\rangle \\ &= \sum_i \alpha_i^{(\tau)} \langle \phi(\mathbf{x}_i), \phi(\mathbf{x}) \rangle \\ &= \sum_i \alpha_i^{(\tau)} k(\mathbf{x}_i, \mathbf{x})\end{aligned}$$

Relaxing Centering Assumption

Consider a kernel function k with its associated feature map ϕ and kernel matrix K
Let $\tilde{\phi}(\mathbf{x}) = \phi(\mathbf{x}) - \frac{1}{N} \sum_i \phi(\mathbf{x}_i)$ be the **centered feature map**

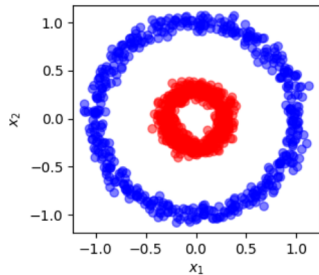
Kernel Matrix \tilde{K} Corresponding to Centered Feature Map

$$\tilde{K} = K - \mathbf{1}_N K - K \mathbf{1}_N + \mathbf{1}_N K \mathbf{1}_N,$$

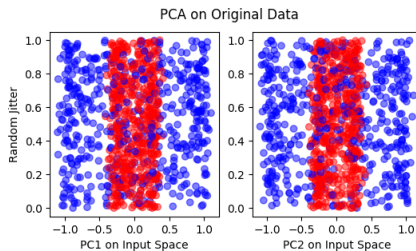
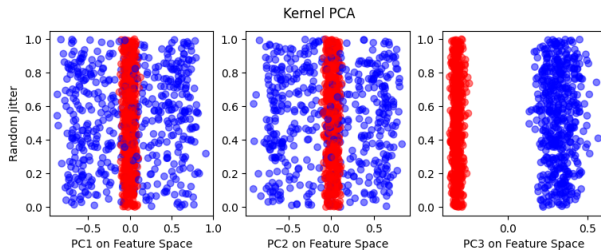
where $\mathbf{1}_N \in \{\frac{1}{N}\}^{N \times N}$.

Derivation Sketch¹: Expand $\langle \tilde{\phi}(\mathbf{x}), \tilde{\phi}(\mathbf{x}') \rangle$ and write each term with $k(\cdot, \cdot)$

¹ See Bishop 2007, Chapter 12.3 for the complete derivation



Input Space

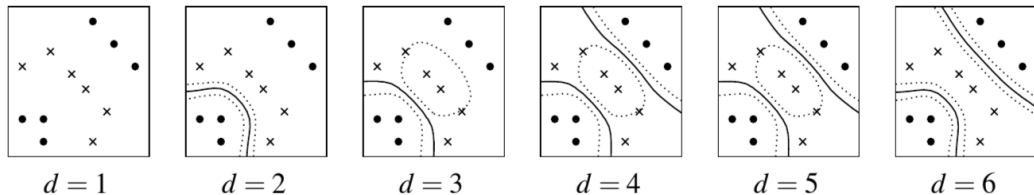


Remarks on Aspects of PCA and Kernel PCA

	PCA	Kernel PCA
Eigenvectors of	$\sum_i \mathbf{x}_i \mathbf{x}_i^\top$	$\sum_i \tilde{\phi}(\mathbf{x}_i) \tilde{\phi}(\mathbf{x}_i)^\top$
τ th Principal Component	$\mathbf{u}_\tau^\top \mathbf{x}$	$\sum_i \alpha_i^{(\tau)} k(\mathbf{x}_i, \mathbf{x})$
Reconstruction	$\sum_\tau \mathbf{u}_\tau \mathbf{u}_\tau^\top \mathbf{x}$	<i>not straightforward</i> ¹

¹ See Mika et al. 1998

Observation: Only Few Components Needed to Preserve Decision Boundary



(Montavon et al. 2011)

Projecting Label Information on Kernel PCA Components

Consider a vector $\mathbf{y} \in \mathbb{R}^N$ (each entry corresponding to x_i)

Recall $K = U\Lambda U^\top \in \mathbb{R}^{N \times N}$ whose τ th column is denoted as \mathbf{u}_τ

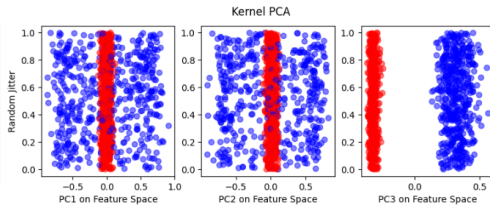
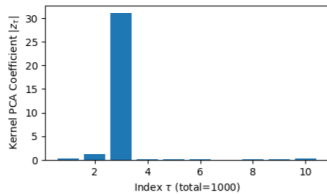
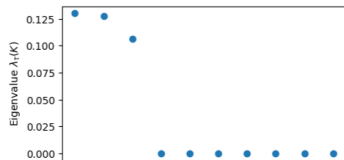
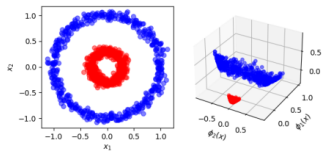
Lemma 1¹ Projection of Y onto Leading d Kernel PCA Components

$$\sum_{\tau=1}^d \mathbf{u}_\tau \mathbf{u}_\tau^\top \mathbf{y}$$

Definition τ th Kernel PCA Coefficient

$$z_\tau = \mathbf{u}_\tau^\top \mathbf{y}$$

¹ Braun et al. 2008



Relevant Information

Assumption Label Corrupted with Additive Noise N_i

$$Y_i = g(X_i) + N_i$$

Definition Relevant Information g

$$\begin{aligned} g(X) &= \mathbb{E}[Y|X] && \text{(Population)} \\ G &= (g(X_1), \dots, g(X_N)) && \text{(Sample)} \end{aligned}$$

Remarks

- ▶ (regression): $\mathbb{E}[Y|X]$ is the true function.
- ▶ (binary classification): $\mathbb{E}[Y|X] = P(Y = 1|X) - P(Y = -1|X)$

Theoretical Result

Theorem 1 (Braun et al. 2008) Upperbound on Contribution of Eigenvector on Relevant Information

If the learning problem can be represented by the kernel asymptotically, then

$$\frac{1}{N} |\mathbf{u}_\tau^\top G| \leq \lambda_\tau C(N) + E(N),$$

where $C(N), E(N)$ are some constants and $E(N) \rightarrow 0$ as $N \rightarrow \infty$.

Interpretation: the relevant information about Y is contained in the leading kernel PCA directions up to a small error.

Relevant Dimension Estimation (RDE)

Formulation

Estimate number of dimensions $d \ll N$ such that

$$\forall i > d : \quad |\mathbf{u}_i^\top G| \approx 0$$

Implication: $\forall i > d$, $\mathbf{u}_i^\top Y$ captures the noise.

RDE by Fitting a Two-Component Model

Assumption¹: Modeling Kernel Coefficients with Two Zero-Mean Gaussians

Let d be a cut-off point that splits kernel coefficients z_τ 's into two parts (*relevant information and noise*)

$$z_\tau \sim \begin{cases} \mathcal{N}(0, \sigma_1^2) & 1 \leq \tau \leq d, \\ \mathcal{N}(0, \sigma_2^2) & d < \tau \leq N \end{cases}$$

where $\sigma_1^2 = \frac{1}{N} \sum_{i=1}^d z_i^2$ and $\sigma_2^2 = \frac{1}{N-d} \sum_{i=d+1}^N z_i^2$.

¹ See Braun et al. 2008 for a more general approach via Leave-One-Out Cross-Validation

RDE by Fitting a Two-Component Model (cont.)

The negative log-likelihood $\ell(d)$ is proportional to

$$-\log \ell(d) \propto \frac{d}{n} \log \sigma_1^2 + \frac{n-d}{d} \log \sigma_2^2$$

Estimated Relevant Dimension via Maximum Likelihood

$$\hat{d} = \arg \min_{1 \leq d < N} -\log \ell(d)$$

Estimating Noise Level

Project labels \mathbf{y} to onto the first \hat{d} kernel PCA components

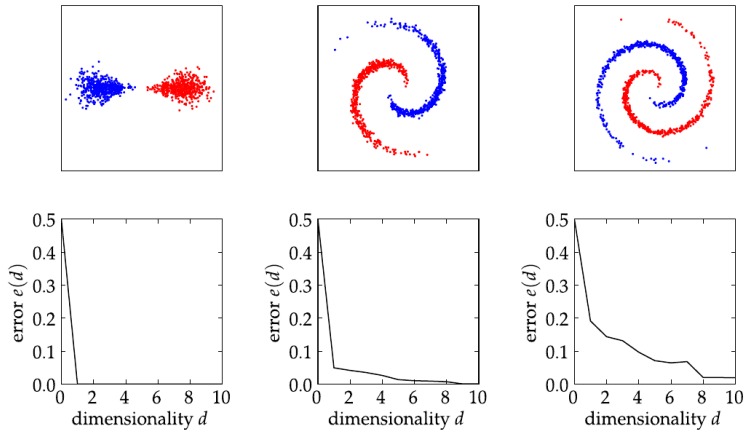
$$\hat{\mathbf{y}} = \sum_{\tau=1}^{\hat{d}} \mathbf{u}_{\tau} \mathbf{u}_{\tau}^{\top} \mathbf{y}$$

Estimate Noise Level

$$e(\hat{d}) = \sum_{i=1}^N L(\hat{y}_i, y_i),$$

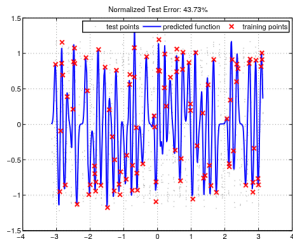
where L is a loss function

Relationship between Data Complexity and Noise Level

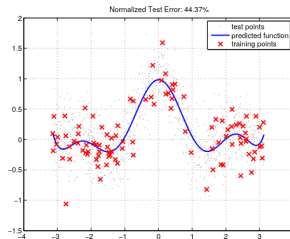


(Montavon et al. 2011)

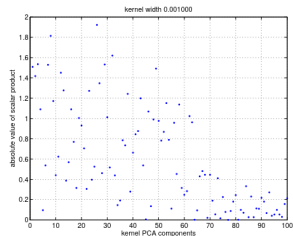
Applications of RDE: Dataset Assessment



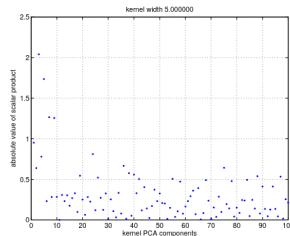
(a) A complex data set.



(b) A noisy data set.



(c) Kernel PCA coefficients for the complex data set.



(d) Kernel PCA coefficients for the noisy data set.

Applications of RDE: Dataset Assessment (cont.)

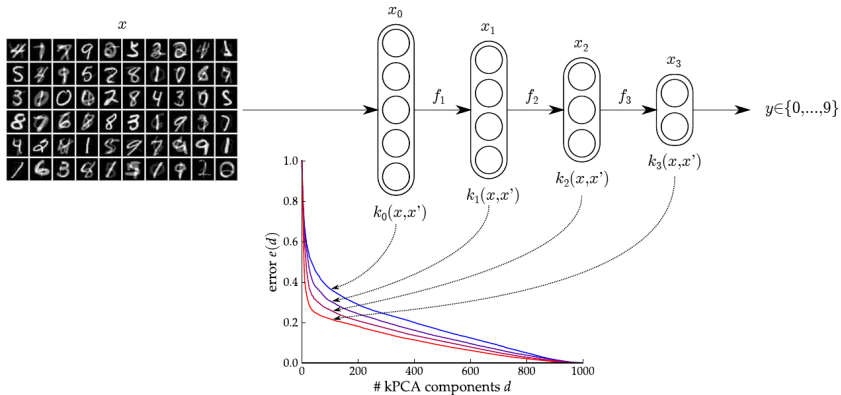
Dataset	RDE Method	Estimated Dimension \hat{d}	Noise Level $e(\hat{d})$
Complex Dataset	TCM	50	16.07%
	LOO-CV	25	40.59 %
Noisy Dataset	TCM	9	40.71%
	LOO-CV	9	40.71 %

Applications of RDE: Dataset Assessment (cont.)

Data Scenario	\hat{d}	$e(\hat{d})$	Possible Mitigation Strategy
Noisy	small	large	find ways to reduce label noise
Complex	large	small	model selection ¹ and/or acquire more data

¹ e.g., incorporate domain knowledge in order to determine a more appropriate kernel function and its parameters

Application: Layer-wise Analysis of Neural Network Representation



layer-wise reduction of both noise and dimensionality
(Montavon 2013)

Summary

- ▶ Kernel PCA
- ▶ Relevant Dimension Estimation and its applications on
 - Data Assessment
 - Analysis of Neural Network Representation

References

- Bishop, Christopher M. (2007). *Pattern recognition and machine learning, 5th Edition*. Information science and statistics. Springer. Chap. 12.3.
- Braun, Mikio L. et al. (2008). "On Relevant Dimensions in Kernel Feature Spaces". In: *J. Mach. Learn. Res.* 9, pp. 1875–1908.
- Mika, Sebastian et al. (1998). "Kernel PCA and de-noising in feature spaces". In: *Advances in neural information processing systems* 11.
- Montavon, Grégoire (2013). "On layer-wise representations in deep neural networks". PhD thesis. Technische Universität Berlin.
- Montavon, Grégoire et al. (2011). "Kernel Analysis of Deep Networks.". In: *Journal of Machine Learning Research* 12.9.
- Schölkopf, Bernhard et al. (1998). "Nonlinear Component Analysis as a Kernel Eigenvalue Problem". In: *Neural Computation* 10.5, pp. 1299–1319.
DOI: 10.1162/089976698300017467.