

HOW ARE THE CENTERED KERNEL PRINCIPAL COMPONENTS RELEVANT TO REGRESSION TASK? — AN EXACT ANALYSIS

Masahiro Yukawa^{1,2*} Klaus-Robert Müller^{3,4,5*} Yuto Ogino¹

1. Department of Electronics and Electrical Engineering, Keio University, Japan
2. Center for Advanced Intelligence Project, RIKEN, Japan
3. Department of Machine Learning, Berlin Institute of Technology, Germany
4. Department of Brain and Cognitive Engineering, Korea University, Republic of Korea
5. Max Planck Institute for Informatics, Germany

ABSTRACT

We present an exact analytic expression of the contributions of the kernel principal components to the relevant information in a nonlinear regression problem. A related study has been presented by Braun, Buhmann, and Müller in 2008, where an upper bound of the contributions was given for a general supervised learning problem but with “uncentered” kernel PCAs. Our analysis clarifies that the relevant information of a kernel regression under explicit centering operation is contained in a finite number of leading kernel principal components, as in the “uncentered” kernel-PCA case, if the kernel matches the underlying nonlinear function so that the eigenvalues of the centered kernel matrix decay quickly. We compare the regression performances of the least-square-based methods with the centered and uncentered kernel PCAs by simulations.

Index Terms— nonlinear regression, kernel PCA, reproducing kernel Hilbert space, spectral decomposition

1. INTRODUCTION

Kernel principal component analysis (PCA) [1] is a nonlinear technique to find the principal components in a so-called *feature space* to which the sample vectors are mapped “nonlinearly” with a reproducing kernel. It can be used for dimensionality reduction, noise reduction, pre-processing of classification/regression/clustering, etc. (see, e.g., [2]). Although kernel PCA only looks at the samples, it is actually relevant to a task (a supervised learning problem such as classification/regression). This has been shown theoretically by Braun, Buhmann, and Müller [3]. More specifically, it has been shown that the *relevant information* about the task is typically contained in the subspace generated by a small number of leading kernel PCA components. Here, the relevant information is a noise-free version of the output (label) vector. The derivation therein is rather complicated since the integral operator associated with the reproducing kernel is considered first and its spectral decomposition is then truncated to attain a bound. This makes the analysis applicable only to an “uncentered” kernel PCA. We break this limitation in this paper. Uncentered (linear) PCAs have been used in different fields such as climatology [4], neuroimaging data [5], microarrays [6], among many others. Its use in reduced-rank signal processing has been studied, e.g., by Scharf [7]. Cadima and Jolliffe have studied the relationships between uncentered and centered (i.e., standard) PCAs, in particular the relationships between the eigenvalues of the

two PCA variants [8]. However, the relationships between uncentered and centered (linear/kernel) PCAs have not been well studied in the context of regression.

In this paper, we focus on a nonlinear regression problem, and study how much the eigenfunctions of the covariance operator of kernel feature vectors contribute to the relevant information. Here, the contribution of an eigenfunction is the quantity of which an upper bound has been given in [3, 9] for an uncentered kernel PCA. We present its exact analytic expression (instead of its bound) for both centered and uncentered kernel PCAs. The contribution essentially decays at the same rate as the eigenvalues. A numerical example shows that the centered and uncentered kernel PCAs have a similar tendency in test-error-decay characteristics.

2. CONTRIBUTIONS OF PRINCIPAL EIGENVECTORS IN LINEAR REGRESSION

In this section, we discuss the linear PCA of which the results are useful to discuss kernel PCA. Let $\{\mathbf{x}_i\}_{i=1}^n$ be a set of n sample vectors $\mathbf{x}_i \in \mathbb{R}^N$. We consider a simple linear regression model

$$y_i := \mathbf{x}_i^T \mathbf{w} + \epsilon_i, \quad i = 1, 2, \dots, n, \quad (1)$$

where $\mathbf{w} \in \mathbb{R}^N$ is the regression vector, ϵ_i is the zero-mean additive noise, and y_i is the output corresponding to the sample vector \mathbf{x}_i ($(\cdot)^T$ stands for transposition). The regression model in (1) can be written in a vector form as

$$\mathbf{y} := [y_1, y_2, \dots, y_n]^T = \mathbf{X}^T \mathbf{w} + \boldsymbol{\epsilon}, \quad (2)$$

where $\mathbf{X} := [\mathbf{x}_1 \ \mathbf{x}_2 \ \dots \ \mathbf{x}_n] \in \mathbb{R}^{N \times n}$ and $\boldsymbol{\epsilon} := [\epsilon_1, \epsilon_2, \dots, \epsilon_n]^T$. Define the *relevant information vector* [3]

$$\mathbf{g} := [g_1, g_2, \dots, g_n] := [E(y_1|\mathbf{x}_1), E(y_2|\mathbf{x}_2), \dots, E(y_n|\mathbf{x}_n)]^T, \quad (3)$$

which is a noise-free version of \mathbf{y} . In the present regression case, it holds that

$$\mathbf{y} = \mathbf{g} + \boldsymbol{\epsilon}. \quad (4)$$

Given any closed linear subspace \mathcal{S} of an arbitrary real Hilbert space \mathcal{X} , we denote by $P_{\mathcal{S}}$ the orthogonal projection operator onto \mathcal{S} that maps any point $x \in \mathcal{X}$ to its closest point in \mathcal{S} .

2.1. Uncentered linear PCA case

An uncentered linear PCA is carried out through the singular value decomposition of the uncentered sample matrix, which is denoted as

$$\mathbf{X} := [\mathbf{x}_1 \ \mathbf{x}_2 \ \dots \ \mathbf{x}_n] = \mathbf{U}_u \boldsymbol{\Sigma}_u \mathbf{V}_u^T, \quad (5)$$

*This work was supported in part by JSPS Grants-in-Aids (15K06081, 15K13986, 15H02757), in part by the National Research Foundation of Korea funded by the Ministry of Education, Science and Technology in the BK21 program, in part by Institute for Information & Communications Technology Promotion (IITP) grant funded by the Korea government (No. 2017-0-00451), and in part by the German Research Foundation under Grant DFG MU 987/6-1, Grant SPP 1527, and Grant MU 987/14-1 as well as BMBF (BDDC).

for the coefficients $\alpha_j, j = 1, 2, \dots, n$. By (15) and (16), it follows that

$$\psi(\mathbf{x}_i) = \sum_{j=1}^n \alpha_j \kappa(\mathbf{x}_i, \mathbf{x}_j) = \boldsymbol{\alpha}^\top \mathbf{k}_i^{[u]}, \quad (17)$$

where $\mathbf{k}_i^{[u]} := [\kappa(\mathbf{x}_i, \mathbf{x}_1), \kappa(\mathbf{x}_i, \mathbf{x}_2), \dots, \kappa(\mathbf{x}_i, \mathbf{x}_n)]^\top$ and $\boldsymbol{\alpha} := [\alpha_1, \alpha_2, \dots, \alpha_n]^\top$. By (13) and (17), the output vector can be written as

$$\mathbf{y} = \mathbf{K}_u \boldsymbol{\alpha} + \boldsymbol{\epsilon}, \quad (18)$$

where the kernel matrix $\mathbf{K}_u := [\mathbf{k}_1^{[u]} \mathbf{k}_2^{[u]} \dots \mathbf{k}_n^{[u]}]$ is symmetric and positive semi-definite.

3.2. Uncentered kernel PCA case

We consider an uncentered kernel PCA based on the eigenvalue decomposition of the kernel matrix:

$$\mathbf{K}_u = \mathbf{Q}_u \boldsymbol{\Lambda}_u \mathbf{Q}_u^\top, \quad (19)$$

where $\mathbf{Q}_u := [\mathbf{q}_1^{[u]} \mathbf{q}_2^{[u]} \dots \mathbf{q}_n^{[u]}] \in \mathbb{R}^{n \times n}$ is an orthogonal matrix and $\boldsymbol{\Lambda}_u := \text{diag}(\lambda_1^{[u]}, \lambda_2^{[u]}, \dots, \lambda_n^{[u]}) \in \mathbb{R}^{n \times n}$ with $\lambda_1^{[u]} \geq \lambda_2^{[u]} \geq \dots \geq \lambda_n^{[u]} \geq 0$. Comparing (18) with (2) under (19), we immediately obtain from (8) the following relation:

$$|\mathbf{g}^\top \mathbf{q}_i^{[u]}| = \lambda_i^{[u]} |\boldsymbol{\alpha}^\top \mathbf{q}_i^{[u]}|. \quad (20)$$

Here, $\mathbf{g} = [\psi(\mathbf{x}_1), \psi(\mathbf{x}_2), \dots, \psi(\mathbf{x}_n)]^\top$ in this case (see (3) for its definition).

3.3. Centered kernel PCA case

We consider the feature vector centered in the RKHS \mathcal{H} :

$$\phi_i := \Phi(\mathbf{x}_i) - \bar{\phi}, \quad i = 1, 2, \dots, n, \quad (21)$$

where $\bar{\phi} := \frac{1}{n} \sum_{i=1}^n \Phi(\mathbf{x}_i)$. It then holds that

$$P_{\mathcal{M}}(\psi) = \eta \bar{\phi} + \sum_{j=1}^n \alpha_j \phi_j, \quad (22)$$

where $\eta := \sum_{j=1}^n \alpha_j$. Note that $\bar{\phi} \notin \mathcal{M}_\phi := \text{span}\{\phi_i\}_{i=1}^n$ in general (see Appendix). Kernel PCA is based on the spectral decomposition of the sample covariance operator¹

$$\mathbf{C} := \sum_{i=1}^n \phi_i \otimes \bar{\phi}_i, \quad (23)$$

where \otimes denotes the Schatten product,² i.e., $(\phi_i \otimes \bar{\phi}_i)f := \langle \phi_i, f \rangle \phi_i$ for any $f \in \mathcal{H}$. Since \mathbf{C} is a finite-dimensional operator (i.e., its image has a finite dimension at most n), it is a compact operator (or a completely continuous operator).³ Hence, the spectral representation theorem admits the following representation:

$$\mathbf{C} = \sum_{i=1}^r \lambda_i p_i \otimes \bar{p}_i, \quad (24)$$

¹The sample covariance operator is also expressed in the following way [2]: $\mathbf{C} := \frac{1}{n} \sum_{i=1}^n (\Phi(\mathbf{x}_i) - \bar{\phi})(\Phi(\mathbf{x}_i) - \bar{\phi})^\top$.

²The bar on ϕ_i is a part of the Schatten product.

³A linear operator is said to be a compact operator (or a completely continuous operator) if it maps a bounded set to a compact set.

where $\{p_i\}_{i=1}^r \subset \mathcal{H}$ is an orthonormal basis of the r -dimensional subspace \mathcal{M}_ϕ , and $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_r > 0$ (see Appendix).

It is known that the eigenvalue problem of \mathbf{C} is a dual problem of the eigenvalue problem of $\mathbf{K} \in \mathbb{R}^{n \times n}$ with its (i, j) entry $k_{i,j} := \langle \phi_i, \phi_j \rangle$. The eigenvalue decomposition of \mathbf{K} is denoted as

$$\mathbf{K} =: [\mathbf{k}_1 \mathbf{k}_2 \dots \mathbf{k}_n]^\top = \mathbf{Q} \boldsymbol{\Lambda} \mathbf{Q}^\top, \quad (25)$$

where $\mathbf{Q} := [\mathbf{q}_1 \mathbf{q}_2 \dots \mathbf{q}_n] \in \mathbb{R}^{n \times n}$ is an orthogonal matrix, and $\boldsymbol{\Lambda} := \text{diag}(\lambda_1, \lambda_2, \dots, \lambda_n) \in \mathbb{R}^{n \times n}$, where $\lambda_i := 0$ for $i = r + 1, \dots, n$. Define a bounded linear operator

$$A : \mathcal{H} \rightarrow \mathbb{R}^n, f \mapsto [\langle \phi_1, f \rangle, \langle \phi_2, f \rangle, \dots, \langle \phi_n, f \rangle]^\top. \quad (26)$$

Its spectral representation is then given by

$$A = \sum_{i=1}^r \sqrt{\lambda_i} \mathbf{q}_i \otimes \bar{p}_i. \quad (27)$$

Denoting by A^* the adjoint operator of A , one can verify that $\mathbf{C} = A^* \circ A$ and $\mathbf{K} = A \circ A^*$, where the latter equation holds when one regards \mathbf{K} as a linear operator from \mathbb{R}^n to \mathbb{R}^n (the symbol \circ denotes composition of operators). Since $\langle A f, \mathbf{x} \rangle = \sum_{i=1}^n x_i \langle \phi_i, f \rangle = \langle f, \sum_{i=1}^n x_i \phi_i \rangle$ for any $f \in \mathcal{H}$ and $\mathbf{x} := [x_1, x_2, \dots, x_n]^\top \in \mathbb{R}^n$, it holds that $A^* \mathbf{x} = \sum_{i=1}^n x_i \phi_i$. It then follows immediately that $\mathbf{K} \mathbf{1} = A(A^* \mathbf{1}) = A(\sum_{i=1}^n \phi_i) = A(0) = \mathbf{0}$, where $\mathbf{1} := [1, 1, \dots, 1]^\top \in \mathbb{R}^n$. This implies that the matrix \mathbf{K} is rank-deficient, and thus $r \leq n - 1$ (cf. Appendix).⁴

In the present nonlinear case, the centered output defined in (10) is given by

$$\tilde{y}_i = \langle \psi, \phi_i \rangle + \tilde{\epsilon}_i, \quad i = 1, 2, \dots, n. \quad (28)$$

Accordingly, each component of $\tilde{\mathbf{g}}$ (see (11)) is given by

$$\tilde{g}_i = \langle \psi, \phi_i \rangle = \langle P_{\mathcal{M}_\phi}(\psi), \phi_i \rangle, \quad i = 1, 2, \dots, n, \quad (29)$$

where the second equality is verified with $\phi_i \in \mathcal{M}_\phi$. We are now ready to show the following theorem.

Theorem 1 Let $P_{\mathcal{M}_\phi}(\psi) = \sum_{i=1}^n \gamma_i \phi_i = A^* \boldsymbol{\gamma} \in \mathcal{M}_\phi$, where $\boldsymbol{\gamma} := [\gamma_1, \gamma_2, \dots, \gamma_n]^\top \in \mathbb{R}^n$. Then, the contribution of \mathbf{q}_i to the essential relevant information vector $\tilde{\mathbf{g}}$ is given by

$$|\tilde{\mathbf{g}}^\top \mathbf{q}_i| = \lambda_i |\boldsymbol{\gamma}^\top \mathbf{q}_i|. \quad (30)$$

Proof. Substituting $P_{\mathcal{M}_\phi}(\psi) = A^* \boldsymbol{\gamma}$ into (29), we obtain $\tilde{g}_i = \langle A^* \boldsymbol{\gamma}, \phi_i \rangle$ and thus

$$\tilde{\mathbf{g}} = A \circ A^*(\boldsymbol{\gamma}) = \mathbf{K} \boldsymbol{\gamma}. \quad (31)$$

By (25) and (31), one can verify the assertion. \square

Theorem 1 states that the contribution of \mathbf{q}_i to the essential relevant information vector $\tilde{\mathbf{g}}$ decays at the same rate as λ_i , as in the case of the uncentered kernel PCA.

⁴Let $p = \sum_{i=1}^n q_i \phi_i \in \mathcal{H}$, $q_i \in \mathbb{R}$, $i = 1, 2, \dots, n$. Then, $\lambda p = \mathbf{C} p$, $\lambda \in \mathbb{R}$, if and only if $\lambda \mathbf{K} \mathbf{q} = \mathbf{K}^2 \mathbf{q}$, where $\mathbf{q} := [q_1, q_2, \dots, q_n]^\top$. Therefore, any solution \mathbf{q} of the dual eigenvalue problem $\lambda \mathbf{q} = \mathbf{K} \mathbf{q}$ gives p such that $\lambda p = \mathbf{C} p$ (see [2]). However, \mathbf{q} corresponding to $\lambda = 0$ may give $p = 0$, as is the case of $\mathbf{q} = \mathbf{1}$. In such a case, the solution \mathbf{q} does not give a solution to the primal eigenvalue problem, since $p = 0$ is not an eigenvector by definition.

3.4. Discussion

Relation to prior work: We discuss the relation between (20) and the result in [3, Theorem 1]. Referring to (20), one can see that the right hand side contains a sole term which decays at the same rate as the eigenvalue λ_i . The truncation error contained in the result in [3] does not appear in our analysis due to the simple derivation without using the integral operator. The absence of the truncation error results in the exact equality appearing in (20) (unlike the result in [3] which shows an upper bound of $|\mathbf{g}^\top \mathbf{q}_i|$ essentially). One may concern a situation when the values of $|\alpha_i^\top \mathbf{q}_i|$ are small for the leading eigenvectors compared to those for minor eigenvectors. This situation however corresponds to the case where many of the sample vectors \mathbf{x}_i are filtered out by the nonlinear function ψ and accordingly many of the outputs y_i contain little signals. In such a situation, the learning problem itself is very challenging, and it would be advised to attain some extra sample vectors that are not filtered out by ψ . It should be mentioned that the decay rate of λ_i depends on the choice of the kernel (see [3]).

4. NUMERICAL EXAMPLE

We conduct a computer simulation using toy data of dimension $N = 2$ to illustrate performance of the variants of kernel PCA: the uncentered kernel PCA and the centered kernel PCA. Training dataset is generated as follows:

1. specify five different cluster centers $\mathbf{c}_1 := (0, 0)$, $\mathbf{c}_2 := (5, 4)$, $\mathbf{c}_3 := (-3, 2)$, $\mathbf{c}_4 := (1, 4)$, $\mathbf{c}_5 := (10, -5)$, and
2. generate 100 data points randomly around one of the cluster centers from the i. i. d. normal distribution $\mathcal{N}(\mathbf{c}_i, 0.04\mathbf{I})$, where the cluster center is chosen randomly with equal probability for each data point.

Test dataset is generated in exactly the same way as (but independently from) the training data. The noise ϵ obeys the i.i.d. normal distribution $\mathcal{N}(0, 0.01)$. The nonlinear function ψ and the training data are depicted in Figure 2. In kernel PCA, the Gaussian kernel is used for $\sigma := 0.2$.

Figure 3 plots (a) the singular values of the kernel matrix \mathbf{K}_u (Uncentered) and \mathbf{K} (Centered), and (b) the estimation error $\sqrt{(\mathbf{y}_{\text{test}} - \hat{\mathbf{y}})^\top (\mathbf{y}_{\text{test}} - \hat{\mathbf{y}}) / \mathbf{y}_{\text{test}}^\top \mathbf{y}_{\text{test}}}$, where \mathbf{y}_{test} is the output vector for the test data and $\hat{\mathbf{y}}$ is its estimate. One can see that there is no notable difference in the singular-value decay between the two cases. Referring to Figure 3(b), the centered and uncentered kernel PCAs show a similar tendency.

5. CONCLUSION

We studied the contributions of the kernel principal components to the relevant information in nonlinear regression for both centered and uncentered cases. Our simple derivation enabled to deal with the centering operation explicitly and led to the exact analytic expression of the contribution (containing no truncation errors). It turned out that the essential contribution decayed at the same rate as the eigenvalues. The numerical example showed that the centered and uncentered kernel PCAs had a similar tendency in test-error-decay characteristics. The present study will be useful to estimate the relevant dimension in nonlinear regression under the centering operation, as well as improving the performance of kernel adaptive filtering [11–17].

APPENDIX

By definition of \mathcal{C} , the number r of nonzero eigenvalues of \mathcal{C} equals to the dimension of \mathcal{M}_ϕ ($:= \text{span}\{\phi_i\}_{i=1}^n$). We can verify the following lemmas. (All the lemmas shown below apply to linear PCA. We omit the proof of the second lemma).

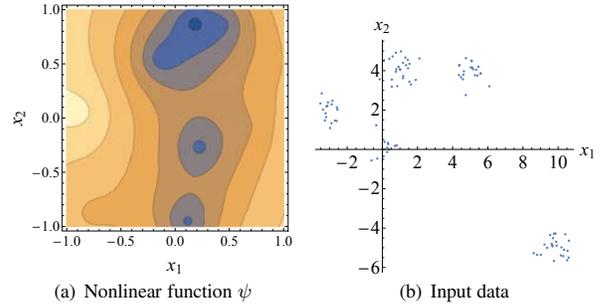
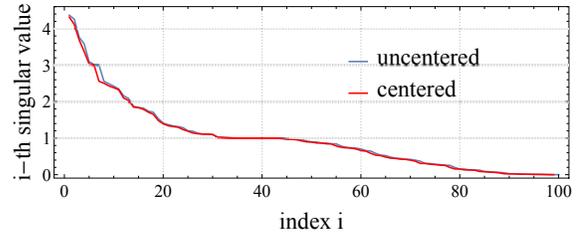
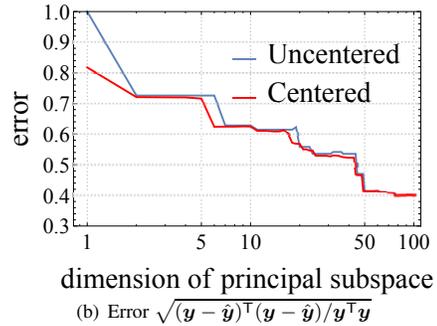


Fig. 2. Simulation data.



(a) Singular-value decay



(b) Error $\sqrt{(\mathbf{y} - \hat{\mathbf{y}})^\top (\mathbf{y} - \hat{\mathbf{y}}) / \mathbf{y}^\top \mathbf{y}}$

Fig. 3. Simulation results.

Lemma 1 Assume that $\bar{\phi} \notin \mathcal{M}_\phi$. Then, the following statements hold.

1. $\mathcal{M} = \mathcal{M}_\phi + \text{span}\{\bar{\phi}\} := \{f + g \mid f \in \mathcal{M}_\phi, g \in \text{span}\{\bar{\phi}\}\}$
2. $\dim(\mathcal{M}_\phi) = \dim(\mathcal{M}) - 1$.

Proof of 1: It is clear that $\Phi(\mathbf{x}_i) = \phi_i + \bar{\phi} \in \mathcal{M}_\phi + \text{span}\{\bar{\phi}\}$, which implies that $\mathcal{M} \subset \mathcal{M}_\phi + \text{span}\{\bar{\phi}\}$. On the other hand, it holds that $\phi_i, \bar{\phi} \in \mathcal{M}$, which implies that $\mathcal{M} \supset \mathcal{M}_\phi + \text{span}\{\bar{\phi}\}$. This completes the proof.

Proof of 2: Clear from Lemma 1.1 and the assumption. \square

Lemma 2 $\bar{\phi} \in \mathcal{M}_\phi$ if and only if $0 \in \text{aff}\{\Phi(\mathbf{x}_i)\}_{i=1}^n$; i.e., there exist $\alpha_i \in \mathbb{R}$, $i = 1, 2, \dots, n$, such that $\sum_{i=1}^n \alpha_i = 1$ and $\sum_{i=1}^n \alpha_i \Phi(\mathbf{x}_i) = 0$. Here, $\text{aff}(\cdot)$ stands for the affine hull.

Lemma 3 If $\{\Phi(\mathbf{x}_i)\}_{i=1}^n$ is linearly independent, then $\bar{\phi} \notin \mathcal{M}_\phi$.

Proof. It is clear that $0 \in \text{aff}\{\Phi(\mathbf{x}_i)\}_{i=1}^n \Leftrightarrow [\sum_{i=1}^n \alpha_i \Phi(\mathbf{x}_i) = 0$ for some $\alpha_i \in \mathbb{R}$ such that $\sum_{i=1}^n \alpha_i = 1] \Rightarrow \{\{\Phi(\mathbf{x}_i)\}_{i=1}^n$ is linearly dependent]. Hence, the linear independence assumption implies $0 \notin \text{aff}\{\Phi(\mathbf{x}_i)\}_{i=1}^n$ and thus $\bar{\phi} \notin \mathcal{M}_\phi$ by Lemma 2. \square

Lemma 3 suggests that, when the Gaussian kernel is employed, $\bar{\phi} \notin \mathcal{M}_\phi$ and hence the statements of Lemma 1 hold true.

6. REFERENCES

- [1] B. Schölkopf, A. Smola, and K.-R. Müller, "Nonlinear component analysis as a kernel eigenvalue problem," *Neural Computation*, vol. 10, pp. 1299–1319, 1998.
- [2] K.-R. Müller, S. Mika, G. Rätsch, K. Tsuda, and B. Schölkopf, "An introduction to kernel-based learning algorithms," *IEEE Trans. Neural Networks*, vol. 12, no. 2, pp. 181–201, Mar. 2001.
- [3] M. L. Braun, J. M. Buhmann, and K.-R. Müller, "On relevant dimensions in kernel feature spaces," *Journal of Machine Learning Research*, vol. 9, pp. 1875–1908, 2008.
- [4] H. van den Dool, *Empirical methods in short-term climate prediction*, Oxford University Press, 2007.
- [5] F. Gwadry, C. Berenstein, J. van Horn, and A. Braun, "Implementation and application of principal component analysis on functional neuroimaging data," Tech. Rep., Institute for Systems Research, 2001.
- [6] I. T. Jolliffe, *Principal Component Analysis*, Springer Series in Statistics. Springer-Verlag, New York, 2nd edition, 2002.
- [7] L. L. Scharf, "The SVD and reduced rank signal processing," *Signal Processing*, vol. 25, pp. 113–133, 1991.
- [8] J. Cadima and I. Jolliffe, "On relationships between uncentred and column-centred principal component analysis," *Pakistan J. Statistics*, vol. 25, no. 4, pp. 473–503, 2009.
- [9] G. Montavon, M. L. Braun, T. Krueger, and K.-R. Müller, "Analyzing local structure in kernel-based learning: Explanation, complexity, and reliability assessment," *IEEE Signal Processing Magazine*, vol. 30, no. 4, pp. 62–74, July 2013.
- [10] D. G. Luenberger, *Optimization by Vector Space Methods*, New York: Wiley, 1969.
- [11] W. Liu, J. Principe, and S. Haykin, *Kernel Adaptive Filtering*, Wiley, New Jersey, 2010.
- [12] J. Kivinen, A. J. Smola, and R. C. Williamson, "Online learning with kernels," *IEEE Trans. Signal Process.*, vol. 52, no. 8, pp. 2165–2176, Aug. 2004.
- [13] Y. Engel, S. Mannor, and R. Meir, "The kernel recursive least-squares algorithm," *IEEE Trans. Signal Process.*, vol. 52, no. 8, pp. 2275–2285, Aug. 2004.
- [14] C. Richard, J. Bermudez, and P. Honeine, "Online prediction of time series data with kernels," *IEEE Trans. Signal Process.*, vol. 57, no. 3, pp. 1058–1067, Mar. 2009.
- [15] K. Slavakis, S. Theodoridis, and I. Yamada, "Adaptive constrained learning in reproducing kernel Hilbert spaces: the robust beamforming case," *IEEE Trans. Signal Process.*, vol. 57, no. 12, pp. 4744–4764, Dec. 2009.
- [16] M. Yukawa, "Multikernel adaptive filtering," *IEEE Trans. Signal Processing*, vol. 60, no. 9, pp. 4672–4682, Sept. 2012.
- [17] M. Yukawa, "Adaptive learning in Cartesian product of reproducing kernel Hilbert spaces," *IEEE Trans. Signal Processing*, vol. 63, no. 22, pp. 6037–6048, Nov. 2015.