

## Kernel Principal Component Regression with Application to Nonlinear Prediction

**Dachapak, Chooleewan**

Department of Electrical and Electronic Systems Engineering, Graduate School of Information Science and Electrical Engineering, Kyushu University : Graduate Student

**Kanae, Shunshoku**

Department of Electrical and Electronic Systems Engineering, Faculty of Information Science and Electrical Engineering, Kyushu University, Kyushu University

**Yang, Zi-Jiang**

Department of Electrical and Electronic Systems Engineering, Faculty of Information Science and Electrical Engineering, Kyushu University, Kyushu University

**Wada, Kiyoshi**

Department of Electrical and Electronic Systems Engineering, Faculty of Information Science and Electrical Engineering, Kyushu University, Kyushu University

<https://doi.org/10.15017/1515812>

---

出版情報 : 九州大学大学院システム情報科学紀要. 8 (1), pp.19-23, 2003-03-26. 九州大学大学院システム情報科学研究所

バージョン :

権利関係 :

## Kernel Principal Component Regression with Application to Nonlinear Prediction

Chooleewan DACHAPAK\*, Shunshoku KANAE\*\*, Zi-Jiang YANG\*\* and Kiyoshi WADA\*\*

(Received December 20, 2002)

**Abstract:** In this study, Kernel Principal Component Analysis (KPCA) is applied as feature selection in a high-dimensional feature space which is nonlinearly related to an input space. By using Mercer Kernels, we can compute principal components in a high dimensional feature space. Then, the extracted features by KPCA method are employed as a new kind of regressors in an ordinary least square regression in the feature space which is Reproducing Kernel Hilbert Space (RKHS). In the experiment, KPCR method was applied to predict the Mackey-Glass time series.

**Keywords:** Kernel Principal Component Analysis, Kernel Principal Component Regression

### 1. Introduction

In the present study, we are interested in the feature selection method as the appropriate preprocessing step for an ordinary least squares regression, since the suitable feature selection can increase the overall performance of algorithms<sup>1)</sup>. Kernel Principal Component Analysis is based on the computation of the standard linear PCA in the feature space mapped from input space by some nonlinear function<sup>2)</sup>. We have applied KPCA method for feature selection in a high dimensional or possibly infinite feature  $F$  (with dimension  $M \leq \infty$ ). This allows us to extract the nonlinear principal components up to the number of data points  $n$  ( $n \leq M$ ). We compute a dot product in the feature space by means of kernel functions in input space, i.e.  $k(\mathbf{x}, \mathbf{y}) = (\Phi(\mathbf{x}) \cdot \Phi(\mathbf{y}))$ , and we never need the mapped patterns  $\Phi(\mathbf{x})$  explicitly. It is often the case that a small number of Kernel Principal Component is sufficient to account for most of the structure in the data. Also, the Kernel principal components are uncorrelated each other, and it is possible that KPCA can decrease the effect of noise<sup>2)</sup>. Then, the features extracted by KPCA method are employed as a regressors of the ordinary least squares regression. This means it performs a linear regression in the feature space  $F$ .

### 2. Kernel Principal Component Analysis

First, mapping the data nonlinearly into the fea-

ture space  $F$  by the nonlinear function as following,

$$\Phi: \mathbf{R}^N \rightarrow F, \quad \mathbf{x} \mapsto \mathbf{X}. \quad (1)$$

Note that the feature space  $F$  could have an arbitrarily large, possibly infinite, dimensionality<sup>3)</sup>. Assume that our data mapped into feature space is centered,  $\sum_{k=1}^n \Phi(\mathbf{x}_k) = 0$ . Then, in  $F$  space, a covariance matrix can be expressed as

$$\bar{C} = \frac{1}{n} \sum_{j=1}^n \Phi(\mathbf{x}_j) \Phi(\mathbf{x}_j)^T. \quad (2)$$

Define  $(n \times M)$  matrix  $\Phi$ ,

$$\Phi^T \triangleq [\Phi(\mathbf{x}_1) \ \Phi(\mathbf{x}_2) \ \dots \ \Phi(\mathbf{x}_n)]. \quad (3)$$

So  $(M \times M)$  matrix  $\bar{C}$  can be rewritten as

$$\bar{C} = \frac{1}{n} \Phi^T \Phi. \quad (4)$$

Let  $(n \times n)$  matrix  $K$  denote  $K = \Phi \Phi^T$ . Noting  $n\bar{C}$  and  $K$ 's nonzero eigenvalues are the same. Let  $\lambda_k, \mathbf{V}^k$  be  $\bar{C}$ 's  $k^{th}$  eigenvalues and eigenvectors respectively, where  $k = 1, 2, \dots, M$ .

$$n\lambda \mathbf{V}^k = n\bar{C} \mathbf{V}^k = \Phi^T \Phi \mathbf{V}^k. \quad (5)$$

Pre-multiply  $\Phi$ ,

$$n\lambda \Phi \mathbf{V}^k = \Phi \Phi^T \Phi \mathbf{V}^k = K \Phi \mathbf{V}^k. \quad (6)$$

Let

$$\gamma^k = \Phi \mathbf{V}^k, \quad (7)$$

then,

$$K \gamma^k = \tilde{\lambda} \gamma^k. \quad (8)$$

Note that  $\tilde{\lambda} = n\lambda$ . Moreover from,

$$\Phi^T \gamma = \Phi^T \Phi \mathbf{V} = \tilde{\lambda} \mathbf{V}, \quad (9)$$

We obtain  $\mathbf{V}^k = \frac{1}{\tilde{\lambda}} \Phi^T \gamma^k$ . Since it is required that

\* Department of Electrical and Electronic Systems Engineering, Graduate Student

\*\* Department of Electrical and Electronic Systems Engineering

eigenvector  $\mathbf{V}$  in  $F$  space must be normalized. We obtain the normalized eigenvector as following form,

$$\mathbf{V}^k = \frac{1}{\sqrt{\lambda}} \Phi^T \gamma^k. \quad (10)$$

Let  $\alpha^k = \frac{1}{\sqrt{\lambda}} \gamma^k$ . Then,

$$\mathbf{V}^k = \Phi^T \alpha^k. \quad (11)$$

Finally, project  $\Phi(\mathbf{x})$  onto the  $k^{th}$  nonlinear principal component by the dot product between  $\mathbf{V}^k$  and  $\Phi(\mathbf{x})$  as below,

$$\mathbf{V}^k \cdot \Phi(\mathbf{x}) = \alpha^{kT} \Phi \Phi(\mathbf{x}). \quad (12)$$

The above  $1^{st}, \dots, M^{th}$  nonlinear principal components,  $\Phi \mathbf{V}$ , are employed as a new regressors for the regression model in the feature space. We used only first  $p$  components in order to reduce the noise effect since the components whose variances are very small can be assumed as the variance from noise.

### 3. Kernel Principal Component Regression

Consider the standard regression model in the feature space

$$\mathbf{y} = \Phi \boldsymbol{\xi} + \boldsymbol{\epsilon}, \quad (13)$$

where  $\mathbf{y}$  is a vector of  $n$  observations,  $\Phi$  is an  $(n \times M)$  matrix of regressors whose  $i^{th}$  row is the vector  $\Phi(\mathbf{x}_i)$  of the mapped  $\mathbf{x}_i$  observation into  $M \leq \infty$  dimensional feature space  $F$ ,  $\boldsymbol{\xi}$  is a vector of regression coefficients and  $\boldsymbol{\epsilon}$  is a vector of error terms. We assume that regressors are centered. Now, let  $(M \times M)$  matrix  $V$  denote

$$V = [\mathbf{V}^1 \ \mathbf{V}^2 \ \dots \ \mathbf{V}^M] \quad (14)$$

then,

$$\Phi^T \Phi V = V \Lambda. \quad (15)$$

Noting  $\Lambda$  is a diagonal matrix of  $\{\tilde{\lambda}_1, \tilde{\lambda}_2, \dots, \tilde{\lambda}_M\}$ .  $V$  is an orthogonal matrix so the linear regression model can be expressed in the term of new regressors as

$$\mathbf{y} = \Phi V V^T \boldsymbol{\xi} + \boldsymbol{\epsilon} = \mathbf{B} \mathbf{w} + \boldsymbol{\epsilon}, \quad (16)$$

where  $\mathbf{B} = \Phi V$ ,  $\mathbf{w} = V^T \boldsymbol{\xi}$ .

Here we can see that the new regressors  $B$  are transformed in the term of Kernel Principal Components which are the extracted features we are interested in. And

$$\mathbf{B}^T \mathbf{B} = V^T \Phi^T \Phi V = V^T V \Lambda = \Lambda \quad (17)$$

So the least squares estimator of coefficient  $\mathbf{w}$  can be expressed as

$$\hat{\mathbf{w}} = (\mathbf{B}^T \mathbf{B})^{-1} \mathbf{B}^T \mathbf{y} = \Lambda^{-1} \mathbf{B}^T \mathbf{y}. \quad (18)$$

From above, we can see that it is difficult to find  $\mathbf{B}$  directly. The easier way to find  $\mathbf{B}$  by using kernel function will be shown.

Let  $y_k$  be the  $k^{th}$  component of  $\mathbf{y}$  and  $\epsilon_k$  be the  $k^{th}$  component of  $\boldsymbol{\epsilon}$ , from the linear regression model, we obtain

$$y_k = \Phi(\mathbf{x}_k)^T \boldsymbol{\xi} + \epsilon_k = \Phi(\mathbf{x}_k)^T V \mathbf{w} + \epsilon_k. \quad (19)$$

Here, consider  $V^T \Phi(\mathbf{x}_k)$ , from the definition of  $V$ , we can see that

$$V^T \Phi(\mathbf{x}_k) = \begin{bmatrix} \mathbf{V}^{1T} \Phi(\mathbf{x}_k) \\ \mathbf{V}^{2T} \Phi(\mathbf{x}_k) \\ \vdots \\ \mathbf{V}^{M^T} \Phi(\mathbf{x}_k) \end{bmatrix} = \begin{bmatrix} \beta_1(\mathbf{x}_k) \\ \beta_2(\mathbf{x}_k) \\ \vdots \\ \beta_M(\mathbf{x}_k) \end{bmatrix}. \quad (20)$$

Note that  $\beta_i(\mathbf{x}_k) = \mathbf{V}^{iT} \Phi(\mathbf{x}_k)$ , because  $\mathbf{V}^i$  and  $\alpha^i$  have the relationship as

$$\mathbf{V}^i = \Phi^T \alpha^i \quad (21)$$

so  $\beta_i(\mathbf{x}_k)$  can be represented by

$$\begin{aligned} \beta_i(\mathbf{x}_k) &= \alpha^{iT} \Phi \Phi(\mathbf{x}_k) \\ &= \alpha^{iT} \begin{bmatrix} \Phi(\mathbf{x}_1)^T \Phi(\mathbf{x}_k) \\ \Phi(\mathbf{x}_2)^T \Phi(\mathbf{x}_k) \\ \vdots \\ \Phi(\mathbf{x}_n)^T \Phi(\mathbf{x}_k) \end{bmatrix} \\ &= \alpha^{iT} \begin{bmatrix} k(\mathbf{x}_1, \mathbf{x}_k) \\ k(\mathbf{x}_2, \mathbf{x}_k) \\ \vdots \\ k(\mathbf{x}_n, \mathbf{x}_k) \end{bmatrix}. \end{aligned} \quad (22)$$

We compute the dot products  $\Phi(\mathbf{x}_i)^T \Phi(\mathbf{x}_k)$  in the feature space by means of kernel functions in input space  $k(\mathbf{x}_i, \mathbf{x}_k) = \Phi(\mathbf{x}_i)^T \Phi(\mathbf{x}_k)$ . Moreover, we do not need to map  $\Phi(\mathbf{x})$  explicitly throughout the study. Then, define  $\mathbf{k}(\mathbf{x}_k)$  as following,

$$\mathbf{k}(\mathbf{x}_k) = \begin{bmatrix} k(\mathbf{x}_1, \mathbf{x}_k) \\ k(\mathbf{x}_2, \mathbf{x}_k) \\ \vdots \\ k(\mathbf{x}_n, \mathbf{x}_k) \end{bmatrix} \quad (23)$$

$\beta_i(\mathbf{x}_k)$  becomes

$$\beta_i(\mathbf{x}_k) = \alpha^{iT} \mathbf{k}(\mathbf{x}_k). \quad (24)$$

So

$$V^T \Phi(\mathbf{x}_k) = \begin{bmatrix} \beta_1(\mathbf{x}_k) \\ \beta_2(\mathbf{x}_k) \\ \vdots \\ \beta_n(\mathbf{x}_k) \end{bmatrix} = \begin{bmatrix} \boldsymbol{\alpha}^1{}^T \\ \boldsymbol{\alpha}^2{}^T \\ \vdots \\ \boldsymbol{\alpha}^n{}^T \end{bmatrix} \mathbf{k}(\mathbf{x}_k) = A^T \mathbf{k}(\mathbf{x}_k). \quad (25)$$

Note that

$$A = [\boldsymbol{\alpha}^1 \ \boldsymbol{\alpha}^2 \ \dots \ \boldsymbol{\alpha}^n]. \quad (26)$$

We can rewrite (19) as

$$\mathbf{y}_k = \mathbf{k}(\mathbf{x}_k)^T A \mathbf{w} + \epsilon_k = \mathbf{k}(\mathbf{x}_k)^T \boldsymbol{\theta} + \epsilon_k \quad (27)$$

where

$$\boldsymbol{\theta} = A \mathbf{w} = A V^T \boldsymbol{\xi}. \quad (28)$$

From the definition of  $B$ ,

$$B = \Phi V \quad (29)$$

$$B = \begin{bmatrix} \Phi(\mathbf{x}_1)^T \\ \Phi(\mathbf{x}_2)^T \\ \vdots \\ \Phi(\mathbf{x}_n)^T \end{bmatrix} V = \begin{bmatrix} \Phi(\mathbf{x}_1)^T V \\ \Phi(\mathbf{x}_2)^T V \\ \vdots \\ \Phi(\mathbf{x}_n)^T V \end{bmatrix}$$

$$B = \begin{bmatrix} \mathbf{k}(\mathbf{x}_1)^T \\ \mathbf{k}(\mathbf{x}_2)^T \\ \vdots \\ \mathbf{k}(\mathbf{x}_n)^T \end{bmatrix} A, \quad (30)$$

where,

$$\begin{bmatrix} \mathbf{k}(\mathbf{x}_1)^T \\ \mathbf{k}(\mathbf{x}_2)^T \\ \vdots \\ \mathbf{k}(\mathbf{x}_n)^T \end{bmatrix} = \begin{bmatrix} k(\mathbf{x}_1, \mathbf{x}_1) & k(\mathbf{x}_2, \mathbf{x}_1) & \dots & k(\mathbf{x}_n, \mathbf{x}_1) \\ k(\mathbf{x}_1, \mathbf{x}_2) & k(\mathbf{x}_2, \mathbf{x}_2) & \dots & k(\mathbf{x}_n, \mathbf{x}_2) \\ \vdots & \vdots & \ddots & \vdots \\ k(\mathbf{x}_1, \mathbf{x}_n) & k(\mathbf{x}_2, \mathbf{x}_n) & \dots & k(\mathbf{x}_n, \mathbf{x}_n) \end{bmatrix}. \quad (31)$$

We can rewrite  $B$  as

$$B = K A. \quad (32)$$

So we can formulate the KPCR model as following expression,

$$\mathbf{y} = B \mathbf{w} + \boldsymbol{\epsilon} = K A \mathbf{w} + \boldsymbol{\epsilon} = K \boldsymbol{\theta} + \boldsymbol{\epsilon}. \quad (33)$$

We prefer (32) to (29) since we can compute  $B$  by means of kernel functions in input space without mapping  $\Phi(\mathbf{x})$  explicitly. As we mentioned above, we employed the first  $p$  nonlinear principal components to create a KPCR model in order to decrease the effect of noise and eliminate large variances of the estimate due to multicollinearities. In this case

the first  $p$  nonlinear principal components mean the first  $p$  column vectors of  $A$ .

#### 4. Experiment

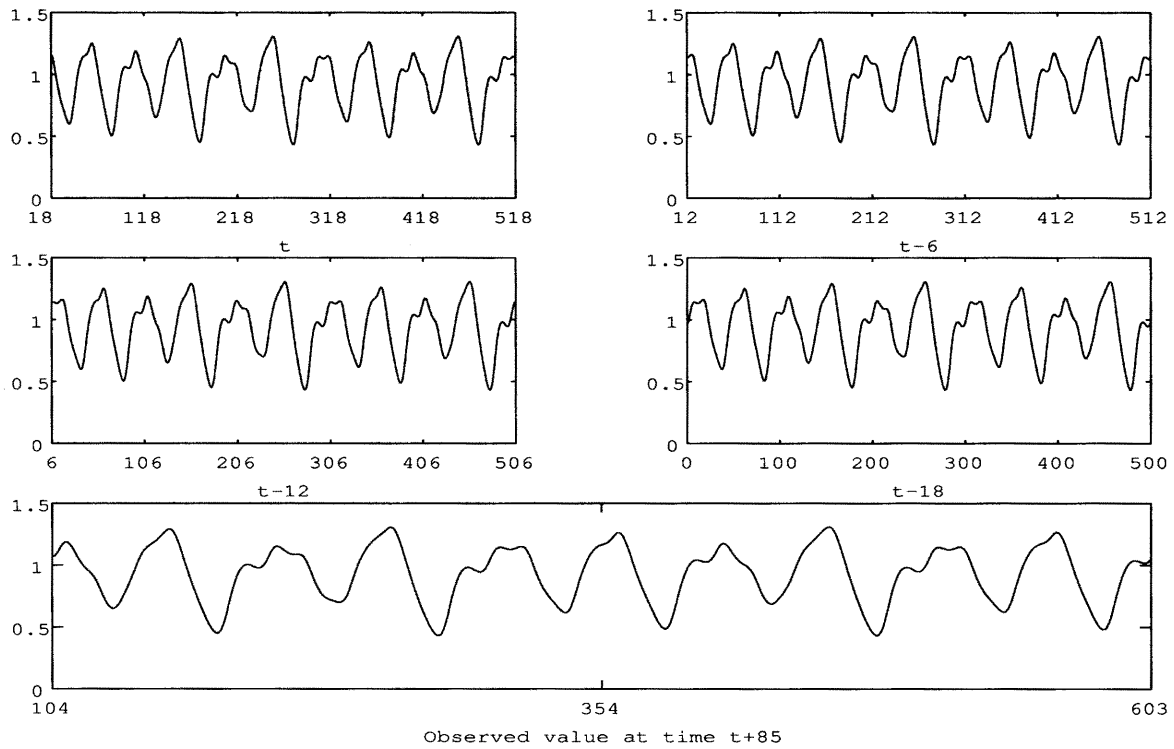
In this present study, we employed Gaussian kernel function,  $k(\mathbf{x}, \mathbf{y}) = e^{-\left(\frac{\|\mathbf{x}-\mathbf{y}\|^2}{L}\right)}$ , where  $L$  determines the width of the Gaussian function. We utilized the noisy Chaotic Mackey-Glass Time-Series applying to KPCR method which the regressors are extracted by Kernel PCA. This KPCR method was trained to predict the value at time  $t + 85$  from inputs at time  $t, t - 6, t - 12, t - 18$ . The training data partitions were constructed over 3000 training samples in steps of 500 samples. So we got six partitions. The simulations were repeated for the width  $L$  from range  $(0.2\hat{\sigma}, 20\hat{\sigma})$ , where  $\hat{\sigma}$  is the variance of the overall clean training set, using the step size 0.01. A fixed test set of size 500 data points was employed through out the experiments. The performance of the regression models to predict a clean Mackey-Glass time series was evaluated in terms of the normalized root mean squared error (NRMSE). The best result on the test set averaged over all individual runs are summarized in **Table 1**.

**Table 1** The comparison of the approximation errors (NRMSE) for prediction of Mackey-Glass training set.  $N/S$  represents the ratio between the standard deviation of the added Gaussian noise and the underlying time-series.

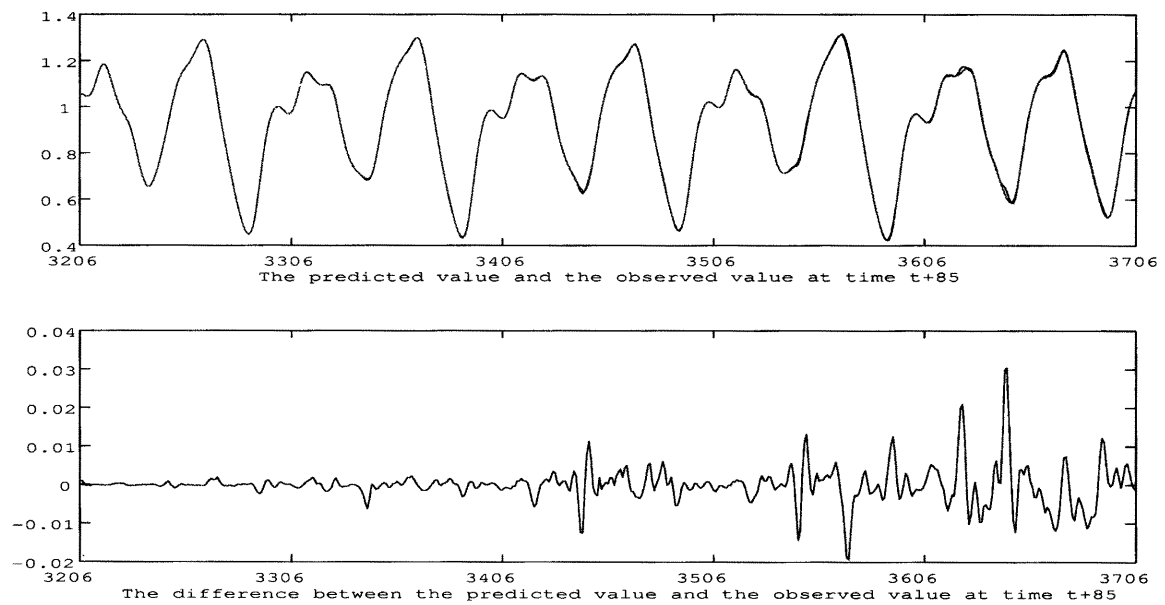
Method	$N/S = 0.0$	$N/S = 0.11$	$N/S = 0.22$
KPCR	0.038	0.307	0.443

#### 5. Conclusions

The Kernel PCA method has been employed for the feature extraction and noise reduction in the preprocessing step for regression problem. With KPCR method, there are some problems occurred when we remove some Kernel principal components with small eigenvalues in order to reduce the noise present because those Kernel principal components may have a significant contribution for the prediction.



**Fig.1** Training data set: 500 samples from 3000 samples of Mackey-Glass time series. First 4 figures show the inputs at time  $t, t - 6, t - 12, t - 18$  from  $t = 19$  to 518. They are employed as the data set  $x$  which are trained by KPCR method to predict the value at time  $t + 85$  (The bottom shows the observed value at time  $t + 85$ ).



**Fig.2** Testing data set: the above shows the predicted value and the observed value at time  $t + 85$  ( $t = 3122$  to 3621). And the bottom shows the different value between the predicted value and the observed value.

## References

- 1) R. Rosipal, M. Girolami, L.J. Trejo, and A. Cichocki. Kernel PCA for Feature Extraction and De-noising in Non-linear Regression. *Neural Computing & Applications*, 10(3), 2001.
  - 2) B. Schölkopf, C.J.C. Burges, and A.J. Smola. Advances in Kernel Methods Support Vector Learning. *The MIT Press, Cambridge, Massachusetts, London England*, 1998.
  - 3) B. Schölkopf, A.J. Smola, and K.R. Muller. Nonlinear Component Analysis as a Kernel Eigenvalue Problem. *Neural Computation*, 10:1299–1319, 1998.
  - 4) R. Rosipal and L.J. Trejo. Kernel Partial Least Squares Regression in Reproducing Kernel Hilbert Space. *Journal of Machine Learning Research*, pages 97–123, December 2(2001).
  - 5) B. Schölkopf. Support Vector Learning. *PhD thesis, Universitat Berlin, Berlin, Germany*, 1997.
  - 6) M.G. Genton. Classes of Kernels for Machine Learning : A Statistics Perspective. *Journal of Machine Learning Research*, pages 299–312, 2(2001).
  - 7) I.T. Jolliffe. Principal Component Analysis. *Springer-Verlag, New York*, 1986.
  - 8) H.D. Vinod and A. Ullah. Recent Advances in Regression Methods. *Marcel Dekker, INC.*, 1981.
-