



Seppo Pulkkinen

# A Probabilistic Approach to Nonlinear Principal Component Analysis With Ap- plications

TURKU CENTRE *for* COMPUTER SCIENCE

TUCS Technical Report  
No 1091, November 2013





# A Probabilistic Approach to Nonlinear Principal Component Analysis With Applications

Seppo Pulkkinen

University of Turku, Department of Mathematics and Statistics

FI-20014 Turku, Finland

`seppo.pulkkinen@utu.fi`

TUCS Technical Report

No 1091, November 2013

## **Abstract**

Principal component analysis (PCA) is a well-established tool for identifying the main sources of variation in multivariate data. However, as a linear method it cannot describe complex nonlinear structures. To overcome this limitation, a novel nonlinear generalization of PCA is developed in this paper. The method obtains the nonlinear principal components from ridges of the underlying density of the data. The density is estimated by using Gaussian kernels. Projection onto a ridge of such a density estimate is formulated as a solution to a differential equation, and a predictor-corrector method is developed for this purpose. The method is further extended to time series data by applying it to the phase space representation of the time series. This extension can be viewed as a nonlinear generalization of singular spectrum analysis (SSA). Ability of the nonlinear PCA to capture complex nonlinear shapes and its SSA-based extension to identify periodic patterns from time series are demonstrated on climate data.

**Keywords:** principal component analysis; singular spectrum analysis; nonlinear; kernel density; nonparametric methods; ridge; predictor-corrector method

**TUCS Laboratory**  
Turku Optimization Group (TOpGroup)

# 1 Introduction

In practical applications, one is often dealing with high-dimensional data that is confined to some low-dimensional subspace. Since its introduction by Pearson [28], *principal component analysis* (PCA, e.g. [18]) has become a ubiquitous tool for identifying such subspaces. The method uses an orthogonal transformation to separate the directions of maximal variance. PCA and its variants have appeared in various contexts such as *empirical orthogonal functions* (EOF) in climate analysis (e.g. [39]), *proper orthogonal decomposition* (POD) in fluid mechanics (e.g. [3]) and the *Karhunen-Loève transform* (KLT) in the theory of stochastic processes (e.g. [21]).

However, as a linear method, PCA is insufficient for describing complex nonlinear data. Several nonlinear extensions have been developed to overcome this limitation. The most prominent of these are the neural network-based nonlinear PCA (NLPCA, e.g. [15, 20, 24, 35]) and kernel PCA (KPCA, e.g. [36]). These methods, however, have shortcomings. NLPCA requires a large number of user-supplied parameters that need to be carefully tuned for the application at hand. Furthermore, the transformation of the input data into the high-dimensional kernel space in KPCA incurs a significant computational cost. A careful choice of kernel function is also needed when using KPCA.

Some variants of PCA, where the principal components are obtained by restricting the analysis to local neighbourhoods of the data points, have been developed (e.g. [8, 9, 19]). However, this approach leads to the problem of determining a global coordinate system. A well-known approach to this problem is local tangent space alignment (LTSA, [42]) that determines a coordinate system by solving an eigenvalue problem constructed from the local principal component coordinates. However, this method and other neighbourhood-based methods are in general sensitive to noise and the choice of the neighbourhoods.

The contribution of this paper is the development of *kernel density principal component analysis* (KDPCA). The proposed method builds on the idea of using *ridges* of the underlying density of the data to estimate nonlinear structures [27]. This idea has later been refined in [29] and [30]. In the proposed approach, the ridges are interpreted as nonlinear counterparts of principal component hyperplanes. The density is estimated by using Gaussian *kernels*.

In the linear PCA, principal component *scores* (i.e. coordinates) of a given sample point are obtained as projections along principal component axes. Generalizing the concept of a principal component axis, the projections in KDPCA are done along curvilinear trajectories onto ridges of a Gaussian kernel density estimate. Based on the theory of ridge sets, it is shown that such projections can be done in a well-defined coordinate system. A projection trajectory is formulated as a solution to a differential equation, and a predictor-corrector algorithm is developed for tracing its solution curve.

A strategy for choosing the kernel *bandwidth* is critical for the practical ap-

plicability of KDPCA. To this end, it will be shown that the nonlinear principal components converge to the linear ones when the kernel bandwidth approaches infinity. Consequently, the robustness of the linear PCA is always attained by choosing a sufficiently large bandwidth, but there is a tradeoff between robustness and ability to describe nonlinear structure.

Finally, KDPCA is extended to time series analysis. In analogy with the well-known *singular spectrum analysis* (SSA, e.g. [11, 38]), it is applied to the *phase space* representation of the time series. This approach addresses the main shortcoming of the linear SSA. That is, being based on the linear PCA, it cannot separate different components of a time series when its trajectory in the phase space forms a closed loop. This is the case for quasiperiodic (i.e. approximately periodic) time series that form an important special class appearing in many applications. Examples include climate analysis (e.g. [15, 16]) and medical applications such as electrocardiography and electroencephalography (e.g. [31]).

The remaining of this paper is organized as follows. In Section 2 we recall the linear PCA. Section 3 is devoted to development of KDPCA, and in Section 4 it is extended to time series data. Test results on a simulated climate data set and an atmospheric time series are given in Section 5. The computational complexity of KDPCA is also analyzed and a comparison with related methods is given. Finally, Section 6 concludes this paper. The more involved proofs are deferred to Appendix A.

## 2 The linear PCA

As the proposed method is a generalization of the linear PCA (e.g. [18]), we briefly recall the theoretical background of this method in this section.

The linear PCA attempts to capture the variability of a given data

$$\mathbf{Y} = [\mathbf{y}_1 \quad \mathbf{y}_2 \quad \cdots \quad \mathbf{y}_n]^T \in \mathbb{R}^{n \times d}$$

by transforming the data into a new coordinate system via an orthogonal transformation. In the new coordinate system, the axes point along directions of maximal variance.

For the formulation of PCA, we denote the mean-centered samples by

$$\tilde{\mathbf{y}}_i = \mathbf{y}_i - \hat{\boldsymbol{\mu}}, \quad \text{where } \hat{\boldsymbol{\mu}} = \frac{1}{n} \sum_{i=1}^n \mathbf{y}_i. \quad (1)$$

Assume that the mean-centered samples  $\tilde{\mathbf{y}}_i$  are transformed into an  $m$ -dimensional space via the mapping

$$\boldsymbol{\theta}_i(\mathbf{A}) = \mathbf{A}^T \tilde{\mathbf{y}}_i,$$

where  $\mathbf{A}$  is a  $d \times m$  matrix with  $0 < m < d$  and with orthonormal columns. Conversely, for the given coordinates  $\boldsymbol{\theta}_i$  in the  $m$ -dimensional space, the corresponding *reconstruction* (i.e. projection onto the hyperplane spanned by the  $m$

first principal components) of  $\mathbf{y}_i$  in the input space is obtained as

$$\hat{\mathbf{y}}_i(\mathbf{A}) = \hat{\boldsymbol{\mu}} + \mathbf{A}\boldsymbol{\theta}_i. \quad (2)$$

With the above definitions, it can be shown that finding the matrix  $\mathbf{A}$  that minimizes the reconstruction error is equivalent to maximizing the variance in the transformed coordinate system [18]. That is,

$$\min_{\mathbf{A} \in O(d,m)} \sum_{i=1}^n \|\hat{\mathbf{y}}_i(\mathbf{A}) - \hat{\boldsymbol{\mu}} - \tilde{\mathbf{y}}_i\|^2 = \max_{\mathbf{A} \in O(d,m)} \sum_{i=1}^n \|\boldsymbol{\theta}_i(\mathbf{A})\|^2,$$

where  $O(d, m)$  denotes the set of  $d \times m$  matrices having orthonormal columns. Furthermore, any  $i$ -th principal component corresponds to the direction of the  $i$ -th largest variance, and these directions form an orthogonal set.

The solution to the above optimization problems is the matrix

$$\mathbf{V}_m = [\mathbf{v}_1 \quad \mathbf{v}_2 \quad \cdots \quad \mathbf{v}_m]$$

where the column vectors  $\mathbf{v}_i$  are the (normalized) eigenvectors of the  $d \times d$  sample covariance matrix

$$\hat{\boldsymbol{\Sigma}}_{\mathbf{Y}} = \frac{1}{n-1} \sum_{i=1}^n (\mathbf{y}_i - \hat{\boldsymbol{\mu}})(\mathbf{y}_i - \hat{\boldsymbol{\mu}})^T \quad (3)$$

corresponding to the  $m$  largest eigenvalues. Thus, projection of the mean-centered sample set  $\tilde{\mathbf{Y}}$  onto the  $m$ -dimensional subspace corresponding to the directions of largest variance is given by

$$\boldsymbol{\Theta} = \mathbf{V}_m^T \tilde{\mathbf{Y}}. \quad (4)$$

In statistical literature, the coordinates  $\boldsymbol{\Theta} \in \mathbb{R}^{d \times m}$  obtained in this way are called principal component *scores* (e.g. [18]).

### 3 Nonlinear kernel density PCA

In this section we develop the kernel density principal component analysis (KD-PCA). The method is based on estimation of the underlying density of the data by Gaussian kernels. It is shown that the nonlinear principal component scores of given sample points can be obtained one by one by successively projecting them onto ridges of the density estimate. The projection curves are defined as a solution to a differential equation, and predictor-corrector method is developed for this purpose.

### 3.1 Ridge definition

We adapt the definition of a ridge set from [30]. An  $r$ -dimensional ridge point of a probability density is a local maximum in a subspace spanned by a subset of the eigenvectors of its Hessian matrix. These eigenvectors correspond to the  $d - r$  algebraically smallest eigenvalues. The one-dimensional ridge set (i.e. ridge curve) of the density of a point set is illustrated in Figure 1.

**Definition 3.1.** A point  $\mathbf{x} \in \mathbb{R}^d$  belongs to the  $r$ -dimensional ridge set  $\mathcal{R}_p^r$ , where  $0 \leq r < d$ , of a twice differentiable probability density  $p : \mathbb{R}^d \rightarrow \mathbb{R}$  if

$$\nabla p(\mathbf{x})^T \mathbf{v}_i(\mathbf{x}) = 0, \quad i > r, \quad (5a)$$

$$\lambda_{r+1}(\mathbf{x}) < 0, \quad (5b)$$

$$\lambda_1(\mathbf{x}) > \lambda_2(\mathbf{x}) > \dots > \lambda_{r+1}(\mathbf{x}), \quad \text{if } r > 0, \quad (5c)$$

where  $\lambda_1(\mathbf{x}) \geq \lambda_2(\mathbf{x}) \geq \dots \geq \lambda_d(\mathbf{x})$  and  $\{\mathbf{v}_i(\mathbf{x})\}_{i=1}^d$  denote the eigenvalues and the corresponding eigenvectors of  $\nabla^2 p(\mathbf{x})$ , respectively.

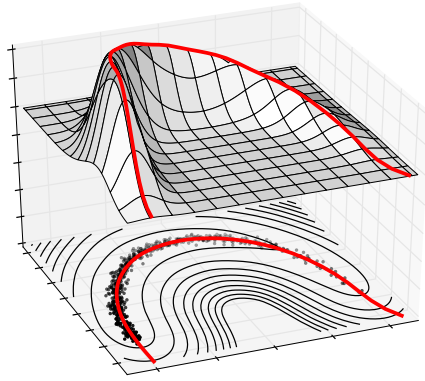


Figure 1: Ridge curve of the density of a point set that is distributed around a curve.

The following result shows a connection between the ridge set and the linear principal components when the underlying density of the data is normal. This result follows trivially from the following lemma (see [27]) and the fact that the logarithm of a normal density with mean  $\boldsymbol{\mu}$  and covariance  $\boldsymbol{\Sigma}$  is a quadratic function whose gradient and Hessian are

$$\nabla \log p(\mathbf{x}) = -\boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu}) \quad \text{and} \quad \nabla^2 \log p = -\boldsymbol{\Sigma}^{-1},$$

respectively.

**Lemma 3.1.** If  $p : \mathbb{R}^d \rightarrow \mathbb{R}$  is twice differentiable, then  $\mathcal{R}_{\log p}^r = \mathcal{R}_p^r$  for all  $r = 0, 1, 2, \dots, d - 1$ .



**Proposition 3.1.** *Let  $p : \mathbb{R}^d \rightarrow \mathbb{R}$  be a  $d$ -variate normal density with mean  $\boldsymbol{\mu}$  and positive definite covariance matrix  $\boldsymbol{\Sigma}$ . Denote the eigenvalues of  $\boldsymbol{\Sigma}$  by  $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_d$  and the corresponding eigenvectors by  $\{\mathbf{v}_i\}_{i=1}^d$ . Then for any  $0 \leq r < d$  such that  $\lambda_1 > \lambda_2 > \dots > \lambda_{r+1}$  we have*

$$\mathcal{R}_p^r = \begin{cases} \{\boldsymbol{\mu}\}, & r = 0, \\ \{\boldsymbol{\mu}\} + \text{span}(\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_r), & r = 1, 2, \dots, d-1. \end{cases}$$

Proposition 3.1 suggests an approach for estimating the principal component scores  $\boldsymbol{\theta}$  of a given point having an underlying density  $p$ . The idea is to project the point onto  $\mathcal{R}_{\log p}^m$  in the subspace spanned by the eigenvectors  $\{\mathbf{v}_i\}_{i=m+1}^d$  of  $\nabla^2 \log p$  for some  $0 < m \leq d$  and then obtain projection coordinates along the first  $m$  eigenvectors. The remaining  $d - m$  components, that are interpreted as noise, are discarded. The point in  $\mathcal{R}_{\log p}^0$ , that is the maximum of  $\log p$ , is chosen as the origin of the coordinate system. This idea will be generalized to the nonlinear case in Section 3.3.

## 3.2 Density estimation and choice of bandwidth

Here we use Gaussian kernels to estimate the density from the given data. In what follows, we establish a connection between ridge sets of this density and the linear principal components. Furthermore, we show that a linear PCA hyperplane is obtained as a special case of such a ridge set.

**Definition 3.2.** *The Gaussian kernel density estimate  $\hat{p}_h$  obtained by drawing a set of samples  $\mathbf{Y} = \{\mathbf{y}_i\}_{i=1}^n \subset \mathbb{R}^d$  from a probability density  $p : \mathbb{R}^d \rightarrow \mathbb{R}$  is*

$$\hat{p}_h(\mathbf{x}) = \frac{1}{n} \sum_{i=1}^n K_h(\mathbf{x} - \mathbf{y}_i), \quad (6)$$

where the kernel  $K_h : \mathbb{R}^d \rightarrow ]0, \infty[$  is the Gaussian function

$$K_h(\mathbf{x}) = \frac{1}{(2\pi)^{\frac{d}{2}} h^d} \exp\left(-\frac{\|\mathbf{x}\|^2}{2h^2}\right) \quad (7)$$

with bandwidth  $h > 0$ .

The following result establishes a connection between linear principal components and ridges of a Gaussian kernel density. It essentially shows that a ridge point lies on a locally defined principal component hyperplane. This hyperplane is determined by a weighted sample mean and the eigenvectors of a weighted sample covariance matrix, where the weights are Gaussian functions.

**Theorem 3.1.** Let  $\hat{p}_h : \mathbb{R}^d \rightarrow \mathbb{R}$  be a Gaussian kernel density estimate, let  $0 < r < d$  and denote the eigenvectors of  $\nabla^2 \log \hat{p}_h(\cdot)$  corresponding to the  $r$  greatest eigenvalues by  $\{\mathbf{v}_i(\cdot)\}_{i=1}^r$ . Define

$$\tilde{\boldsymbol{\mu}}(\mathbf{x}) = \sum_{i=1}^n c_i(\mathbf{x}) \mathbf{y}_i, \quad (8)$$

$$\tilde{\boldsymbol{\Sigma}}(\mathbf{x}) = \sum_{i=1}^n c_i(\mathbf{x}) [\mathbf{y}_i - \tilde{\boldsymbol{\mu}}(\mathbf{x})][\mathbf{y}_i - \tilde{\boldsymbol{\mu}}(\mathbf{x})]^T, \quad (9)$$

where

$$c_i(\mathbf{x}) = \frac{\exp\left(-\frac{\|\mathbf{x} - \mathbf{y}_i\|}{2h^2}\right)}{\sum_{j=1}^n \exp\left(-\frac{\|\mathbf{x} - \mathbf{y}_j\|}{2h^2}\right)}, \quad i = 1, 2, \dots, n.$$

Assume that the eigenvalues of  $\nabla^2 \log \hat{p}_h(\mathbf{x})$  satisfy the condition  $\lambda_1(\mathbf{x}) > \lambda_2(\mathbf{x}) > \dots > \lambda_{r+1}(\mathbf{x})$ . Then

$$\nabla \log \hat{p}_h(\mathbf{x})^T \mathbf{v}_i(\mathbf{x}) = 0 \quad \text{for all } i > r$$

if and only if

$$\mathbf{x} - \tilde{\boldsymbol{\mu}}(\mathbf{x}) \in \text{span}(\tilde{\mathbf{v}}_1(\mathbf{x}), \tilde{\mathbf{v}}_2(\mathbf{x}), \dots, \tilde{\mathbf{v}}_r(\mathbf{x})),$$

where  $\{\tilde{\mathbf{v}}_i(\mathbf{x})\}_{i=1}^r$  denote the eigenvectors of  $\tilde{\boldsymbol{\Sigma}}(\mathbf{x})$  corresponding to the  $r$  greatest eigenvalues.

*Proof.* First, we note the formulae

$$\nabla \log \hat{p}_h(\mathbf{x}) = \frac{\nabla \hat{p}_h(\mathbf{x})}{\hat{p}_h(\mathbf{x})}$$

and

$$\nabla^2 \log \hat{p}_h(\mathbf{x}) = \frac{\nabla^2 \hat{p}_h(\mathbf{x})}{\hat{p}_h(\mathbf{x})} - \frac{\nabla \hat{p}_h(\mathbf{x}) \nabla \hat{p}_h(\mathbf{x})^T}{\hat{p}_h(\mathbf{x})^2}.$$

By a straightforward calculation we then obtain that (cf. proof of Lemma A.1 in Appendix A)

$$h^2 \nabla \log p_h(\mathbf{x}) = -[\mathbf{x} - \tilde{\boldsymbol{\mu}}(\mathbf{x})] \quad (10)$$

and

$$h^4 \nabla^2 \log p_h(\mathbf{x}) + h^2 \mathbf{I} = \tilde{\boldsymbol{\Sigma}}(\mathbf{x}). \quad (11)$$

By equation (11), the matrices  $\nabla^2 \log \hat{p}_h$  and  $\tilde{\boldsymbol{\Sigma}}(\mathbf{x})$  have the same eigenvectors. Hence, by equation (10) the condition that

$$[\mathbf{x} - \tilde{\boldsymbol{\mu}}(\mathbf{x})]^T \tilde{\mathbf{v}}_i(\mathbf{x}) = 0 \quad \text{for all } i > r$$

is equivalent to

$$\nabla \log \hat{p}_h(\mathbf{x})^T \mathbf{v}_i(\mathbf{x}) = 0 \quad \text{for all } i > r,$$

from which the claim follows by the orthogonality of the eigenvectors  $\tilde{\mathbf{v}}_i(\mathbf{x})$ .  $\square$

Ridges of a Gaussian kernel density can be used in an exploratory fashion by adjusting the bandwidth  $h$ . As suggested by Theorem 3.1, this parameter determines the scale of the structures sought from the data.

An important special case arises when  $h$  approaches infinity. At this limit, the  $r$ -dimensional ridge set of the density approaches the  $r$ -dimensional PCA hyperplane, which can be readily observed from equations (8)–(11). A rigorous proof of this property is deferred to Appendix A. Thus, by choosing a large  $h$  we achieve the robustness of PCA but, on the other hand, compromise the ability to describe nonlinear structure in the data.

**Assumption 3.1.** *The  $r + 1$  greatest eigenvalues of the sample covariance matrix  $\hat{\Sigma}_{\mathbf{Y}}$  defined by equation (3) satisfy the conditions  $\lambda_1 > \lambda_2 > \dots > \lambda_{r+1} > 0$ .*

**Theorem 3.2.** *Let  $\hat{p}_h : \mathbb{R}^d \rightarrow \mathbb{R}$  be a Gaussian kernel density estimate, let  $0 \leq r < d$  and let Assumption 3.1 be satisfied. Define the set*

$$S_\infty^r = \left\{ \hat{\boldsymbol{\mu}} + \sum_{i=1}^r \alpha_i \mathbf{v}_i \mid \boldsymbol{\alpha} \in \mathbb{R}^r \right\},$$

where  $\hat{\boldsymbol{\mu}}$  denotes the sample mean (1) and  $\{\mathbf{v}_i\}_{i=1}^r$  denote the eigenvectors of the sample covariance matrix  $\hat{\Sigma}_{\mathbf{Y}}$  corresponding to the eigenvalues  $\{\lambda_i\}_{i=1}^r$ . Then for any compact set  $U \subset \mathbb{R}^d$  such that  $U \cap S_\infty^r \neq \emptyset$  and  $\varepsilon > 0$  there exists  $h_0 > 0$  such that

$$\left. \begin{array}{l} \text{dist}(\mathcal{R}_{\hat{p}_{h^2 \mathbf{I}}}^r \cap U, S_\infty^r) < \varepsilon, \\ \text{dist}(S_\infty^r \cap U, \mathcal{R}_{\hat{p}_{h^2 \mathbf{I}}}^r) < \varepsilon \end{array} \right\} \text{ for all } h \geq h_0,$$

where

$$\text{dist}(S_1, S_2) = \sup_{\mathbf{x} \in S_1} \inf_{\mathbf{y} \in S_2} \|\mathbf{x} - \mathbf{y}\|.$$

### 3.3 Obtaining principal component scores from ridge sets

Based on Proposition 3.1, we now develop the theoretical basis for estimating the first  $m$  nonlinear principal component scores of a given point set. The idea is to obtain the scores one by one by successively projecting the points onto lower-dimensional ridge sets of the underlying density that is estimated by Gaussian kernels. The projections are done along eigenvector curves that are defined by a differential equation. The arc lengths of the curves are interpreted as the principal component scores. As a special case of this approach, we obtain an orthogonal projection onto a linear PCA hyperplane.

For now, we assume that a given point has already been projected onto an  $m$ -dimensional ridge set of its underlying density  $p$  with some  $m \leq d$ . For  $r = 1, 2, \dots, m$ , we define a projection curve  $\gamma_r : \mathbb{R} \rightarrow \mathbb{R}^d$  onto the  $r - 1$ -dimensional

ridge set as a solution to the initial value problem

$$\begin{aligned} \frac{d}{dt} [\mathbf{P}_r(\boldsymbol{\gamma}_r(t)) \nabla \log p(\boldsymbol{\gamma}_r(t))] &= \mathbf{0}, \quad t \geq 0, \\ \boldsymbol{\gamma}_r(0) &= \mathbf{x}_0, \quad \mathbf{x}_0 \in \mathcal{R}_{\log p}^r \setminus \mathcal{R}_{\log p}^{r-1}, \end{aligned} \quad (12)$$

where  $\mathbf{P}_r(\cdot) = \mathbf{I} - \mathbf{v}_r(\cdot) \mathbf{v}_r(\cdot)^T$  and  $\{\mathbf{v}_i(\cdot)\}_{i=1}^d$  denote the eigenvectors corresponding to the eigenvalues  $\lambda_1(\cdot) \geq \lambda_2(\cdot) \geq \dots \geq \lambda_d(\cdot)$  of  $\nabla^2 \log p$ .

We begin with a special case that motivates the above definition and shows its connection to the linear PCA projection. Namely, for any  $d$ -dimensional normal density  $p$ , a ridge point  $\mathbf{x}_0 \in \mathcal{R}_{\log p}^r$ , where  $1 \leq r \leq m$ , can be projected onto the lower-dimensional ridge set  $\mathcal{R}_{\log p}^{r-1}$  by following the solution curve of (12) that is a straight line parallel to the eigenvector  $\mathbf{v}_r$ . This property follows trivially from the definitions of the normal density and the ridge set because  $\log p$  is in this case a quadratic function.

**Proposition 3.2.** *Let  $p$  be a  $d$ -variate normal density with symmetric and positive definite covariance matrix  $\boldsymbol{\Sigma}$  and let  $1 \leq r \leq d$ . If the eigenvalues of  $\boldsymbol{\Sigma}$  satisfy the condition  $\lambda_1 > \lambda_2 > \dots > \lambda_{r+1}$ , then for any solution curve  $\boldsymbol{\gamma}_r$  of the initial value problem (12) we have*

$$\boldsymbol{\gamma}'_r(t) / \|\boldsymbol{\gamma}'_r(t)\| = \pm \mathbf{v}_r$$

for all  $t \geq 0$ . Furthermore, if the sign of  $\boldsymbol{\gamma}'_r$  is chosen such that

$$\boldsymbol{\gamma}'_r(t)^T \nabla \log p(\boldsymbol{\gamma}_r(t)) > 0 \quad \text{for all } t \geq 0,$$

then  $\log p$  has a unique maximum point  $\mathbf{x}^* \in \mathcal{R}_{\log p}^{r-1}$  along the curve  $\boldsymbol{\gamma}_r$ .

When the density  $p$  is not normal, obtaining an expression for the tangent vector  $\boldsymbol{\gamma}'_r(t)$  is nontrivial. However, by utilizing the formula for the derivatives of eigenvectors (e.g. [22]), equation (12) can after some calculation be rewritten as

$$\mathbf{A}_r(\boldsymbol{\gamma}_r(t)) \boldsymbol{\gamma}'_r(t) = \mathbf{0}, \quad (13)$$

where

$$\mathbf{A}_r(\mathbf{x}) = \mathbf{P}_r(\mathbf{x}) \nabla^2 \log p(\mathbf{x}) - \mathbf{F}_r(\mathbf{x}), \quad (14)$$

$$\mathbf{F}_r(\mathbf{x}) = \mathbf{v}_r(\mathbf{x})^T \nabla \log p(\mathbf{x}) \nabla \mathbf{v}_r(\mathbf{x})^T + \mathbf{v}_r(\mathbf{x}) \nabla \log p(\mathbf{x})^T \nabla \mathbf{v}_r(\mathbf{x}) \quad (15)$$

and

$$\nabla \mathbf{v}_r(\mathbf{x}) = [\lambda_r(\mathbf{x}) \mathbf{I} - \nabla^2 \log p(\mathbf{x})]^+ \nabla^3 \log p(\mathbf{x}) \mathbf{v}_r(\mathbf{x}),$$

and the operator “ $+$ ” denotes the Moore-Penrose pseudoinverse (e.g. [10]).

For a general density  $p$ , projection onto the ridge set  $\mathcal{R}_{\log p}^{r-1}$  can still be done by maximizing  $\log p$  along the curve  $\boldsymbol{\gamma}_r$ , but this requires additional justification. To this end, we first show that when  $\boldsymbol{\gamma}_r$  approaches a ridge point  $\mathbf{x}^* \in \mathcal{R}_{\log p}^{r-1}$ , the tangent vector  $\boldsymbol{\gamma}'_r$  becomes parallel to the eigenvector  $\mathbf{v}_r$ . Here we need a technical assumption that will be justified later in this subsection.

**Assumption 3.2.** *The eigenvalues of  $\nabla^2 \log p$  satisfy the conditions*

$$(i) \lambda_1(\boldsymbol{\gamma}_r(t)) > \lambda_2(\boldsymbol{\gamma}_r(t)) > \cdots > \lambda_{r+1}(\boldsymbol{\gamma}_r(t)),$$

$$(ii) \lambda_1(\boldsymbol{\gamma}_r(t)) < 0$$

for all  $t \geq 0$ .

**Proposition 3.3.** *Let  $1 \leq r \leq d$  and let  $\boldsymbol{\gamma}'_r$  denote the normalized tangent vector of a solution curve of (12). If Assumption 3.2 is satisfied and*

$$\lim_{t \rightarrow t^*} \mathbf{v}_r(\boldsymbol{\gamma}_r(t))^T \nabla \log p(\boldsymbol{\gamma}_r(t)) = 0 \quad (16)$$

for some  $t^* > 0$ , then

$$\lim_{t \rightarrow t^*} |\boldsymbol{\gamma}'_r(t)^T \mathbf{v}_r(\boldsymbol{\gamma}_r(t))| = 1. \quad (17)$$

*Proof.* Define the set

$$U = \{\mathbf{x} \in \mathbb{R}^d \mid \lambda_1(\mathbf{x}) < 0 \text{ and} \quad (18)$$

$$\lambda_1(\mathbf{x}) > \lambda_2(\mathbf{x}) > \cdots > \lambda_{r+1}(\mathbf{x})\}. \quad (19)$$

The range of the matrix in the second term of  $\mathbf{F}_r(\mathbf{x})$  defined by equation (15), that is

$$\mathbf{G}(\mathbf{x}) = \mathbf{v}_r(\mathbf{x}) \nabla \log p(\mathbf{x})^T \nabla \mathbf{v}_r(\mathbf{x}),$$

is clearly spanned by the vector  $\mathbf{v}_r(\mathbf{x})$  for all  $\mathbf{x} \in U$ . Furthermore,  $\mathbf{v}_r(\mathbf{x})$  is uniquely determined by condition (19). We also note that the range of the first term of the matrix  $\mathbf{A}_r(\mathbf{x})$  defined by equation (14), that is

$$\mathbf{B}(\mathbf{x}) = \mathbf{P}_r(\mathbf{x}) \nabla^2 \log p(\mathbf{x}),$$

is the set  $\{\mathbf{w} \in \mathbb{R}^d \mid \mathbf{w}^T \mathbf{v}_r(\mathbf{x}) = 0\}$  for all  $\mathbf{x} \in U$ . This follows from the definition of the matrix  $\mathbf{P}_r(\mathbf{x})$ , the eigendecomposition of  $\nabla^2 \log p(\mathbf{x})$  and condition (18) that guarantees nonsingularity of  $\nabla^2 \log p(\mathbf{x})$ .

On the other hand, by the limit (16) the first term of the matrix  $\mathbf{F}_r(\boldsymbol{\gamma}_r(t))$  defined by equation (15), that is

$$\mathbf{v}_r(\boldsymbol{\gamma}_r(t))^T \nabla \log p(\boldsymbol{\gamma}_r(t)) \nabla \mathbf{v}_r(\boldsymbol{\gamma}_r(t))^T$$

converges to zero as  $t$  approaches  $t^*$ . In view of the above observation that the ranges of the matrices  $\mathbf{B}(\mathbf{x})$  and  $\mathbf{G}(\mathbf{x})$  are orthogonal for all  $\mathbf{x} \in U$ , equations (13)–(15) and Assumption 3.2 thus imply that

$$\lim_{t \rightarrow t^*} \mathbf{B}(\boldsymbol{\gamma}_r(t)) \boldsymbol{\gamma}'_r(t) = \mathbf{0} \quad \text{and} \quad \lim_{t \rightarrow t^*} \mathbf{G}(\boldsymbol{\gamma}_r(t)) \boldsymbol{\gamma}'_r(t) = \mathbf{0}.$$

The claim follows from the first of the above limits because the range of the symmetric matrix  $\mathbf{B}(\mathbf{x})$  is orthogonal to its null space.  $\square$

Proposition 3.3 implies the following properties that motivate seeking for a lower-dimensional ridge point by maximizing  $\log p$  along the curve  $\gamma_r$ .

**Proposition 3.4.** *If  $\gamma_r$  is a solution to (12) for some  $1 \leq r \leq d$  and Assumption 3.2 is satisfied, then either  $\gamma_r(t) \in \mathcal{R}_{\log p}^r \setminus \mathcal{R}_{\log p}^{r-1}$  for all  $t \geq 0$  or  $\lim_{t \rightarrow t^*} \gamma_r(t) \in \mathcal{R}_{\log p}^{r-1}$  for some  $t^* > 0$ . In the latter case,  $\log p$  attains its local maximum along  $\gamma_r$  at the limit point  $\gamma_r(t^*)$ .*

*Proof.* By equation (12), the choice of  $\mathbf{x}_0$  and the definition of the matrix  $\mathbf{P}_r(\cdot)$ , for all  $i \neq r$  and  $t \geq 0$  we have

$$\mathbf{v}_i(\gamma_r(t))^T \nabla \log p(\gamma_r(t)) = c_i$$

for some constants  $c_i \neq 0$ . By Assumption 3.2 and Definition 3.1 this implies that either  $\gamma_r(t) \in \mathcal{R}_{\log p}^r \setminus \mathcal{R}_{\log p}^{r-1}$  for all  $t \geq 0$  or  $\lim_{t \rightarrow t^*} \gamma_r(t) \in \mathcal{R}_{\log p}^{r-1}$  for some  $t^* > 0$ . In the latter case we have

$$\mathbf{v}_r(\gamma_r(t^*))^T \nabla \log p(\gamma_r(t^*)) = 0.$$

Thus, the limit (17) implies that

$$\lim_{t \rightarrow t^*} \frac{d}{dt} \log p(\gamma_r(t)) = \lim_{t \rightarrow t^*} \nabla \log p(\gamma_r(t))^T \gamma_r'(t) = 0.$$

Furthermore, by condition (5b) the point  $\gamma_r(t^*)$  is a local maximum of  $\log p$  along  $\gamma_r$ .  $\square$

Recall that our aim is to use projection curves  $\gamma_r$  defined by equation (12) to obtain the first  $m$  nonlinear principal component scores of the given sample points  $\mathbf{y}_i$ . This will be done by using the kernel density  $\log \hat{p}_h$  defined by (6) and (7) as the objective function. Differently to the normal density in Proposition 3.2, this density is not guaranteed to be unimodal or have connected ridge sets. For instance, when  $h$  is too small, the density becomes multimodal.

Unimodality of the density and connectedness of its ridge sets are essential here. This is because as in the linear PCA, our aim is to describe the data in a single well-defined coordinate system. Hence, we assume the following.

**Assumption 3.3.** *Define the set  $U_h = \bigcup_{i=1}^n \mathcal{L}_i^h$ , where*

$$\mathcal{L}_i^h = \{\mathbf{x} \in \mathbb{R}^d \mid \log \hat{p}_h(\mathbf{x}) \geq \log \hat{p}_h(\mathbf{y}_i)\}.$$

*Let  $\lambda_1(\cdot) \geq \lambda_2(\cdot) \geq \dots \geq \lambda_d(\cdot)$  denote the eigenvalues of  $\nabla^2 \log \hat{p}_h$ . Assume that for all  $\mathbf{x} \in U_h$  we have*

$$0 > \lambda_1(\mathbf{x}) > \lambda_2(\mathbf{x}) > \dots > \lambda_{m+1}(\mathbf{x}) \quad (20)$$

*and that set  $U_h$  is connected.*

Damon [5] and Miller [23] give a rigorous treatment of ridge sets of  $C^\infty$ -functions in a differential geometric framework. Under the above assumption, their results guarantee that the  $r$ -dimensional ridge set of the density  $\log \hat{p}_h$  forms a connected manifold in the set  $U_h$  for any  $1 \leq r \leq m$ . Furthermore,  $\log \hat{p}_h$  is unimodal in  $U_h$ . In addition, Assumption 3.3 implies continuity of the Hessian eigenvectors (e.g. [26], Theorem 3.1.3), which is essential for the definition of the initial value problem (12). Assumption 3.3 also entails Assumption 3.2 when the curves  $\gamma_r$  lie in  $U_h$ .

Assumption 3.3 can be satisfied by choosing a sufficiently large  $h$ . The following result is proven in Appendix A.

**Theorem 3.3.** *Under Assumption 3.1 for  $r = m$ , for any Gaussian kernel density estimate  $\hat{p}_h$  there exists  $h_0 > 0$  such that Assumption 3.3 is satisfied for all  $h \geq h_0$ .*

Finally, the arc length of a curve  $\gamma_r$  gives the (curvilinear) distance of its starting point to the ridge set  $\mathcal{R}_{\log \hat{p}_h}^{r-1}$ . Assume that we have projected a given sample point  $\mathbf{y}_i$  onto the ridge set  $\mathcal{R}_{\log \hat{p}_h}^m$ . Starting from such a point, computing the arc lengths successively for  $r = m, m-1, \dots, 1$  then yields the first  $m$  principal component scores of  $\mathbf{y}_i$ . When Assumption 3.3 is satisfied, imposing the conditions (cf. Proposition 3.2)

$$\left. \begin{array}{l} \gamma_r'(t)^T \nabla \log \hat{p}_H(\gamma_r(t)) > 0, \\ \|\gamma_r'(t)\| = 1 \end{array} \right\} \text{ for all } \begin{cases} r = m, m-1, \dots, 1, \\ t \geq 0 \end{cases} \quad (21)$$

guarantees that the curves  $\gamma_r$  lie in the set  $U_h$ . This ensures that the ridge projections are well-defined.

Denote the projection of a given sample point  $\mathbf{y}_i$  onto the set  $\mathcal{R}_{\log \hat{p}_h}^m$  as  $\tilde{\mathbf{y}}_i$  and the starting points of the curves  $\gamma_r$  as  $\mathbf{x}_0^r$ . The  $m$  first principal component scores of  $\mathbf{y}_i$  are then obtained recursively as

$$\theta_r = s_r^* \int_0^{t_r^*} \|\gamma_r'(t)\| dt, \quad \mathbf{x}_0^r = \begin{cases} \tilde{\mathbf{y}}_i, & r = m, \\ \gamma_{r+1}(t_{r+1}^*), & 1 \leq r < m \end{cases} \quad (22)$$

for  $r = m, m-1, \dots, 1$ . Here we assume that for each  $r$  there exists  $t_r^* \geq 0$  such that  $\gamma_r(t_r^*) \in \mathcal{R}_{\log \hat{p}_h}^{r-1}$ . The multiplier  $s_r^* = \lim_{t \rightarrow t_r^* -} s_r(t)$ , where

$$s_r(t) = \begin{cases} 1, & \text{if } \gamma_r'(t)^T \mathbf{v}_r(\gamma_r(t)) > 0, \\ -1, & \text{otherwise,} \end{cases} \quad (23)$$

is introduced to ensure that the principal component score  $\theta_r$  has the correct sign.

### 3.4 Algorithm for estimating principal component scores

Based on the theory developed in Subsection 3.3, we are now ready to develop the algorithm for estimating the nonlinear principal component scores

$$\Theta = [\boldsymbol{\theta}_1 \quad \boldsymbol{\theta}_2 \quad \dots \quad \boldsymbol{\theta}_m]^T \in \mathbb{R}^{n \times m}$$

of a given sample set

$$\mathbf{Y} = [\mathbf{y}_1 \quad \mathbf{y}_2 \quad \cdots \quad \mathbf{y}_n]^T \in \mathbb{R}^{n \times d}$$

for a given  $0 < m \leq d$ . This amounts to first projecting the samples  $\mathbf{y}_i$  onto the ridge set  $\mathcal{R}_{\log \hat{p}_h}^m$  and then successively projecting them onto the lower-dimensional ridge sets  $\mathcal{R}_{\log \hat{p}_h}^r$  until  $r = 0$ . The latter projections are done by tracing the curves  $\gamma_r$  by using a predictor-corrector method. As a by-product, the principal component scores are obtained from a numerical approximation of the integral (22).

A pseudocode of the algorithm is listed as Algorithm 1. It involves the initial projection onto the ridge set  $\mathcal{R}_{\log \hat{p}_h}^m$  (lines 2 and 3), and after that  $m \times n$  loops. Each iteration for  $r = m, m-1, \dots, 1$  projects each of the  $n$  sample points onto the ridge set  $\mathcal{R}_{\log \hat{p}_h}^{r-1}$ . The intermediate projections are stored in the variables  $\{\mathbf{x}_i^*\}_{i=1}^n$ . For the initial ridge projection and the corrector steps, the algorithm utilizes the trust region Newton method developed in [30] (the GTRN algorithm). This method is briefly described at the end of this subsection.

In the following, we describe the steps for carrying out one ridge projection (i.e. one iteration of the loop over the index  $i$ ) for a given  $r$ . The starting point  $\mathbf{x}_0$  for  $\gamma_r$  is chosen as  $\mathbf{x}_i^*$  representing the projection of the sample point  $\mathbf{y}_i$  onto the set  $\mathcal{R}_{\log \hat{p}_h}^r$ . Assuming that there exists a monotonously increasing sequence  $\{t_k\}$  such that  $\gamma_r(t_{k^*}) \in \mathcal{R}_{\log \hat{p}_h}^{r-1}$  for some  $k^*$ , we introduce the notation  $\mathbf{x}_k = \gamma_r(t_k)$  for the iterates along the curve  $\gamma_r$ . With this notation, an approximation to the integral (22) is given by

$$\int_0^{t_r^*} \|\gamma_r'(t)\| dt \approx \sum_{k=1}^{k^*} \|\gamma_r(t_k) - \gamma_r(t_{k-1})\| = \sum_{k=1}^{k^*} \|\mathbf{x}_k - \mathbf{x}_{k-1}\|.$$

The algorithm uses a predictor-corrector method to generate the iterates  $\mathbf{x}_k$ . At the predictor step (line 18), the algorithm proceeds along a tangent vector  $\mathbf{u}_k = \gamma_r'(t_k)$  solved from equation (13). That is,

$$\tilde{\mathbf{x}}_k = \mathbf{x}_k + \tau s_k \mathbf{u}_k,$$

where  $\tau > 0$  is some user-supplied step size,  $\|\mathbf{u}_k\| = 1$  and the multiplier

$$s_k = \begin{cases} 1, & \text{if } \mathbf{u}_k^T \nabla \log \hat{p}_h(\mathbf{x}_k) > 0, \\ -1, & \text{otherwise} \end{cases}$$

(line 9) is introduced to impose conditions (21). To project the predictor estimate  $\tilde{\mathbf{x}}_k$  back to the ridge set  $\mathcal{R}_{\log \hat{p}_h}^r$ , the algorithm takes a corrector step (line 19).

A stopping criterion is imposed to terminate the tracing of the curve  $\gamma_r$  when a maximum of  $\log \hat{p}_h$  along  $\gamma_r$  is encountered (line 11). For  $k > 0$ , the condition

$$s_{k-1} \mathbf{u}_{k-1}^T \mathbf{u}_k s_k < 0$$



---

**Algorithm 1:** Nonlinear principal component scores

---

**input** : sample points  $\mathbf{Y} = [\mathbf{y}_1 \ \mathbf{y}_2 \ \dots \ \mathbf{y}_n]^T \in \mathbb{R}^{n \times d}$   
Gaussian kernel density estimate  $\hat{p}_h : \mathbb{R}^d \rightarrow \mathbb{R}$   
ridge dimension  $0 < m \leq d$   
step size  $\tau > 0$

**output:** principal component scores  $\Theta = [\theta_1 \ \theta_2 \ \dots \ \theta_n]^T \in \mathbb{R}^{n \times m}$

```
1  $\Theta \leftarrow \mathbf{0}$ 
2 for  $i = 1, 2, \dots, n$  do
3    $\mathbf{x}_i^* \leftarrow \text{GTRN}(\hat{p}_h, m, \mathbf{y}_i, \tau, 10^{-5})$ 
4   for  $r = m, m - 1, \dots, 1$  do
5     for  $i = 1, 2, \dots, n$  do
6        $\mathbf{x}_0 \leftarrow \mathbf{x}_i^*$ 
7       for  $k = 0, 1, \dots$  do
8         Obtain the tangent vector  $\mathbf{u}_k$  from (13) such that  $\|\mathbf{u}_k\| = 1$ .
9         if  $\mathbf{u}_k^T \nabla \log \hat{p}_h(\mathbf{x}_k) > 0$  then  $s_k \leftarrow 1$  else  $s_k \leftarrow -1$ 
10        if  $k > 0$  then
11          if  $s_{k-1} \mathbf{u}_{k-1}^T \mathbf{u}_k s_k < 0$  then
12             $\mathbf{x}_i^* \leftarrow \text{GTRN}(\hat{p}_h, r - 1, (\mathbf{x}_{k-1} + \mathbf{x}_k)/2, 0.5\tau, 10^{-5})$ 
13             $\theta_{i,r} \leftarrow \theta_{i,r} + \|\mathbf{x}_i^* - \mathbf{x}_{k-1}\|$ 
14            if  $(\mathbf{x}_i^* - \mathbf{x}_{k-1})^T \mathbf{v}_r(\mathbf{x}_i^*) < 0$  then  $\theta_{i,r} \leftarrow -\theta_{i,r}$ 
15            Terminate the inner loop.
16          else
17             $\theta_{i,r} \leftarrow \theta_{i,r} + \|\mathbf{x}_k - \mathbf{x}_{k-1}\|$ 
18           $\tilde{\mathbf{x}}_k \leftarrow \mathbf{x}_k + \tau s_k \mathbf{u}_k$ 
19           $\mathbf{x}_{k+1} \leftarrow \text{GTRN}(\hat{p}_h, r, \tilde{\mathbf{x}}_k, 0.5\tau, 10^{-5})$ 
```

---

tests whether the gradient changes sign along the curve. When this condition is met, the algorithm projects the midpoint of the current and previous iterate onto a nearby ridge point  $\mathbf{x}_i^* \in \mathcal{R}_{\log \hat{p}_h}^{r-1}$  (line 12). At line 14, the algorithm computes the sign  $s_r^*$  for the integral (22) by approximately testing condition (23). The inner iteration (i.e. iteration of the loop over the index  $k$ ) is then terminated, and the point  $\mathbf{x}_i^*$  is retained as a starting point for projection onto a lower-dimensional ridge set.

The GTRN algorithm developed in [30] implements a Newton-type method for projecting a  $d$ -dimensional point onto an  $r$ -dimensional ridge set of a probability density. The method successively maximizes a quadratic model of the objective function. The maximization is constrained within a *trust region* to guarantee convergence. To obtain a ridge projection, it is done in the subspace spanned by the Hessian eigenvectors corresponding to the  $d - r$  smallest eigenvalues. That is, at each iteration  $l$  the next iterate  $\mathbf{z}_{l+1} = \mathbf{z}_l + \mathbf{s}_l$  is obtained by solving the

subproblem

$$\max_{\mathbf{s}} Q_l(\mathbf{s}) \quad \text{s.t.} \quad \begin{cases} \|\mathbf{s}\| \leq \Delta_l, \\ \mathbf{s} \in \text{span}(\mathbf{v}_{r+1}(\mathbf{z}_l), \mathbf{v}_{r+2}(\mathbf{z}_l), \dots, \mathbf{v}_d(\mathbf{z}_l)), \end{cases}$$

where  $Q_l$  denotes the model function at the current iterate  $\mathbf{z}_l$ ,  $\{\mathbf{v}_i(\mathbf{z}_l)\}_{i=r+1}^d$  denote the eigenvectors and  $\Delta_l \leq \Delta_{\max}$  denotes the current trust region radius that is updated after each iteration. For each call of GTRN, Algorithm 1 uses the experimentally chosen  $\Delta_{\max} = 0.5\tau$  ( $\tau$  for the initial projection) and stopping criterion threshold  $\varepsilon_{pr} = 10^{-5}$ .

**Remark 3.1.** *The GTRN algorithm can be viewed as an approximate solution method to an initial value problem of the form (12), where  $\mathbf{P}_r(\cdot) = \sum_{i=1}^r \mathbf{v}_i(\cdot)\mathbf{v}_i(\cdot)^T$ . As Algorithm 1, GTRN yields an orthogonal projection when applied to the logarithm of a normal density. Differently to Algorithm 1, projection of a  $d$ -dimensional point onto an  $r$ -dimensional ridge set by this algorithm only requires continuity of the first  $r$  Hessian eigenvectors. That is, when the  $r + 1$  greatest eigenvalues are distinct in  $U_h$ .*

**Remark 3.2.** *Tests for connectedness of the ridge sets are not included in Algorithm 1 for the sake of simplicity. This can be done by testing if the traced curve crosses a point  $\mathbf{x}$  where  $\lambda_1(\mathbf{x}) = 0$  or  $\lambda_i(\mathbf{x}) = \lambda_j(\mathbf{x})$  for some  $i, j = 1, 2, \dots, r + 1$  such that  $i \neq j$ , where  $\lambda_i(\cdot)$  denote the eigenvalues of  $\nabla^2 \log \hat{p}_h$  (see [23, 29]). When either of these conditions is met, the algorithm can be restarted with a larger  $h$  or smaller initial ridge dimension  $m$ .*

**Remark 3.3.** *A consistent orientation of the eigenvectors  $\mathbf{v}_r(\mathbf{x}_i^*)$  at the projected points is necessary for the principal component scores  $\theta_{i,r}$  to have correct signs. However, in practice the signs of the eigenvectors depend on the numerical algorithm for computing them. Therefore, the implementation of Algorithm 1 uses an Euclidean minimum spanning tree (e.g. [17]) to align the eigenvectors after each iteration of the outer loop.*

## 4 Nonlinear extension of SSA to time series data

In this section, the KDPCA method developed in Section 3 is extended to time series data. The method, that we call KDSSA, is based on the singular spectrum analysis (SSA) that is an extension of the linear PCA. In SSA, a time series is embedded in a multidimensional *phase space*. This is done by constructing a *trajectory matrix* from time-lagged copies of the time series. Formally, the trajectory

matrix of a time series  $\mathbf{x} = (x_1, x_2, \dots, x_n)$  is defined as

$$\mathbf{Y}_{\mathbf{x},L} = \begin{bmatrix} x_1 & x_2 & x_3 & \cdots & x_L \\ x_2 & x_3 & x_4 & \cdots & x_{L+1} \\ x_3 & x_4 & x_5 & \cdots & x_{L+2} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ x_{n-L+1} & x_{n-L+2} & x_{n-L+3} & \cdots & x_n \end{bmatrix}, \quad (24)$$

where  $L$  is some user-supplied time window length.

Applying the linear PCA to the above matrix, one can obtain the principal components and the reconstructed time series by using the formulae given by Vautard et al. [38]. Generalizing their approach, we minimize the reconstruction error

$$E(\mathbf{x}) = \sum_{i=1}^{n-L+1} \sum_{j=1}^L (\tilde{y}_{i,j} - x_{i+j-1})^2 \quad (25)$$

using the first  $m$  nonlinear principal components, where  $m \leq L$ . Here the vectors  $\tilde{\mathbf{y}}_i$  denote the projections of the row vectors  $\mathbf{y}_i$  of  $\mathbf{Y}_{\mathbf{x},L}$  onto the  $m$ -dimensional ridge set of their Gaussian kernel density estimate.

A straightforward calculation shows that by equating the gradient  $\nabla E(\mathbf{x})$  to zero, we obtain the formulae

$$x_i^* = \begin{cases} \frac{1}{L} \sum_{j=1}^L \tilde{y}_{i-j+1,j}, & L \leq i \leq n - L + 1 \\ \frac{1}{i} \sum_{j=1}^i \tilde{y}_{i-j+1,j}, & 1 \leq i \leq L - 1 \\ \frac{1}{n - i + 1} \sum_{j=i-n+L}^L \tilde{y}_{i-j+1,j}, & n - L + 2 \leq i \leq n \end{cases} \quad (26)$$

for the elements of the reconstructed time series such that  $E(\mathbf{x}^*)$  minimizes the reconstruction error (25).

**Remark 4.1.** *SSA can be extended to multivariate time series [38] and to time-series in a two-dimensional grid [12]. As these modifications differ from the univariate SSA only by the construction of the trajectory matrix and the formulae for obtaining the reconstructed time series, we do not consider the multivariate case here.*

In this paper the nonlinear SSA is applied to quasiperiodic time series (i.e. noisy time series having some underlying periodic pattern). The motivation is as follows. Assuming that a time series follows the model

$$X(t) = f(t) + \varepsilon(t)$$

for some periodic function  $f$  and  $\varepsilon$  representing the noise, it is reasonable to model the trajectory samples (i.e. the rows of the matrix  $\mathbf{Y}_{x,L}$ ) as a point set that is randomly distributed around a closed curve (cf. Figure 5 in Section 5).

When the aim is to obtain a noise-free time series from a reconstructed phase space trajectory, only an approximate projection onto the ridge curve (i.e. one-dimensional ridge set) of the trajectory density suffices. The GTRN algorithm developed in [30] is appropriate for this purpose. On the other hand, a parametrization of the reconstructed trajectory can be obtained by the algorithm developed in [29]. Differently to Algorithm 1, this algorithm yields a continuous parametrization even when the trajectory density is multimodal, provided that the set of its ridge curves forms a single closed loop. Both of the aforementioned approaches are demonstrated in the next section.

## 5 Practical applications

This section is devoted to practical applications of the proposed KDPCA method. The method is applied to a climate model output that exhibits a highly nonlinear behaviour. In addition, its SSA-based extension is applied to an atmospheric time series. Finally, computational complexity analysis and comparison to related methods is given.

### 5.1 Test setup

All numerical tests were done on a machine with a 3.0GHz Core 2 Duo processor and 6GB system memory. Algorithm 1 as well as the algorithms developed in [29] and [30] used in the tests were implemented in Fortran 95. Algorithm 1 was run with  $m = d$  and  $\tau = 0.05h$ . For the nonlinear SSA, the algorithms of [29] and [30] were run with their default parameters, except for GTRN the parameters  $\Delta_{\max}$  and  $\varepsilon_{pr}$  were chosen as  $0.25h$  and  $10^{-4}$ , respectively.

### 5.2 Application to simulated climate model data

In the first test, KDPCA was applied to a simulated sea surface temperature dataset. This dataset is provided by the National Oceanic and Atmospheric Administration (NOAA). The data was obtained from the Coupled Model Intercomparison Project phase 3 (CMIP3) simulations of the GFDL-CM2.1 climate model [1, 6]. The chosen simulation output is from the pre-industrial control (`picntrl`) simulation experiment. All preprocessing steps were done as in [33] and [34], where this dataset has been analyzed in detail.

The simulated sea surface temperatures are computed on a latitude-longitude grid. The subregion  $20^\circ$  S- $20^\circ$  N,  $125^\circ$  E- $65^\circ$  W representing the pacific ocean was chosen from the original data. The temperatures are monthly values, and

the length of the simulation is 500 years, starting from the first century AD. This yields a dataset of 6000 samples. As a preprocessing step, seasonality was removed from the data by subtracting the monthly mean values of the raw data for each month. In climatological literature, such mean values are called *anomalies*.

To make estimation of the nonlinear principal components computationally feasible, the high-dimensional data ( $d = 10073$ , one dimension for each ocean grid point) was first projected onto the first ten principal components obtained by PCA. As these principal components explain 87.3% of variance, a significant amount of information was not lost by carrying out this preprocessing step. The kernel bandwidth was chosen as  $h = 40$ .

The GFDL-CM2.1 dataset and its first principal component obtained by KDPCA are plotted in Figure 2. This figure shows cross-sections of the data and the principal component curve along the first linear principal component axes. Projection of the GFDL-CM2.1 data onto the surface spanned by the first two principal components obtained by KDPCA is shown in Figure 3. The corresponding principal component scores are plotted in Figure 4.

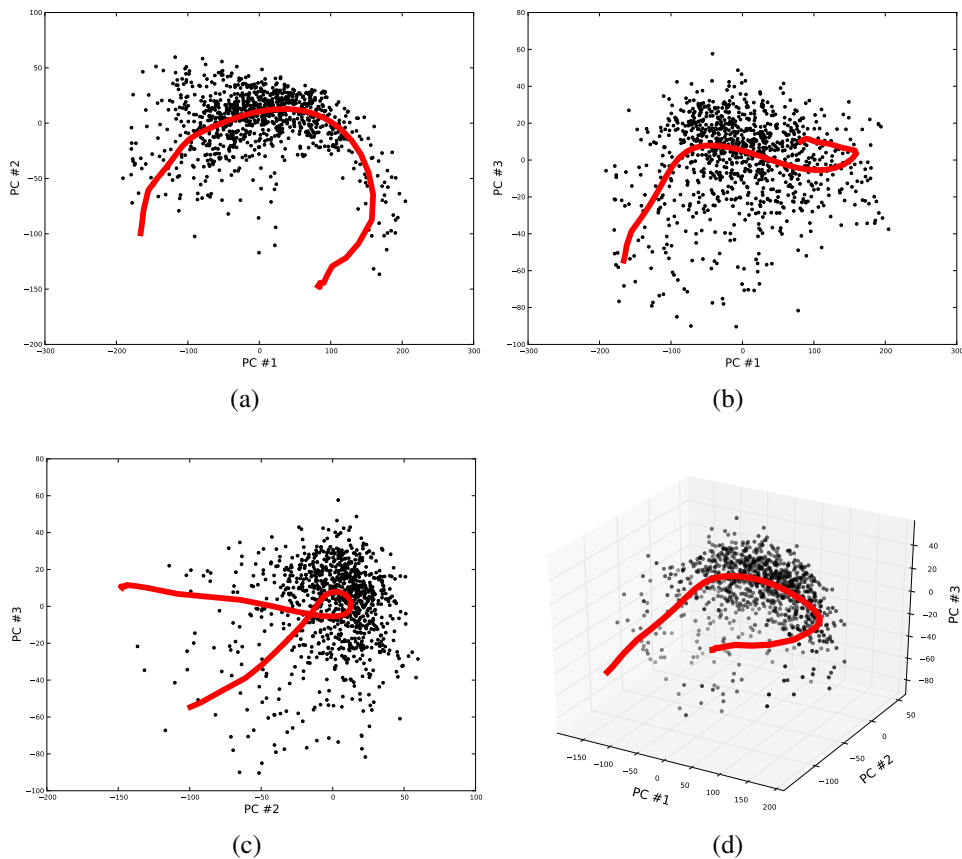


Figure 2: The first nonlinear principal component obtained from the GFDL-CM2.1 dataset (only a subset of the curve is drawn).

Compared to the linear principal component projection shown in Figure 2a, it is clear that the nonlinear principal components represent the "unfolded" dataset and they are better able to capture the variance in the data. Comparison of explained variances of the first eight linear and nonlinear principal components listed in Table 1 also supports this claim. The variance explained by KDPCA is more concentrated towards the first principal component than the variance explained by PCA. Here the explained variances for the nonlinear principal components were obtained from the covariance matrix of the corresponding scores.

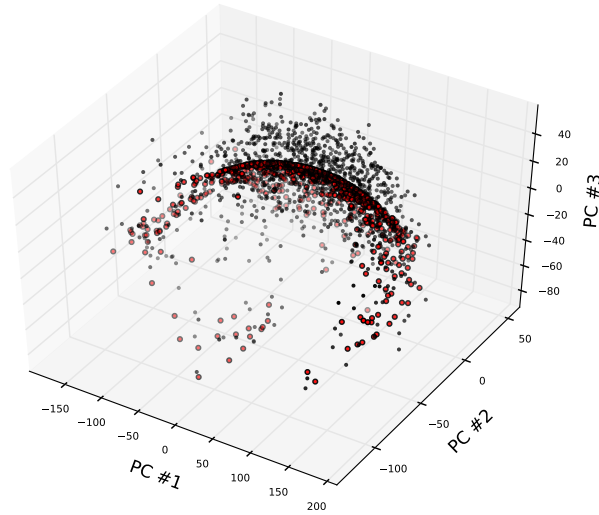


Figure 3: Projection of the GFDL-CM2.1 dataset onto the surface spanned by its first two nonlinear principal components.

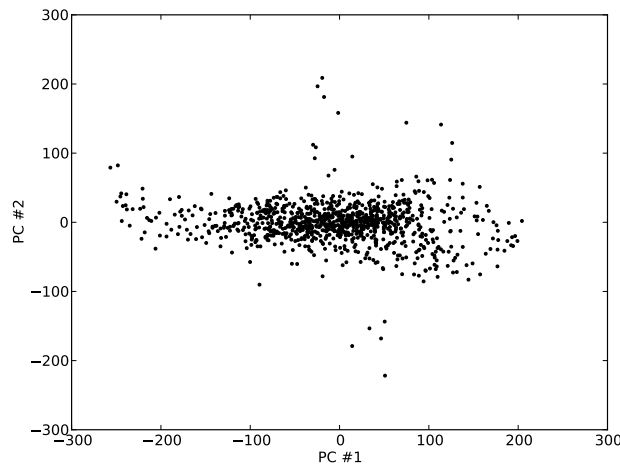


Figure 4: Two first nonlinear principal component scores obtained from the GFDL-CM2.1 dataset.

	1	2	3	4	5	6	7	8
PCA	60.0 %	10.6 %	6.1 %	2.5 %	2.1 %	1.7 %	1.3 %	1.1 %
KDPCA	66.2 %	10.4 %	3.9 %	2.1 %	1.4 %	1.1 %	0.9 %	0.8 %

Table 1: Explained variances of the eight first linear and nonlinear principal components, GFDL-CM2.1 dataset.

A typical application of principal component analysis (and its nonlinear extensions) is to explain the variance in the given data by some small set of variables. This has been done in [33] and [34] for the GFDL-CM2.1 data, and the two main sources of variation were identified. The first principal component correlates with the so-called NINO3 index that is related to the El Niño Southern Oscillation (ENSO) phenomenon. The second one correlates with the Pacific warm water volume. The analysis done here could be carried out further, but we do not attempt repeat the earlier experiments by [33] and [34], as using KDPCA would yield similar results than the earlier nonlinear PCA extensions. Of more interest are the differences between KDPCA and the previously proposed methods. A discussion of potential advantages of using KDPCA is given in Subsection 5.5.

### 5.3 Application to atmospheric time series

The quasi-biennial oscillation (QBO) is one of the most well-studied atmospheric phenomena. The QBO is a quasiperiodic oscillation of the equatorial zonal wind between easterlies and westerlies in the tropical stratosphere with a mean period of 28 to 29 months. Motivated by an earlier neural network-based nonlinear SSA approach in [16], the nonlinear SSA (KDSSA) described in Section 4 was applied to a QBO time series. The time series is provided by the institute of meteorology at the University of Berlin [2]. It consists of monthly mean zonal winds between 1953-2013 constructed from balloon observations at seven different pressure levels corresponding to the altitude range 20-30 km. Here we use a simplified test setup and analyze only the observations from the 30 Hpa level, resulting to a univariate time series.

The trajectory matrix (24) was obtained from the QBO time series with  $L = 18$  months. The linear PCA was first applied in the 18-dimensional phase space so that the first four principal components were retained. These principal components explain 95.2 % of the variance, and thus a significant amount of information was not lost by doing this step. The resulting samples were then projected onto the kernel density ridges by using the GTRN algorithm described in [30]. The bandwidth was chosen as  $h = 200$ . The reconstructed time series was obtained by transforming the projected samples from the four-dimensional space back to the 18-dimensional phase space and using the formulae (26).

The trajectory samples and their kernel density ridge projections in the reduced four-dimensional phase space are plotted in Figure 5. This figure shows a

cross-section along the first three linear principal components. Due to the underlying periodic structure present in the time series, its reconstructed phase space trajectory forms a closed loop that passes through the middle of the point cloud. The QBO time series and the reconstructed time series obtained by using the reconstructed phase space trajectory are plotted in Figure 6. For comparison, the reconstructed time series obtained by using the first linear SSA component and the first and second linear components combined are also plotted in this figure.

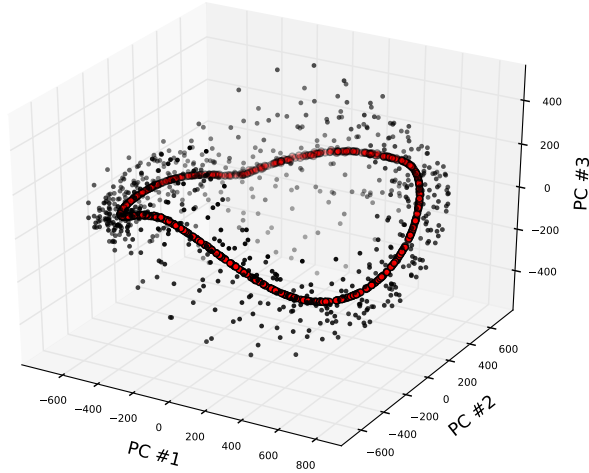


Figure 5: Phase space trajectory of the QBO time series and the reconstructed trajectory curve obtained by kernel density ridge projection.

The conclusion from Figures 5 and 6 is that the nonlinear SSA is able to capture the underlying periodic structure in the QBO time series. It is clear that the closed loop found by the nonlinear approach, as shown in Figure 5, cannot be described by any combination of linear principal components. Consequently, it can be seen from Figure 6 that the linear SSA reconstruction by using only the first principal component is inadequate to describe the structure of the time series. On the other hand, by adding more principal components in the analysis, the linear SSA only captures noise and not the underlying periodic pattern.

In Subsection 5.3, the principal component scores (i.e. the coordinates along the nonlinear principal components) were of main interest. Also, in the nonlinear SSA tracking the coordinates of a time series along its reconstructed trajectory curve in the phase space may provide useful information. Namely, when the time series is close to periodic, anomalously short or long cycles can be identified by carrying out such analysis. For the QBO time series, this has been done in [16] by using the neural network-based NLPCA.

Obtaining the coordinates of a time series along its reconstructed phase space trajectory is also possible by using the ridge-based approach. In order to demonstrate this, an approximate parametrization of the trajectory was obtained by us-



ing the algorithm developed in [29]. The projection coordinates were obtained for each sample by finding the nearest point along line segments connecting the trajectory points and computing its distance along the approximate curve to a point fixed as the origin.

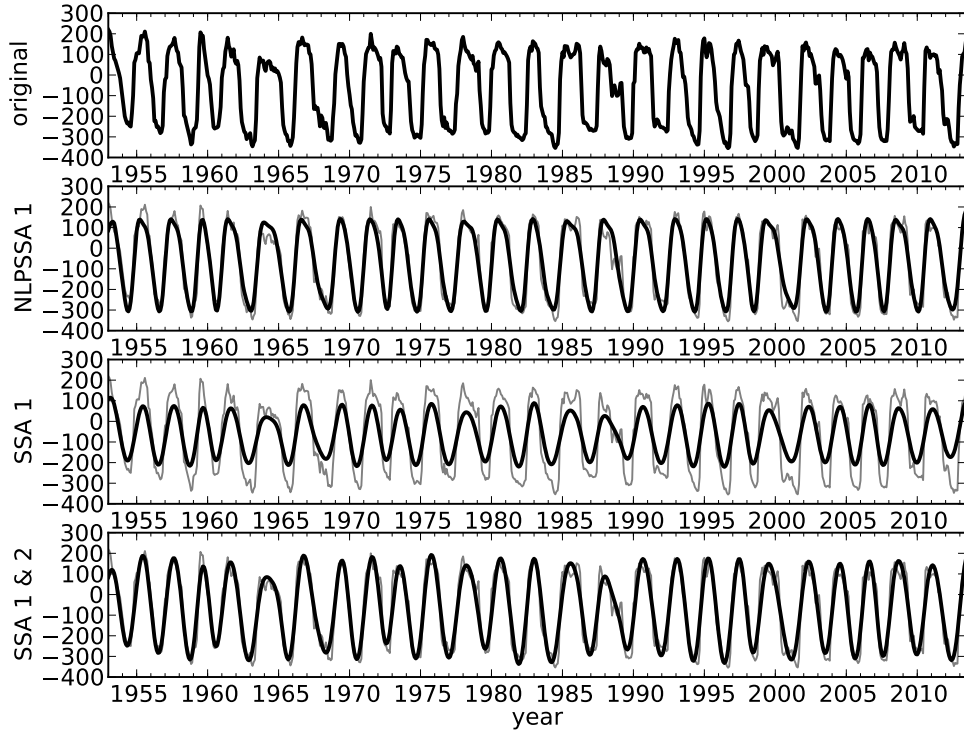


Figure 6: The QBO time series at the 30 Hpa level and the reconstructed time series obtained by using the first KDSSA component, the first linear SSA component and the two first linear SSA components combined. The original time series is plotted in gray in the lower figures.

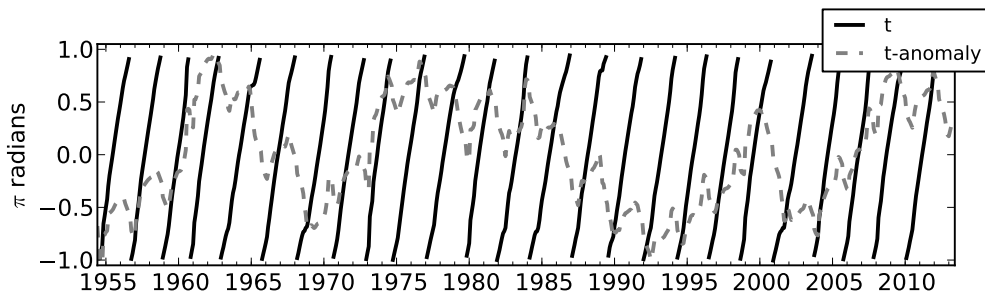


Figure 7: The first nonlinear principal component coordinate of the QBO time series ( $t$ ) and the deviation from the fitted regression line ( $t$ -anomaly).

Due to its very regular period, the QBO time series progresses along its re-

constructed phase space trajectory at a nearly constant rate. This can be seen from Figure 7 showing the trajectory coordinate  $t$  scaled to the interval  $[-\pi, \pi]$  as a function of time. In addition, following the approach of [16], anomalies (i.e. deviations from the constant rate) were calculated. This was done by fitting a regression line to the  $t$ -time series obtained by concatenating the individual cycles and then subtracting the regression line from the concatenated time series. The normalized  $t$ -anomaly time series obtained in this way is also plotted in Figure 7.

Comparison of the  $t$ -anomaly time series to the  $t$ -time series and Figure 6 shows its relation to fluctuations from the mean period length. Namely, up- and downward trends in the  $t$ -anomalies correspond to abnormally short and long cycles, respectively. This can be seen, for instance, by comparing the long periods during 1962-1969, 1984-1993 and 2000-2002 and the short periods during 1955-1962, 1969-1975 and 2004-2009 to the  $t$ -anomaly time series.

## 5.4 Complexity Analysis

This subsection is devoted to discussion of computational complexity of KDPCA and comparison with existing nonlinear dimensionality reduction methods. After the initial projection step by using the GTRN algorithm having computational cost  $\mathcal{O}(n^2d^2 + nd^3)$  (see [30]), the computational cost of KDPCA is  $\mathcal{O}(n^2d^3m + d^3nm + n^2dm)$ , which is explained in the following paragraphs.

The main source of computational cost is the evaluation of the Gaussian kernel density and its derivatives. For each of the  $m$  projection steps, this needs to be done for all  $n$  sample points a number of times that depends on the chosen step size  $\tau$ . For the third derivative that dominates the computational cost, the cost of a single evaluation is  $\mathcal{O}(nd^3)$ . This makes the total complexity of evaluations  $\mathcal{O}(n^2d^3m)$ . When  $d$  is small, this cost can be reduced by order of  $n$  by using the fast Gauss transform [13] or related techniques.

Computation of the tangent vector in Algorithm 1 and obtaining the trust region step in the corrector involve eigendecomposition of a  $d \times d$  matrix [30]. The cost of this operation is  $\mathcal{O}(d^3)$ , and this is done  $\mathcal{O}(nm)$  times in the algorithm, making the total cost of eigendecompositions  $\mathcal{O}(d^3nm)$ . Finally, the cost of traversing the Euclidean minimum spanning tree by using a basic implementation is  $\mathcal{O}(n^2d)$ . This is done  $m$  times in the algorithm, after each projection of all the sample points, and thus the total cost of traversal of such trees is  $\mathcal{O}(n^2dm)$ .

Computational efficiency of Algorithm 1 can be improved by replacing the projection curve tangent by a Hessian eigenvector (cf. Propositions 3.2 and 3.3). In practice, this leads to slightly worse approximations for the higher-dimensional principal component scores (the first principal component is not affected). This approximation reduces the evaluation cost by order of  $d$  since third derivatives are not needed.

When only the first nonlinear principal component (i.e. principal curve) is sought, a significant speedup can be achieved by using a specialized algorithm

developed in [29]. Using this algorithm requires choosing the kernel bandwidth so that the ridge curve set of the density consists of one connected curve. Under this assumption, it suffices to use one starting point, and the total computational cost of tracing the ridge curve is  $\mathcal{O}(nd^3)$ . The principal component scores can be obtained from projections onto the line segments forming the approximate curve as in Subsection 5.3 at a cost of  $\mathcal{O}(ndk)$ , where  $k$  is the number of line segments.

## 5.5 Comparison to Other Methods

Though the neural network-based NLPCA (e.g. [15, 20, 24, 35]) has become popular particularly in climate analysis, it has several shortcomings. Some of them are discussed below.

- NLPCA involves minimization of a complicated cost function that generally has a large number of local minima. This problem is typically addressed by using a large number of starting points, which may incur a high computational cost. KDPCA is not affected by this issue because it does not attempt to minimize a single cost function. Instead, it performs local maximizations from each sample point. The projection curves are uniquely defined when the ridge sets are connected.
- The principal components obtained by KDPCA have a probabilistic interpretation. This is not the case for NLPCA that is based on an artificially constructed neural network. In fact, the NLPCA principal curves and surfaces are not guaranteed to follow regions of high concentration of the data points. Examples of this are given in [4]. Due to this issue, drawing statistical inferences from the NLPCA output should be done with extreme caution.
- NLPCA uses a number of artificial penalty terms to avoid overfitting. Despite this, the density of the data along the first nonlinear principal component can exhibit spurious multimodality [4]. This can occur even when the underlying density of the data is close to normal. On the other hand, KDPCA always attains the robustness of the linear PCA when the kernel bandwidth  $h$  is chosen sufficiently large.
- When using NLPCA, the type of the curve (open or closed) to be fitted to the data needs to be specified a priori in the neural network structure. KDPCA can determine this automatically when the principal curve is traced by using the algorithm developed in [29].
- The curves fitted by NLPCA are not parametrized by arc length. This may introduce a significant bias to the principal component scores. When drawing statistical inferences from a curve fitted by NLPCA, arc length reparametrization should be done to remove the bias [25]. However, this approach has not been generalized to higher dimensions. On the other hand, KDPCA

produces an arc length parametrization for principal component curves and surfaces of any dimension due to its construction.

KDPCA has also certain advantages compared to other commonly used non-linear dimensionality reduction methods. This is because it seems to perform well in the presence of noise and it operates directly in the input space.

- Graph-based methods such as Isomap [37], Hessian eigenmaps [7] and maximum variance unfolding [40] are based on the assumption that the data lies directly on a low-dimensional manifold. Thus, they are sensitive to noise, and blindly applying such methods to noisy data may lead to undesired results.
- The aforementioned methods and kernel-based methods such as KPCA do not produce a reconstruction of the data in the original input space. This would be a very desired feature, for instance, in climate analysis where plots of reconstructed grid data or time series provide information about the main sources of variation.

KDPCA has also some shortcomings that it shares with NLPCA. Namely, the bandwidth parameter  $h$  plays a similar role to the penalty parameters in NLPCA. KDPCA breaks down when  $h$  is too small. In this case, the ridges of the density estimate no longer form a connected set, and the density becomes multimodal. The issue of choosing a sufficiently large  $h$  is difficult. In the absence of a more sophisticated method, it can be done by visual inspection of the obtained principal components in two or three dimensions (cf. Figures 2 and 5). A possible approach could be to use some computationally cheap plugin estimate as an initial guess for the bandwidth.

## 6 Conclusions and discussion

Principal component analysis (PCA) is a well-established tool for exploratory data analysis. However, as a linear method it cannot describe complex nonlinear structure in the given data. To address this deficiency, a novel nonlinear generalization of the linear PCA was developed in this paper.

The proposed KDPCA method is based on the idea of using ridges of the underlying density of the data to estimate nonlinear structures. It was shown that the principal component coordinates of a given point set can be obtained by successively projecting the points onto lower-dimensional ridge sets of the density. Such a projection was defined as a solution to a differential equation. A predictor-corrector method using a Newton-based corrector was developed for this purpose.

Gaussian kernels were used for estimation of the density from the data. This choice has several advantages. First, a fundamental result was derived showing that by letting  $h$  approach infinity, KDPCA reduces to the linear PCA and achieves

its robustness when desired. Second, the theory of ridge sets ensures that any ridge set of a Gaussian kernel density estimate has a well-defined coordinate system when  $h$  is sufficiently large. Third, this choice allows automatic estimation of an appropriate bandwidth from the data, though this approach was not pursued in this paper.

Based on the linear singular spectrum analysis (SSA), KDPCA was extended to time series data. It was shown that when a time series is (quasi)periodic, the first nonlinear principal component of its phase space representation can be used to reconstruct the underlying periodic pattern from noise. Though the periodicity assumption is restrictive, such time series are relevant for many practical applications. Examples include climate analysis and medical applications such as electrocardiography and electroencephalography.

The proposed KDPCA method and its SSA-based variant were applied to a highly nonlinear dataset obtained from a climate model and to an atmospheric time series. The method is superior to the linear PCA in capturing the complex nonlinear structure of the data. It also has several advantages compared to the existing nonlinear dimensionality reduction methods. In particular, KDPCA requires only one parameter, that is, the kernel bandwidth  $h$ . The bandwidth has an intuitive interpretation as a scale parameter.

While KDPCA showed impressive results on the test datasets, its applicability to real-world data remains to be fully confirmed. When the dataset is noisy and sparse, which is typical for observational data, the additional information obtained by KDPCA might not justify its high computational cost. However, using the fast Gauss transform mentioned in Subsection 5.4 could significantly improve the scalability of KDPCA to large datasets. Computational difficulties due to high dimensionality of the input data can also be circumvented. In many situations, the variance is confined to some low-dimensional linear subspace that can be identified by using a simpler method as a preprocessing step. Examples of this were given in Subsections 5.2 and 5.3, where KDPCA was applied to a low-dimensional projection obtained by the linear PCA.

**Acknowledgements.** The author was financially supported by the TUCS Graduate Programme. He would also like to thank Prof. Marko Mäkelä, Doc. Napsu Karmitsa and Prof. Hannu Oja for their valuable comments.

## References

- [1] CM2.X Coupled Climate Models. <http://nomads.gfdl.noaa.gov/CM2.X>. visited on 18/11/2013.
- [2] The quasi-biennial-oscillation (QBO) data serie. <http://www.geo.fu-berlin.de/en/met/ag/strat/produkte/qbo>. visited on 18/11/2013.

- [3] G. Berkooz, P. Holmes, and J. L. Lumley. The proper orthogonal decomposition in the analysis of turbulent flows. *Annual Review of Fluid Mechanics*, 25:539–575, 1993.
- [4] B. Christiansen. The shortcomings of nonlinear principal component analysis in identifying circulation regimes. *Journal of Climate*, 18(22):4814–4823, 2005.
- [5] J. Damon. Generic structure of two-dimensional images under Gaussian blurring. *SIAM Journal on Applied Mathematics*, 59(1):97–138, 1998.
- [6] T. L. Delworth and coauthors. GFDL’s CM2 global coupled climate models. part I: Formulation and simulation characteristics. *Journal of Climate*, 19(5):643–674, 2006.
- [7] D. L. Donoho and C. Grimes. Hessian eigenmaps: Locally linear embedding techniques for high-dimensional data. *Proceedings of the National Academy of Sciences*, 100(10):5591–5596, 2003.
- [8] J. Einbeck, L. Evers, and C. Bailer-Jones. Representing complex data using localized principal components with application to astronomical data. In A. N. Gorban, B. Kégl, D. C. Wunsch, and A. Y. Zinovyev, editors, *Principal Manifolds for Data Visualization and Dimension Reduction*, volume 58 of *Lecture Notes in Computational Science and Engineering*, pages 178–201. Springer, Berlin, Heidelberg, 2008.
- [9] J. Einbeck, G. Tutz, and L. Evers. Local principal curves. *Statistics and Computing*, 15(4):301–313, 2005.
- [10] G. H. Golub and C. F. Van Loan. *Matrix Computations*. The Johns Hopkins University Press, Baltimore, third edition, 1996.
- [11] N. Golyandina, V. Nekrutkin, and A. A. Zhigljavsky. *Analysis of Time Series Structure: SSA and related techniques*. Chapman and Hall/CRC Press, Boca Raton, Florida, USA, 2001.
- [12] N. E. Golyandina and K. D. Usevich. 2D-extension of singular spectrum analysis: algorithm and elements of theory. In V. Olshevsky and E. Tyrtyshnikov, editors, *Matrix Methods: Theory, Algorithms, Applications*. World Scientific Publishing, Singapore, 2010.
- [13] L. Greengard and J. Strain. The fast Gauss transform. *SIAM Journal on Scientific and Statistical Computing*, 12(1):79–94, 1991.
- [14] N. J. Higham. *Functions of Matrices: Theory and Computation*. SIAM, Philadelphia, 2008.

- [15] W. W. Hsieh. Nonlinear multivariate and time series analysis by neural network methods. *Reviews of Geophysics*, 42(1):1–25, 2004.
- [16] W. W. Hsieh and K. Hamilton. Nonlinear singular spectrum analysis of the tropical stratospheric wind. *Quarterly Journal of the Royal Meteorological Society*, 129(592):2367–2382, 2003.
- [17] J. W. Jaromczyk and G. T. Toussaint. Relative neighborhood graphs and their relatives. *Proceedings of the IEEE*, 80(9):1502–1517, 1992.
- [18] I. T. Jolliffe. *Principal Component Analysis*. Springer-Verlag, Berlin, 1986.
- [19] N. Kambhatla and K. T. Leen. Dimension reduction by local principal component analysis. *Neural Computation*, 9(7):1493–1516, 1997.
- [20] M. A. Kramer. Nonlinear principal component analysis using autoassociative neural networks. *AIChE Journal*, 37(2):233–243, 1991.
- [21] M. Loève. *Probability theory: foundations, random sequences*. van Nostrand, Princeton, New Jersey, USA, 1955.
- [22] J. R. Magnus. On differentiating eigenvalues and eigenvectors. *Economic Theory*, 1(2):179–191, 1985.
- [23] J. Miller. *Relative Critical Sets in  $R^n$  and Applications to Image Analysis*. PhD thesis, University of North Carolina, 1998.
- [24] A. H. Monahan. Nonlinear principal component analysis: Tropical Indo–Pacific sea surface temperature and sea level pressure. *Journal of Climate*, 14(2):219–233, 2001.
- [25] S. C. Newbigging, L. A. Mysak, and W. W. Hsieh. Improvements to the nonlinear principal component analysis method, with applications to ENSO and QBO. *Atmosphere-Ocean*, 41(4):291–299, 2003.
- [26] J. M. Ortega. *Numerical Analysis: A Second Course*. SIAM, Philadelphia, 1990.
- [27] U. Ozertem and D. Erdogmus. Locally defined principal curves and surfaces. *Journal of Machine Learning Research*, 12:1249–1286, April 2011.
- [28] K. Pearson. On lines and planes of closest fit to systems of points in space. *Philosophical Magazine Series 6*, 2(11):559–572, 1901.
- [29] S. Pulkkinen. Ridge curve approach to extraction of curvilinear structures from noisy data. TUCS Technical Report TR1082, Turku Centre for Computer Science (<http://tucs.fi>), Turku, Finland, 2013.

- [30] S. Pulkkinen, M. M. Mäkelä, and N. Karmitsa. A generative model and a generalized trust region Newton method for noise reduction. *Computational Optimization and Applications*, 2013. doi: 10.1007/s10589-013-9581-4.
- [31] R. M. Rangayyan. *Biomedical Signal Analysis: A Case-Study Approach*. IEEE Press/Wiley, New York, 2002.
- [32] M. Renardy and R. C. Rogers. An introduction to partial differential equations. In J. E. Marsden, L. Sirovich, and S. S. Antman, editors, *Texts in Applied Mathematics*, 13. Springer-Verlag, New York, second edition, 2004.
- [33] I. Ross. *Nonlinear Dimensionality Reduction Methods in Climate Data Analysis*. PhD thesis, University of Bristol, United Kingdom, 2008.
- [34] I. Ross, P. J. Valdes, and S. Wiggins. ENSO dynamics in current climate models: an investigation using nonlinear dimensionality reduction. *Nonlinear Processes in Geophysics*, 15:339–363, 2008.
- [35] M. Scholz, M. Fraunholz, and J. Selbig. Nonlinear principal component analysis: Neural network models and applications. In A. N. Gorban, B. Kégl, D. C. Wunsch, and A. Y. Zinovyev, editors, *Principal Manifolds for Data Visualization and Dimension Reduction*, volume 58 of *Lecture Notes in Computational Science and Engineering*, pages 44–67. Springer, Berlin Heidelberg, 2008.
- [36] B. Schölkopf, A. Smola, and K.-R. Müller. Kernel principal component analysis. In W. Gerstner, A. Germond, M. Hasler, and J.-D. Nicoud, editors, *Artificial Neural Networks — ICANN’97*, volume 1327 of *Lecture Notes in Computer Science*, pages 583–588. Springer, Berlin Heidelberg, 1997.
- [37] J. B. Tenenbaum, V. de Silva, and J. C. Langford. A global geometric framework for nonlinear dimensionality reduction. *Science*, 290(5500):2319–2323, 2000.
- [38] R. Vautard, P. Yiou, and M. Ghil. Singular-spectrum analysis: A toolkit for short, noisy chaotic signals. *Physica D: Nonlinear Phenomena*, 58(1-4):95–126, 1992.
- [39] B. C. Weare, A. R. Navato, and E. R. Newell. Empirical orthogonal analysis of Pacific sea surface temperatures. *Journal of Physical Oceanography*, 6(5):671–678, 1976.
- [40] K. Weinberger and L. Saul. Unsupervised learning of image manifolds by semidefinite programming. *International Journal of Computer Vision*, 70(1):77–90, 2006.



- [41] E. F. Whittlesey. Fixed points and antipodal points. *The American Mathematical Monthly*, 70(8):807–821, 1963.
- [42] Z. Zhang and H. Zha. Principal manifolds and nonlinear dimensionality reduction via tangent space alignment. *SIAM Journal on Scientific Computing*, 26(1):313–338, 2004.

## A Proofs of Technical Results

In this appendix we give proofs of Theorems 3.2 and 3.3. The proofs are carried out by making the following simplifying assumption that can be made without loss of generality.

**Assumption A.1.** *The sample points  $\mathbf{y}_i$  satisfy the condition  $\sum_{i=1}^n \mathbf{y}_i = \mathbf{0}$ .*

First, we recall the density estimate defined by equations (6) and (7), that is,

$$\hat{p}_h(\mathbf{x}) = \frac{1}{n} \sum_{i=1}^n K_h(\mathbf{x} - \mathbf{y}_i) = \frac{1}{(2\pi)^{\frac{d}{2}} h^d n} \sum_{i=1}^n \exp\left(-\frac{\|\mathbf{x} - \mathbf{y}_i\|^2}{2h^2}\right). \quad (27)$$

The gradient and Hessian of this function are

$$\nabla \hat{p}_h(\mathbf{x}) = -\frac{1}{h^2 n} \sum_{i=1}^n (\mathbf{x} - \mathbf{y}_i) K_{h^2 \mathbf{I}}(\mathbf{x} - \mathbf{y}_i) \quad (28)$$

and

$$\nabla^2 \hat{p}_h(\mathbf{x}) = \frac{1}{h^2 n} \sum_{i=1}^n \left[ \frac{(\mathbf{x} - \mathbf{y}_i)(\mathbf{x} - \mathbf{y}_i)^T}{h^2} - \mathbf{I} \right] K_{h^2 \mathbf{I}}(\mathbf{x} - \mathbf{y}_i), \quad (29)$$

respectively. For the logarithm of  $\hat{p}_h$  we have the formulae

$$\nabla \log \hat{p}_h(\mathbf{x}) = \frac{\nabla \hat{p}_h(\mathbf{x})}{\hat{p}_h(\mathbf{x})} \quad \text{and} \quad \nabla^2 \log \hat{p}_h(\mathbf{x}) = \frac{\nabla^2 \hat{p}_h(\mathbf{x})}{\hat{p}_h(\mathbf{x})} - \frac{\nabla \hat{p}_h(\mathbf{x}) \nabla \hat{p}_h(\mathbf{x})^T}{\hat{p}_h(\mathbf{x})^2}. \quad (30)$$

In what follows, we establish limiting values for the logarithm of the Gaussian kernel density estimate and its derivatives as the bandwidth  $h$  approaches infinity. Furthermore, we show that convergence to these limits is uniform. To this end, we need uniform boundedness of the following functions and existence of a uniform Lipschitz constant.

**Proposition A.1.** *Let  $\hat{p}_h : \mathbb{R}^d \rightarrow \mathbb{R}$  be a Gaussian kernel density estimate, let  $U \subset \mathbb{R}^d$  be compact and let Assumption A.1 be satisfied. Then there exist constants  $M_i > 0$  such that for any  $h_0 > 0$  the functions*

$$\begin{aligned} F_0(\mathbf{x}; h) &= h^2 \log [(2\pi)^{\frac{d}{2}} h^d \hat{p}_h(\mathbf{x})], \\ F_1(\mathbf{x}; h) &= h^2 \nabla \log \hat{p}_h(\mathbf{x}), \\ F_2(\mathbf{x}; h) &= h^4 \nabla^2 \log \hat{p}_h(\mathbf{x}) + h^2 \mathbf{I} \end{aligned}$$

satisfy the conditions

$$\|F_i(\mathbf{x}; h)\| \leq M_i$$

for all  $\mathbf{x} \in U$ ,  $i = 0, 1, 2$  and  $h \geq h_0$ . Furthermore, there exist constants  $L_i > 0$  such that

$$\|\mathbf{F}_i(\mathbf{x}; h) - \mathbf{F}_i(\mathbf{y}; h)\| \leq L_i \|\mathbf{x} - \mathbf{y}\|$$

for all  $\mathbf{x}, \mathbf{y} \in U$ ,  $i = 0, 1, 2$  and  $h \geq h_0$ .

The proof of the above proposition is omitted due to space constraints. The upper bounds  $M_i$  can be obtained by straightforward but rather tedious calculation. Likewise, the Lipschitz constants  $L_i$  can be obtained by showing that the derivatives of the functions  $\mathbf{F}_i(\cdot; h)$  with respect to  $h$  are uniformly bounded for all  $h \geq h_0 > 0$  in any compact set.

The following result that guarantees uniform convergence of a function sequence is a direct extension of the Arzelà-Ascoli theorem (e.g. [32]). The original formulation of the Arzelà-Ascoli theorem is given for equicontinuous functions, but the result of the theorem also holds when the function sequence is Lipschitz continuous with uniformly bounded Lipschitz constant.

**Theorem A.1.** Let  $\{\mathbf{f}_k\}_{k \in \mathbb{N}}$  be a sequence of functions from some set  $\Omega \subset \mathbb{R}^m$  to  $\mathbb{R}^d$  and assume that  $\{\mathbf{f}_k\}$  converges (pointwise) to a limit function  $\mathbf{f}^*$ . If there exists  $M > 0$  such that

$$\|\mathbf{f}_k(\mathbf{x})\| \leq M$$

and  $L > 0$  such that

$$\|\mathbf{f}_k(\mathbf{x}) - \mathbf{f}_k(\mathbf{y})\| \leq L \|\mathbf{x} - \mathbf{y}\|$$

for all  $k \in \mathbb{N}$ ,  $\mathbf{x} \in \Omega$  and  $\mathbf{y} \in \Omega$ , then the sequence  $\{\mathbf{f}_k\}$  converges uniformly to  $\mathbf{f}^*$  in  $\Omega$ .

Now we are ready to prove the following result that will be utilized in the subsequent proofs.

**Lemma A.1.** Let  $\hat{p}_h : \mathbb{R}^d \rightarrow \mathbb{R}$  be a Gaussian kernel density estimate and assume A.1. Then

$$\lim_{h \rightarrow \infty} h^2 \log [(2\pi)^{\frac{d}{2}} h^d \hat{p}_h(\mathbf{x})] = -\frac{1}{2n} \sum_{i=1}^n \|\mathbf{x} - \mathbf{y}_i\|^2, \quad (31)$$

$$\lim_{h \rightarrow \infty} h^2 \nabla \log \hat{p}_h(\mathbf{x}) = -\frac{1}{n} \sum_{i=1}^n (\mathbf{x} - \mathbf{y}_i) = -\mathbf{x} \quad (32)$$

and

$$\lim_{h \rightarrow \infty} [h^4 \nabla^2 \log \hat{p}_h(\mathbf{x}) + h^2 \mathbf{I}] = \frac{1}{n} \sum_{i=1}^n \mathbf{y}_i \mathbf{y}_i^T \quad (33)$$

for all  $\mathbf{x} \in \mathbb{R}^d$ . Furthermore, convergence to these limits is uniform in any compact set.

*Proof.* Let  $\mathbf{x} \in \mathbb{R}^d$ . For the logarithm of a sum of Gaussians we have the limit (the proof is omitted)

$$\lim_{h \rightarrow \infty} h^2 \log \left[ \frac{1}{n} \sum_{i=1}^n \exp \left( -\frac{\|\mathbf{x} - \mathbf{y}_i\|^2}{2h^2} \right) \right] = -\frac{1}{2n} \sum_{i=1}^n \|\mathbf{x} - \mathbf{y}_i\|^2.$$

Consequently, taking the logarithm of (27) and the limit  $h \rightarrow \infty$  and using the above limit yields (31).

Substitution of equations (27) and (28) into (30) yields

$$\nabla \log \hat{p}_h(\mathbf{x}) = \frac{\nabla \hat{p}_h(\mathbf{x})}{\hat{p}_h(\mathbf{x})} = -\frac{1}{h^2} \frac{\sum_{i=1}^n (\mathbf{x} - \mathbf{y}_i) K_{h^2 \mathbf{I}}(\mathbf{x} - \mathbf{y}_i)}{\sum_{i=1}^n K_{h^2 \mathbf{I}}(\mathbf{x} - \mathbf{y}_i)}.$$

The limit (32) follows from multiplying the above expression by  $h^2$ , taking the limit  $h \rightarrow \infty$  and using the limit

$$\lim_{h \rightarrow \infty} \exp\left(-\frac{r^2}{2h^2}\right) = 1 \quad \text{for all } r \in \mathbb{R}. \quad (34)$$

On the other hand, from equations (27)–(29) we obtain that

$$\begin{aligned} \frac{\nabla^2 \hat{p}_h(\mathbf{x})}{\hat{p}_h(\mathbf{x})} &= \frac{1}{h^2} \frac{\sum_{i=1}^n \left[ \frac{(\mathbf{x} - \mathbf{y}_i)(\mathbf{x} - \mathbf{y}_i)^T}{h^2} - \mathbf{I} \right] K_{h^2 \mathbf{I}}(\mathbf{x} - \mathbf{y}_i)}{\sum_{i=1}^n K_{h^2 \mathbf{I}}(\mathbf{x} - \mathbf{y}_i)} \\ &= \frac{1}{h^4} \frac{\sum_{i=1}^n (\mathbf{x} - \mathbf{y}_i)(\mathbf{x} - \mathbf{y}_i)^T K_{h^2 \mathbf{I}}(\mathbf{x} - \mathbf{y}_i)}{\sum_{i=1}^n K_{h^2 \mathbf{I}}(\mathbf{x} - \mathbf{y}_i)} - \frac{\mathbf{I}}{h^2} \end{aligned}$$

and

$$\frac{\nabla \hat{p}_h(\mathbf{x}) \nabla \hat{p}_h(\mathbf{x})^T}{\hat{p}_h(\mathbf{x})^2} = \frac{1}{h^4} \frac{[\sum_{i=1}^n (\mathbf{x} - \mathbf{y}_i) K_{h^2 \mathbf{I}}(\mathbf{x} - \mathbf{y}_i)] [\sum_{i=1}^n (\mathbf{x} - \mathbf{y}_i) K_{h^2 \mathbf{I}}(\mathbf{x} - \mathbf{y}_i)]^T}{[\sum_{i=1}^n K_{h^2 \mathbf{I}}(\mathbf{x} - \mathbf{y}_i)]^2}.$$

Thus, by the limit (34) we obtain

$$\lim_{h \rightarrow \infty} \left[ h^4 \frac{\nabla^2 \hat{p}_h(\mathbf{x})}{\hat{p}_h(\mathbf{x})} + h^2 \mathbf{I} \right] = \frac{1}{n} \sum_{i=1}^n (\mathbf{x} - \mathbf{y}_i)(\mathbf{x} - \mathbf{y}_i)^T \quad (35)$$

and

$$\lim_{h \rightarrow \infty} h^4 \frac{\nabla \hat{p}_h(\mathbf{x}) \nabla \hat{p}_h(\mathbf{x})^T}{\hat{p}_h(\mathbf{x})^2} = \frac{1}{n^2} \left[ \sum_{i=1}^n (\mathbf{x} - \mathbf{y}_i) \right] \left[ \sum_{i=1}^n (\mathbf{x} - \mathbf{y}_i) \right]^T. \quad (36)$$

It follows from Assumption A.1 that

$$\frac{1}{n} \sum_{i=1}^n (\mathbf{x} - \mathbf{y}_i)(\mathbf{x} - \mathbf{y}_i)^T = \mathbf{x}\mathbf{x}^T + \frac{1}{n} \sum_{i=1}^n \mathbf{y}_i \mathbf{y}_i^T \quad (37)$$

and

$$\frac{1}{n^2} \left[ \sum_{i=1}^n (\mathbf{x} - \mathbf{y}_i) \right] \left[ \sum_{i=1}^n (\mathbf{x} - \mathbf{y}_i) \right]^T = \mathbf{x}\mathbf{x}^T. \quad (38)$$

The limit (33) then follows by substituting (37) and (38) into (35) and (36) and using equation (30). Finally, uniform convergence to the limits (31)–(33) follows from Proposition A.1 and Theorem A.1.  $\square$

The following two lemmata facilitate the proof of Theorem 3.2.

**Lemma A.2.** Let  $\hat{p}_h : \mathbb{R}^d \rightarrow \mathbb{R}$  be a Gaussian kernel density estimate and let Assumptions 3.1 and A.1 be satisfied. Denote the eigenvalues of  $\nabla^2 \log \hat{p}_h$  by  $\lambda_1(\cdot; h) \geq \lambda_2(\cdot; h) \geq \dots \geq \lambda_d(\cdot; h)$  and the corresponding eigenvectors by  $\{\mathbf{w}_i(\cdot; h)\}_{i=1}^d$ . Then for any compact set  $U \subset \mathbb{R}^d$  there exists  $h_0 > 0$  such that

$$\lambda_1(\mathbf{x}; h) < 0, \quad (39)$$

$$\lambda_i(\mathbf{x}; h) \neq \lambda_j(\mathbf{x}; h) \quad (40)$$

for all  $\mathbf{x} \in U$ ,  $h \geq h_0$  and  $i, j = 1, 2, \dots, r+1$  such that  $i \neq j$ . Furthermore, if we define

$$\mathbf{W}(\mathbf{x}; h) = [\mathbf{w}_1(\mathbf{x}; h) \quad \mathbf{w}_2(\mathbf{x}; h) \quad \dots \quad \mathbf{w}_r(\mathbf{x}; h)]$$

and

$$\mathbf{V} = [\mathbf{v}_1 \quad \mathbf{v}_2 \quad \dots \quad \mathbf{v}_r],$$

where  $\{\mathbf{v}_i\}_{i=1}^r$  denote the eigenvectors of the matrix  $\hat{\Sigma}_{\mathbf{Y}}$  defined by equation (3) corresponding to its  $r$  greatest eigenvalues, then for all  $\varepsilon > 0$  there exists  $h_0 > 0$  such that

$$\|\mathbf{W}(\mathbf{x}; h) - \mathbf{V}\| < \varepsilon \quad \text{for all } \mathbf{x} \in U \text{ and } h \geq h_0. \quad (41)$$

*Proof.* Let  $\tilde{\lambda}_1 \geq \tilde{\lambda}_2 \geq \dots \geq \tilde{\lambda}_d$  denote the eigenvalues of the matrix  $\hat{\Sigma}_{\mathbf{Y}}$  and let  $\{h_k\}$  be some sequence such that  $\lim_{k \rightarrow \infty} h_k = \infty$ . By uniform convergence to the limit (33) under Assumption A.1 and continuity of eigenvalues of a matrix as a function of its elements (e.g. [26], Theorem 3.1.2), for all  $\varepsilon > 0$  there exists  $k_0$  such that

$$|h_k^4 \lambda_i(\mathbf{x}; h_k) + h_k^2 - \frac{n-1}{n} \tilde{\lambda}_i| < \varepsilon \quad (42)$$

for all  $i = 1, 2, \dots, r+1$ ,  $\mathbf{x} \in U$  and  $k \geq k_0$ . Consequently, condition (39) holds for all  $\mathbf{x} \in U$  for any sufficiently large  $h$  by Assumption 3.1. It also follows from Assumption 3.1, condition (42) and the reverse triangle inequality that for all  $\varepsilon > 0$  and  $i, j = 1, 2, \dots, r+1$  such that  $i \neq j$  and  $|\tilde{\lambda}_i - \tilde{\lambda}_j| > \varepsilon$  there exists  $k_1$  such that

$$h_k^4 |\lambda_i(\mathbf{x}; h_k) - \lambda_j(\mathbf{x}; h_k)| \geq \left| |h_k^4 \lambda_i(\mathbf{x}; h_k) - h_k^2| - |h_k^4 \lambda_j(\mathbf{x}; h_k) - h_k^2| \right| > \frac{n-1}{n} \varepsilon$$

for all  $\mathbf{x} \in U$  and  $k \geq k_1$ . This implies condition (40). Similarly, condition (41) follows from uniform convergence to the limit (33) under Assumption A.1, condition (40) and continuity of eigenvectors as a function of matrix elements when the eigenvalues are distinct (e.g. [26], Theorem 3.1.3).  $\square$

**Lemma A.3.** Assume 3.1 and A.1 and define the function

$$\tilde{\mathbf{W}}(\mathbf{x}; h) = \mathbf{I} - \mathbf{W}(\mathbf{x}; h) \mathbf{W}(\mathbf{x}; h)^T,$$

where the function  $\mathbf{W}$  is defined as in Lemma A.2, and the set  $S_\infty^r$  as in Theorem 3.2. Then the limit

$$\lim_{h \rightarrow \infty} h^2 \|\tilde{\mathbf{W}}(\mathbf{x}; h) \nabla \log \hat{p}_h(\mathbf{x})\| \quad (43)$$

exists for all  $\mathbf{x} \in \mathbb{R}^d$ . Furthermore,  $\mathbf{x} \in S_\infty^r$  if and only if the limit (43) is zero.

*Proof.* By the limits (32) and (41) the limit (43) exists for all  $\mathbf{x} \in \mathbb{R}^d$ . Furthermore, for any  $\mathbf{x} \in \mathbb{R}^d$ , the condition that the limit (43) is zero is equivalent to the condition that

$$\mathbf{v}_i^T \mathbf{x} = 0 \quad \text{for all } i = r+1, r+2, \dots, d,$$

where the vectors  $\mathbf{v}_i$  are defined as in Lemma A.2. By the orthogonality of the vectors  $\mathbf{v}_i$ , the definition of the set  $S_\infty^r$  and Assumption A.1, this condition is equivalent to the condition that  $\mathbf{x} \in S_\infty^r$ .  $\square$

For the proof of Theorem 3.2, we define the set

$$S_h^r = \{\mathbf{x} \in \mathbb{R}^d \mid \|\tilde{\mathbf{W}}(\mathbf{x}; h) \nabla \log \hat{p}_h(\mathbf{x})\| = 0\}, \quad (44)$$

where the function  $\tilde{\mathbf{W}}$  is defined as in Lemma A.3. Under Assumption A.1, we prove both claims of Theorem 3.2 by using the following two lemmata.

**Lemma A.4.** *Let  $U \subset \mathbb{R}^d$  be a compact set such that  $U \cap S_\infty^r \neq \emptyset$  for some  $0 \leq r < d$ . If Assumptions 3.1 and A.1 are satisfied, then for all  $\varepsilon > 0$  there exists  $h_0 > 0$  such that*

$$\sup_{\mathbf{x} \in S_h^r \cap U} \inf_{\mathbf{y} \in S_\infty^r} \|\mathbf{x} - \mathbf{y}\| < \varepsilon \quad \text{for all } h \geq h_0. \quad (45)$$

*Proof.* The proof is by contradiction. Let  $0 \leq r < d$  and let  $U \subset \mathbb{R}^d$  be a compact set such that  $U \cap S_\infty^r \neq \emptyset$ . Assume that there exists  $\varepsilon_1 > 0$  such that for all  $h_0 > 0$  there exists  $h \geq h_0$  such that condition (45) is not satisfied. This implies that for all  $h_0 > 0$  there exists  $h \geq h_0$  such that

$$\inf_{\mathbf{y} \in S_\infty^r} \|\mathbf{x} - \mathbf{y}\| \geq \varepsilon_1 \quad \text{for some } \mathbf{x} \in S_h^r \cap U. \quad (46)$$

Let  $\{\mathbf{x}_k\}$  denote a sequence of such points  $\mathbf{x}$  with the corresponding sequence  $h_k$ . Since the set  $S_h^r \cap U$  is compact by the compactness of  $U$  and the continuity of  $\tilde{\mathbf{W}}(\cdot, h)$  in  $U$  for any sufficiently large  $h$ , the sequence  $\{\mathbf{x}_k\}$  has a convergent subsequence  $\{z_k\}$  whose limit point we shall denote as  $z^*$ . Clearly  $z^* \notin S_\infty^r$  by condition (46). Thus, by Lemma A.3 we deduce that for some  $c > 0$ ,

$$\lim_{k \rightarrow \infty} \|\mathbf{F}(z^*; h_k)\| = c, \quad \text{where } \mathbf{F}(\mathbf{x}; h) = h^2 \tilde{\mathbf{W}}(\mathbf{x}; h) \nabla \log \hat{p}_h(\mathbf{x}).$$

In view of the definition (44), the above limit implies that there exists  $\varepsilon_2 > 0$  and  $k_0$  such that for all  $k \geq k_0$ ,

$$\|\mathbf{F}(z^*; h_k) - \mathbf{F}(\mathbf{y}; h_k)\| \geq \varepsilon_2 \quad \text{for all } \mathbf{y} \in S_{h_k}^r \cap U. \quad (47)$$

On the other hand, if we define the function  $\mathbf{F}^*(\mathbf{x}) = -(\mathbf{I} - \mathbf{V}\mathbf{V}^T)\mathbf{x}$ , the triangle inequality yields

$$\begin{aligned} \|\mathbf{F}(z^*; h_k) - \mathbf{F}(\mathbf{y}; h_k)\| &\leq \|\mathbf{F}(z^*; h_k) - \mathbf{F}^*(z^*)\| + \|\mathbf{F}(\mathbf{y}; h_k) - \mathbf{F}^*(z^*)\| \\ &\leq \|\mathbf{F}(z^*; h_k) - \mathbf{F}^*(z^*)\| + \|\mathbf{F}(\mathbf{y}; h_k) - \mathbf{F}^*(\mathbf{y})\| + \|\mathbf{F}^*(\mathbf{y}) - \mathbf{F}^*(z^*)\|. \end{aligned}$$

Combining this with the inequality

$$\|\mathbf{F}^*(\mathbf{y}) - \mathbf{F}^*(\mathbf{z}^*)\| \leq \|\mathbf{I} - \mathbf{V}\mathbf{V}^T\| \|\mathbf{y} - \mathbf{z}^*\| = \|\mathbf{y} - \mathbf{z}^*\|$$

and noting the convergence of  $\mathbf{F}(\cdot; h_k)$  to the function  $\mathbf{F}^*$  (that is uniform in  $U$ ) as  $k \rightarrow \infty$  (by Lemmata A.1 and A.2), we deduce from (47) that for all  $\varepsilon_2 > \varepsilon_3 > 0$  there exists  $k_1$  such that

$$\|\mathbf{z}^* - \mathbf{y}\| + \varepsilon_3 \geq \|\mathbf{F}(\mathbf{z}^*; h_k) - \mathbf{F}(\mathbf{y}; h_k)\| \geq \varepsilon_2 \quad (48)$$

for all  $\mathbf{y} \in S_{h_k}^r \cap U$  and  $k \geq k_1$ .

Condition (48) implies that for all  $0 < \varepsilon_3 < \varepsilon_2$  there exists  $k_1$  such that

$$\inf_{\mathbf{y} \in S_{h_k}^r \cap U} \|\mathbf{z}^* - \mathbf{y}\| \geq \varepsilon_2 - \varepsilon_3 \quad \text{for all } k \geq k_1. \quad (49)$$

On the other hand, for all  $\varepsilon > 0$  we have  $\mathbf{z}_k \in B(\mathbf{z}^*; \varepsilon)$  for any sufficiently large  $k$  due to the assumption that  $\mathbf{z}_k$  converges to  $\mathbf{z}^*$ . If we choose  $0 < \varepsilon < \varepsilon_2$ , then the sequence  $\{\mathbf{x}_k\}$ , whose subsequence is  $\{\mathbf{z}_k\}$ , has an element  $\mathbf{x}_k \notin S_{h_k}^r \cap U$  for some  $k$  by condition (49). This leads to a contradiction with the construction of the sequence  $\{\mathbf{x}_k\}$ , which states that  $\mathbf{x}_k \in S_{h_k}^r \cap U$  for all  $k$ .  $\square$

**Lemma A.5.** *Let  $\hat{p}_h$  be a Gaussian kernel density estimate, let  $0 \leq r < d$ , let Assumptions 3.1 and A.1 be satisfied and define the set  $S_\infty^r$  as in Theorem 3.2. Then for any compact set  $U \subset \mathbb{R}^d$  such that  $U \cap S_\infty^r \neq \emptyset$  and  $\varepsilon > 0$  there exists  $h_0 > 0$  such that*

$$\sup_{\mathbf{x} \in S_\infty^r \cap U} \inf_{\mathbf{y} \in \mathcal{R}_{\log \hat{p}_h}^r} \|\mathbf{x} - \mathbf{y}\| < \varepsilon \quad \text{for all } h \geq h_0. \quad (50)$$

*Proof.* Let  $0 \leq r < d$  and let  $\{\mathbf{v}_i\}_{i=r+1}^d$  denote a set of orthonormal eigenvectors of the matrix  $\hat{\Sigma}_{\mathbf{Y}}$  corresponding to the  $d - r$  smallest eigenvalues. The vectors  $\{\mathbf{v}_i\}_{i=r+1}^d$  are uniquely determined up to the choice of basis because the eigenvectors  $\{\mathbf{v}_i\}_{i=1}^r$  spanning their orthogonal complement are uniquely determined by Assumption 3.1. Define the sets

$$D_{\mathbf{x}, \varepsilon} = \left\{ \mathbf{x} + \sum_{i=r+1}^d \alpha_{i-r} \mathbf{v}_i \mid \sum_{i=1}^r \alpha_i^2 \leq \varepsilon^2 \right\}$$

and

$$D_\varepsilon = \bigcup_{\mathbf{x} \in S_\infty^r \cap U} D_{\mathbf{x}, \varepsilon}$$

for some orthonormal eigenvectors  $\{\mathbf{v}_i\}_{i=r+1}^d$  spanning the orthogonal complement of  $\text{span}(\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_r)$ .

Let  $\{\mathbf{u}_i(\cdot; h)\}_{i=1}^d$  denote a set of orthonormal vectors that are orthogonal to the eigenvectors  $\{\mathbf{w}_i(\cdot; h)\}_{i=1}^r$  of  $\nabla^2 \log \hat{p}_h$  corresponding to the  $r$  greatest eigenvalues. Define the functions

$$\mathbf{F}(\mathbf{x}; h) = h^2 \mathbf{U}(\mathbf{x}; h)^T \nabla \log \hat{p}_h(\mathbf{x})$$

and

$$\tilde{\mathbf{F}}_{\mathbf{x}_0}(\mathbf{y}; h) = h^2 \mathbf{U}(\bar{\mathbf{V}}\mathbf{y} + \mathbf{x}_0; h)^T \nabla \log \hat{p}_h(\bar{\mathbf{V}}\mathbf{y} + \mathbf{x}_0),$$

where

$$\mathbf{U}(\mathbf{x}; h) = [\mathbf{u}_1(\mathbf{x}; h) \quad \mathbf{u}_2(\mathbf{x}; h) \quad \cdots \quad \mathbf{u}_{d-r}(\mathbf{x}; h)]$$

and  $\bar{\mathbf{V}} = [\mathbf{v}_{r+1} \quad \mathbf{v}_{r+2} \quad \cdots \quad \mathbf{v}_d]$  assuming that the orientation is chosen so that  $\det(\bar{\mathbf{V}}) = 1$ . To fix the orientation of the vectors  $\mathbf{u}_i(\mathbf{x}; h)$ , we impose the constraint

$$\mathbf{U}(\mathbf{x}; h) = \arg \min_{\mathbf{U}' \in Q_{\mathbf{x}, h}} \|\mathbf{U}' - \bar{\mathbf{V}}\|_F. \quad (51)$$

Here  $\|\cdot\|_F$  denotes the Frobenius norm,

$$Q_{\mathbf{x}, h} = \{\mathbf{U}' \in O(d, d-r) \mid \mathbf{U}'^T \mathbf{W}(\mathbf{x}; h) = \mathbf{0}, \quad \det(\mathbf{U}') = 1\},$$

$O(d, d-r)$  denotes a  $d \times (d-r)$  matrix having orthonormal columns and the matrix  $\mathbf{W}(\mathbf{x}; h)$  is defined as in Lemma A.2. It can be shown that the function  $\mathbf{U}(\cdot; h)$  is well-defined for any  $h > 0$ .<sup>1</sup> Spanning the orthogonal complement of the columns of  $\mathbf{W}(\cdot; h)$ , the columns of  $\mathbf{U}(\cdot; h)$  are also continuous in a given compact set when  $\mathbf{W}(\cdot; h)$  is continuous. That is, when condition (40) is satisfied in such a set by Lemma A.2.

The above definitions and condition (41) in the compact set  $D_\varepsilon$  imply that for all  $\varepsilon_1, \varepsilon_2 > 0$  there exists  $h_0 > 0$  such that

$$\|\mathbf{U}(\mathbf{x}; h) - \bar{\mathbf{V}}\| < \varepsilon_2 \quad \text{for all } \mathbf{x} \in D_{\varepsilon_1} \text{ and } h \geq h_0.$$

Consequently, uniform convergence to the limit (32) as  $h \rightarrow \infty$  by Lemma A.3 together with the property that

$$\bar{\mathbf{V}}^T (\bar{\mathbf{V}} \mathbf{y} + \mathbf{x}_0) = \mathbf{y} \quad \text{for all } \mathbf{y} \in \mathbb{R}^{d-r} \text{ and } \mathbf{x}_0 \in S_\infty^r$$

following from Assumption A.1 implies that for all  $\varepsilon_1, \varepsilon_2 > 0$  there exists  $h_0 > 0$  such that

$$\|\tilde{\mathbf{F}}_{\mathbf{x}_0}(\mathbf{y}; h) - (-\mathbf{y})\| < \varepsilon_2 \quad \text{for all } \mathbf{x}_0 \in S_\infty^r \cap U, \mathbf{y} \in \tilde{D}_{\varepsilon_1} \text{ and } h \geq h_0,$$

where  $\tilde{D}_\varepsilon = \{\mathbf{y} \in \mathbb{R}^{d-r} \mid \|\mathbf{y}\| \leq \varepsilon\}$ .

By the above condition, for any  $0 < \varepsilon_2 < \varepsilon_1$  there exists  $h_0 > 0$  such that for all  $h \geq h_0$  and  $\mathbf{x}_0 \in S_\infty^r \cap U$  we have  $-\tilde{\mathbf{F}}_{\mathbf{x}_0}(\mathbf{y}; h)^T \mathbf{y} > 0$  for all  $\mathbf{y} \in \partial \tilde{D}_{\varepsilon_1}$ , where  $\partial$  denotes the boundary of a set. On the other hand,  $-\mathbf{y}$  is the inward-pointing normal vector of the disk  $\tilde{D}_{\varepsilon_1}$  at any  $\mathbf{y} \in \partial \tilde{D}_{\varepsilon_1}$ . Together with the continuity of  $\tilde{\mathbf{F}}_{\mathbf{x}_0}(\cdot; h)$  in  $\tilde{D}_{\varepsilon_1}$  when  $h$  is sufficiently large, the well-known results from topology (e.g. [41]) then imply that  $\tilde{\mathbf{F}}_{\mathbf{x}_0}(\cdot; h)$  has at least one zero point  $\mathbf{y}^*$  in the interior of  $\tilde{D}_{\varepsilon_1}$  for all  $\mathbf{x}_0 \in S_\infty^r \cap U$  and  $h \geq h_0$ . Clearly, for any such  $\mathbf{y}^*$  and  $\mathbf{x}_0$  the point  $\mathbf{x}^* = \bar{\mathbf{V}} \mathbf{y}^* + \mathbf{x}_0$  lies in the set  $D_{\mathbf{x}_0, \varepsilon}$  and  $\mathbf{F}(\mathbf{x}^*; h) = \tilde{\mathbf{F}}_{\mathbf{x}_0}(\mathbf{y}^*; h) = \mathbf{0}$ .

From the above we conclude that for all  $\varepsilon > 0$  there exists  $h_0 > 0$  such that for all  $\mathbf{x}_0 \in S_\infty^r \cap U$  condition (5a) holds for  $\log \hat{p}_h$  at least at one point in  $D_{\mathbf{x}_0, \varepsilon}$  for all  $h \geq h_0$ . On the other hand, for all  $\varepsilon > 0$  conditions (5b) and (5c) are satisfied in the compact set  $D_\varepsilon$  for all sufficiently large  $h$  by conditions (39) and (40). Hence, we have proven that for all  $\varepsilon > 0$  condition (50) holds for all sufficiently large  $h$ .  $\square$

<sup>1</sup>Problem (51) can be equivalently formulated as an *orthogonal Procrustes problem*. With the matrices defined above, this problem has a unique solution (e.g. [14]).

**Proof of Theorem 3.2** (page 7). Follows directly from Lemmata A.4 and A.5 by the property that  $\mathcal{R}_{\hat{p}_h}^r = \mathcal{R}_{\log \hat{p}_h}^r \subseteq S_h^r$  for all  $0 \leq r < d$  and  $h > 0$  by Lemma 3.1 and Definition 3.1.  $\square$

Next, we prove Theorem 3.3 under Assumption A.1 by using the following lemma.

**Lemma A.6.** Let  $\hat{p}_h$  be a Gaussian kernel density estimate, assume A.1 and define the set

$$U_h = \bigcup_{i=1}^n \mathcal{L}_i^h, \quad \text{where } \mathcal{L}_i^h = \{\mathbf{x} \in \mathbb{R}^d \mid \log \hat{p}_h(\mathbf{x}) \geq \log \hat{p}_h(\mathbf{y}_i)\}.$$

Then for some  $r > \max_{i=1,2,\dots,n} \|\mathbf{y}_i\|$  there exists  $h_0 > 0$  such that  $U_h \subseteq B(\mathbf{0}; r)$  for all  $h \geq h_0$ .

*Proof.* The proof is by contradiction. Assume that for all  $r > r_0 = \max_{i=1,2,\dots,n} \|\mathbf{y}_i\|$  and  $h_0 > 0$  there exists  $h \geq h_0$  such that  $\mathbf{x} \in U_h \setminus B(\mathbf{0}; r)$ . Let  $\{\mathbf{x}_k\}$ ,  $\{r_k\}$  and  $\{h_k\}$  denote sequences satisfying these properties such that  $\{r_k\}$  and  $\{h_k\}$  are monotonously increasing. This implies that

$$\|\mathbf{x}_k\| > r_k > r_0 = \max_{i=1,2,\dots,n} \|\mathbf{y}_i\| \quad \text{for all } k \geq k_0 \quad (52)$$

and also that for all  $k \geq k_0$ ,

$$\log \hat{p}_{h_k}(\mathbf{x}_k) \geq \log \hat{p}_{h_k}(\mathbf{y}_j) \quad \text{for some } j \in \{1, 2, \dots, n\}. \quad (53)$$

By Assumption A.1 and condition (52) we have that  $\|\mathbf{x}_k - \mathbf{y}_i\| \geq \|\mathbf{x}_k\| - r_0$  for all  $k \geq k_0$  and  $i = 1, 2, \dots, n$ . Consequently,

$$h_k^2 \log \left[ \frac{1}{n} \sum_{i=1}^n \exp \left( -\frac{\|\mathbf{x}_k - \mathbf{y}_i\|^2}{2h_k^2} \right) \right] \leq h_k^2 \log \left[ \exp \left( -\frac{(\|\mathbf{x}_k\| - r_0)^2}{2h_k^2} \right) \right] = -\frac{(\|\mathbf{x}_k\| - r_0)^2}{2}$$

for all  $k \geq k_0$ . By equation (27), this implies that

$$h_k^2 [\log \hat{p}_{h_k}(\mathbf{x}_k) + \log [(2\pi)^{\frac{d}{2}} h_k^d]] \leq -\frac{(\|\mathbf{x}_k\| - r_0)^2}{2} \quad (54)$$

for all  $k \geq k_0$ . On the other hand, by the limit (31), Assumption A.1 and the choice of  $r_0$  we have

$$\begin{aligned} \lim_{k \rightarrow \infty} h_k^2 [\log \hat{p}_{h_k}(\mathbf{y}_j) + \log [(2\pi)^{\frac{d}{2}} h_k^d]] &= -\frac{1}{2n} \sum_{i=1}^n \|\mathbf{y}_j - \mathbf{y}_i\|^2 \\ &= -\frac{1}{2n} \left( \sum_{i=1}^n \|\mathbf{y}_j\|^2 - 2 \sum_{i=1}^n \mathbf{y}_j^T \mathbf{y}_i + \sum_{i=1}^n \|\mathbf{y}_i\|^2 \right) \geq -r_0^2 \end{aligned} \quad (55)$$

for all  $j = 1, 2, \dots, n$ . Plugging the limits (54) and (55) into inequality (53) then leads to a contradiction for any sufficiently large  $k$  since  $\lim_{k \rightarrow \infty} \|\mathbf{x}_k\| = \infty$  by condition (52) and the assumption that the sequence  $\{r_k\}$  is monotonously increasing.  $\square$

**Proof of Theorem 3.3** (page 11). By Lemma A.6 there exists  $r > \max_{i=1,2,\dots,n} \|\mathbf{y}_i\|$  such that  $U_h \subseteq B(\mathbf{0}; r)$  for all sufficiently large  $h$ . Thus, condition (20) for all  $\mathbf{x} \in U_h$  and such  $h$  follows from Assumption 3.1, compactness of the set  $B(\mathbf{0}; r)$  and Lemma A.2. Finally, connectedness of the set  $U_h$  for all sufficiently large  $h$  follows from the strict concavity of  $\log \hat{p}_h$  in  $B(\mathbf{0}; r) \supseteq U_h$  by condition (39).  $\square$





TURKU  
CENTRE *for*  
COMPUTER  
SCIENCE

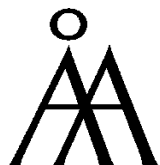
Joukahaisenkatu 3-5 A, 20520 TURKU, Finland | [www.tucs.fi](http://www.tucs.fi)



**University of Turku**

*Faculty of Mathematics and Natural Sciences*

- Department of Information Technology
  - Department of Mathematics
- Turku School of Economics*
- Institute of Information Systems Sciences



**Abo Akademi University**

- Department of Computer Science
- Institute for Advanced Management Systems Research

ISBN 978-952-12-2970-1

ISSN 1239-1891