
Estimating Density Ridges by Direct Estimation of Density-Derivative-Ratios

Hiroaki Sasaki^{1,4}
hsasaki@is.naist.jp

Takafumi Kanamori^{2,4}
kanamori@is.nagoya-u.ac.jp

Masashi Sugiyama^{3,4}
sugi@k.u-tokyo.ac.jp

¹ Graduate School of Information Science, Nara Institute of Science and Technology, Nara, Japan

² Department of Computer Science and Mathematical Informatics, Nagoya University, Aichi, Japan

³ Graduate School of Frontier Sciences, The University of Tokyo, Tokyo, Japan

⁴ Center for Advanced Intelligence Project, RIKEN, Tokyo, Japan

Abstract

Estimation of *density ridges* has been gathering a great deal of attention since it enables us to reveal lower-dimensional structures hidden in data. Recently, *subspace constrained mean shift* (SCMS) was proposed as a practical algorithm for density ridge estimation. A key technical ingredient in SCMS is to accurately estimate the ratios of the density derivatives to the density. SCMS takes a three-step approach for this purpose — first estimating the data density, then computing its derivatives, and finally taking their ratios. However, this three-step approach can be unreliable because a good density estimator does not necessarily mean a good density derivative estimator and division by an estimated density could significantly magnify the estimation error. To overcome these problems, we propose a novel method that directly estimates the ratios without going through density estimation and division. Our proposed estimator has an analytic-form solution and it can be computed efficiently. We further establish a non-parametric convergence bound for the proposed ratio estimator. Finally, based on this direct ratio estimator, we develop a practical algorithm for density ridge estimation and experimentally demonstrate its usefulness on a variety of datasets.

1 Introduction

Estimating the *ridge* of the data density possesses a wide range of real-world applications, including estimation of filamentary structures formed by galaxies in cosmology [Chen et al., 2016], extraction of curvilinear structures (e.g., blood vessels in eye balls) in medical imaging [You et al., 2011], skeletonization of optical characters for feature extraction and compression [Kégl and Krzyżak, 2002], traffic pattern analysis [Einbeck and Dwyer, 2011], and shape analysis in computer vision [Su et al., 2013] (see Pulkkinen [2015] for more applications). For this reason, density ridge estimation has been gathering a great deal of attention recently [Ozertem and Erdogmus, 2011, Genovese et al., 2014, Chen et al., 2015b,a, Ghassabeh et al., 2013].

Extending the classical concept of *principal curves* [Hastie, 1984, Hastie and Stuetzle, 1989], a practical algorithm for density ridge estimation called *subspace constrained mean shift* (SCMS) [Ozertem and Erdogmus, 2011] has been proposed recently. SCMS is essentially a projected gradient ascent algorithm over an estimated data density: At each iteration, a gradient vector of the estimated density is computed as in *mean shift* clustering [Fukunaga and Hostetler, 1975, Cheng, 1995, Comaniciu and Meer, 2002], and then it is projected to the subspace orthogonal to the density ridge. Along this projected gradient vector, data points are updated toward the density ridge until they converge. See Genovese et al. [2014], Chen et al. [2015b], Chen et al. [2015a] and Ghassabeh et al. [2013] for theoretical properties of SCMS.

Technically, the key ingredients to obtain such projected gradient vectors are *the ratios of density derivatives to the density*. SCMS takes a three-step approach to estimate the ratios: First, estimate the data density

by *kernel density estimation*, then compute its derivatives, and finally take their ratios. However, this three-step approach can be unreliable because a good density estimator does not necessarily mean a good density derivative estimator, and division by an estimated density could significantly magnify the estimation error.

To cope with these problems, we propose a novel estimator called the *least-squares density-derivative-ratios* (LSDDR). In stark contrast to the three-step approach in SCMS, LSDDR neither performs density estimation nor involves division by an estimated quantity; rather, it directly estimates the *density-derivative-ratios*. Previously, a method to directly estimate the log-density derivatives (which is equal to the ratio of the first-order density derivative to the density) has been proposed [Cox, 1985, Sasaki et al., 2014]. LSDDR can be regarded as its generalization to higher-order derivatives. LSDDR has an analytic-form solution and it can be computed efficiently. Furthermore, we establish a non-parametric convergence bound for LSDDR.

Based on this LSDDR, we then develop a new algorithm for density ridge estimation called the *least-squares density ridge finder* (LSDRF). We experimentally demonstrate the advantages of LSDRF over SCMS on a variety of datasets.

2 Density Ridge Estimation

In this section, we formulate the problem of density ridge estimation and review an existing method.

2.1 Problem Formulation

Suppose that independent and identically distributed samples $\mathcal{X} = \{\mathbf{x}_i \mid \mathbf{x}_i = (x_i^{(1)}, x_i^{(2)}, \dots, x_i^{(D)})^\top\}_{i=1}^n$ drawn from an unknown probability distribution with density $p(\mathbf{x})$ are available, where $^\top$ denotes the transpose. For positive integer d such that $d < D$, the goal is to estimate from \mathcal{X} the d -dimensional *density ridge* [Eberly, 1996, Ozertem and Erdogmus, 2011, Genovese et al., 2014], which is defined as a collection of points satisfying

$$\mathcal{R} = \{\mathbf{x} \in \mathbb{R}^D \mid \mathbf{V}(\mathbf{x})\mathbf{V}(\mathbf{x})^\top \nabla p(\mathbf{x}) = 0, \eta_{d+1}(\mathbf{x}) < 0\}, \quad (1)$$

where ∇ denotes the differential operator w.r.t. \mathbf{x} , $\mathbf{V}(\mathbf{x}) = (\mathbf{v}_{d+1}, \dots, \mathbf{v}_D)$, and \mathbf{v}_i is the eigenvector associated with the eigenvalue $\eta_i(\mathbf{x})$ of the Hessian matrix $\nabla \nabla p(\mathbf{x})$. We assume that the eigenvalues are sorted in descending order such that $\eta_1(\mathbf{x}) \geq \eta_2(\mathbf{x}) \geq \dots \geq \eta_D(\mathbf{x})$.

2.2 Subspace Constrained Mean Shift (SCMS)

According to definition (1), a practical algorithm for finding density ridges called *subspace constrained mean shift* (SCMS) was proposed in Ozertem and Erdogmus [2011]. SCMS is based on *mean shift* (MS) clustering [Fukunaga and Hostetler, 1975, Cheng, 1995, Comaniciu and Meer, 2002] and the *inverse local-covariance matrix*.

MS is a mode-seeking clustering method based on kernel density estimation (KDE):

$$\hat{p}_{\text{KDE}}(\mathbf{x}) = \frac{1}{nZ} \sum_{i=1}^n K_{\text{KDE}} \left(\frac{\|\mathbf{x} - \mathbf{x}_i\|^2}{h^2} \right),$$

where Z is the normalizing constant, K_{KDE} is a kernel function for KDE, and h denotes its bandwidth. MS updates data points toward the nearest mode (i.e., a local maximum) of $\hat{p}_{\text{KDE}}(\mathbf{x})$ by

$$\mathbf{x} \leftarrow \mathbf{x} + \widehat{\mathbf{m}}_{\text{MS}}(\mathbf{x}) = \frac{\sum_{i=1}^n \mathbf{x}_i G_{\text{KDE}} \left(\frac{\|\mathbf{x} - \mathbf{x}_i\|^2}{h^2} \right)}{\sum_{i=1}^n G_{\text{KDE}} \left(\frac{\|\mathbf{x} - \mathbf{x}_i\|^2}{h^2} \right)},$$

where $G_{\text{KDE}}(t) = -K'_{\text{KDE}}(t) = -\frac{d}{dt} K_{\text{KDE}}(t)$. $\widehat{\mathbf{m}}_{\text{MS}}(\mathbf{x})$ is called the *mean shift vector* which is shown to be parallel to $\nabla \hat{p}_{\text{KDE}}(\mathbf{x})$ [Comaniciu and Meer, 2002]:

$$\begin{aligned} \nabla \hat{p}_{\text{KDE}}(\mathbf{x}) &= \frac{2}{nh^2Z} \sum_{i=1}^n (\mathbf{x}_i - \mathbf{x}) G_{\text{KDE}} \left(\frac{\|\mathbf{x} - \mathbf{x}_i\|^2}{h^2} \right) \\ &= \alpha(\mathbf{x}) \widehat{\mathbf{m}}_{\text{MS}}(\mathbf{x}), \end{aligned}$$

where $\alpha(\mathbf{x}) = \frac{2}{nh^2Z} \sum_{i=1}^n G_{\text{KDE}} \left(\frac{\|\mathbf{x} - \mathbf{x}_i\|^2}{h^2} \right)$. The above equation indicates that MS performs gradient ascent with adaptive step size $1/\alpha(\mathbf{x})$.

The basic idea of SCMS is to perform MS-like gradient ascent on the subspace which is orthogonal to the density ridge. SCMS obtains such a subspace as the span of the eigenvectors of the Hessian matrix of the log-density, which is called the *inverse local-covariance matrix* [Ozertem and Erdogmus, 2011]:

$$\begin{aligned} \Sigma^{-1}(\mathbf{x}) &= -\nabla \nabla \log p(\mathbf{x}) \\ &= -\frac{\nabla \nabla p(\mathbf{x})}{p(\mathbf{x})} + \frac{\nabla p(\mathbf{x}) \nabla p(\mathbf{x})^\top}{p(\mathbf{x})^2}. \quad (2) \end{aligned}$$

As theoretically shown in Genovese et al. [2014], the log-density is used instead of the (non-log) density because it has some advantages: In practice, a projector to the subspace is obtained by applying *principal component analysis* (PCA) to the estimated inverse local-covariance matrix $\widehat{\Sigma}_{\text{KDE}}^{-1}(\mathbf{x})$ where $p(\mathbf{x})$ in (2) is

Input: A data point \mathbf{x} .

Step 1 Initialize $t = 0$ and $\mathbf{y}(t) = \mathbf{x}$.

Step 2 Evaluate the mean shift vector $\widehat{\mathbf{m}}_{\text{MS}}(\mathbf{y}(t))$.

Step 3 Evaluate the inverse local-covariance matrix using $\widehat{p}_{\text{KDE}}(\mathbf{y}(t))$.

Step 4 Construct a projector by applying PCA to $\widehat{\Sigma}_{\text{KDE}}^{-1}(\mathbf{x})$.

Step 5 Update $\mathbf{y}(t)$ as $\mathbf{y}(t+1) = \mathbf{y}(t) + \widehat{\mathbf{V}}_{\text{KDE}} \widehat{\mathbf{V}}_{\text{KDE}}^\top \widehat{\mathbf{m}}_{\text{MS}}(\mathbf{y}(t))$.

Step 6 Stop if $\|\mathbf{y}(t+1) - \mathbf{y}(t)\|_2 < \epsilon$. Otherwise, $t \leftarrow t+1$ and go back to Step 2.

Output: $\widehat{\mathbf{y}} = \mathbf{y}(t)$

Figure 1: An algorithm of subspace constrained mean shift [Ozertem and Erdogmus, 2011, Ghassabeh et al., 2013]. $\|\cdot\|_2$ is the ℓ_2 norm and ϵ denotes a small positive constant. As an input data point, a data sample itself is typically used.

replaced with $\widehat{p}_{\text{KDE}}(\mathbf{x})$. Then, the projected gradient update rule of SCMS is given as

$$\mathbf{x} \leftarrow \mathbf{x} + \widehat{\mathbf{V}}_{\text{KDE}} \widehat{\mathbf{V}}_{\text{KDE}}^\top \widehat{\mathbf{m}}_{\text{MS}}(\mathbf{x}). \quad (3)$$

The entire algorithm of SCMS is summarized in Figure 1. The convergence of the SCMS algorithm is proved in Ghassabeh et al. [2013].

In SCMS, one of the key challenges is to accurately estimate the inverse local covariance matrix (2). SCMS takes a three-step approach, i.e., estimate $p(\mathbf{x})$ by KDE, compute its derivatives, and plug them into (2). However, this approach can perform poorly because a good density estimator does not necessarily mean a good density derivative estimator. In addition, division by an estimated density could significantly magnify the estimation error. A more appropriate way would be to directly estimate the ratios in (2) without going through density estimation and division. Following this idea, we next propose a novel direct estimator of the ratios.

3 Direct Estimation of Density-Derivative-Ratios

In this section, we propose a direct estimator of density-derivative-ratios. Then, the estimator is theoretically analyzed.

3.1 Least-Squares Density-Derivative-Ratios

Here, our tentative goal is to estimate the ratio of the k -th order partial derivative of $p(\mathbf{x})$ to $p(\mathbf{x})$ from $\mathcal{X} = \{\mathbf{x}_i\}_{i=1}^n$:

$$\frac{\partial^{k,\mathbf{j}} p(\mathbf{x})}{p(\mathbf{x})}, \quad (4)$$

where $\partial^{k,\mathbf{j}} = \frac{\partial^k}{\partial(x^{(1)})^{j_1} \partial(x^{(2)})^{j_2} \dots \partial(x^{(D)})^{j_D}}$, $\mathbf{j} = (j_1, j_2, \dots, j_D)$ and $j_1 + j_2 + \dots + j_D = k$ for $j_i \in \{0, 1, \dots, k\}$. For instance, when $k = 1$ (or $k = 2$), $\partial^{k,\mathbf{j}} p(\mathbf{x})/p(\mathbf{x})$ is a single element of $\nabla p(\mathbf{x})/p(\mathbf{x})$ (or of $\nabla \nabla p(\mathbf{x})/p(\mathbf{x})$).

Our main idea is to directly fit a model $r_{k,\mathbf{j}}(\mathbf{x})$ to $\frac{\partial^{k,\mathbf{j}} p(\mathbf{x})}{p(\mathbf{x})}$ under the squared-loss:

$$\begin{aligned} & J_{k,\mathbf{j}}(r_{k,\mathbf{j}}) \\ &= \int \left\{ r_{k,\mathbf{j}}(\mathbf{x}) - \frac{\partial^{k,\mathbf{j}} p(\mathbf{x})}{p(\mathbf{x})} \right\}^2 p(\mathbf{x}) d\mathbf{x} - C_{k,\mathbf{j}} \\ &= \int \{r_{k,\mathbf{j}}(\mathbf{x})\}^2 p(\mathbf{x}) d\mathbf{x} - 2 \int r_{k,\mathbf{j}}(\mathbf{x}) \partial^{k,\mathbf{j}} p(\mathbf{x}) d\mathbf{x}, \end{aligned} \quad (5)$$

where $C_{k,\mathbf{j}} = \int \left\{ \frac{\partial^{k,\mathbf{j}} p(\mathbf{x})}{p(\mathbf{x})} \right\}^2 p(\mathbf{x}) d\mathbf{x}$. The first term in (5) can be naively estimated from samples, but it seems challenging to estimate the second term because it includes the derivative of the unknown density. However, similarly to Sasaki et al. [2015], repeatedly applying *integration by parts* allows us to transform the second term as

$$\begin{aligned} & \int r_{k,\mathbf{j}}(\mathbf{x}) \{ \partial^{k,\mathbf{j}} p(\mathbf{x}) \} d\mathbf{x} \\ &= (-1)^k \int \{ \partial^{k,\mathbf{j}} r_{k,\mathbf{j}}(\mathbf{x}) \} p(\mathbf{x}) d\mathbf{x}, \end{aligned} \quad (6)$$

where we assumed that as for all j , $|x^{(j)}| \rightarrow \infty$, the product of $\partial^{k_1, \mathbf{j}_1} r_{k_1, \mathbf{j}_1}(\mathbf{x})$ and $\partial^{k_2, \mathbf{j}_2} p(\mathbf{x})$ approaches zero for any pairs of k_1 and k_2 satisfying $k_1 + k_2 = k - 1$. As a result, the right-hand side of (6) can be easily estimated from samples. Then, an empirical version of (5) is given by

$$\tilde{J}_{k,\mathbf{j}}(r_{k,\mathbf{j}}) = \sum_{i=1}^n r_{k,\mathbf{j}}(\mathbf{x}_i)^2 - 2(-1)^k \partial^{k,\mathbf{j}} r_{k,\mathbf{j}}(\mathbf{x}_i). \quad (7)$$

To estimate $r_{k,\mathbf{j}}$, we employ a linear-in-parameter model:

$$r_{k,\mathbf{j}}(\mathbf{x}) = \sum_{i=1}^n \theta_{k,\mathbf{j}}^{(i)} \psi_{k,\mathbf{j}}^{(i)}(\mathbf{x}) = \boldsymbol{\theta}_{k,\mathbf{j}}^\top \boldsymbol{\psi}_{k,\mathbf{j}}(\mathbf{x}),$$

where we set $\psi_{k,\mathbf{j}}^{(i)}(\mathbf{x}) = \partial^{k,\mathbf{j}} K(\mathbf{x}, \mathbf{x}_i)$ and $K(\cdot, \cdot)$ is a smooth kernel function such as the Gaussian kernel. Substituting the model into (7) and adding the

ℓ_2 -regularizer yield the following quadratic objective function:

$$\begin{aligned} \tilde{J}_{k,j}(\boldsymbol{\theta}_{k,j}) \\ = \boldsymbol{\theta}_{k,j}^\top \mathbf{G}_{k,j} \boldsymbol{\theta}_{k,j} - 2(-1)^k \boldsymbol{\theta}_{k,j}^\top \mathbf{h}_{k,j} + \lambda_{k,j} \boldsymbol{\theta}_{k,j}^\top \boldsymbol{\theta}_{k,j}, \end{aligned} \quad (8)$$

where

$$\begin{aligned} \mathbf{G}_{k,j} &= \frac{1}{n} \sum_{i=1}^n \boldsymbol{\psi}_{k,j}(\mathbf{x}_i) \boldsymbol{\psi}_{k,j}(\mathbf{x}_i)^\top, \\ \mathbf{h}_{k,j} &= \frac{1}{n} \sum_{i=1}^n \partial^{k,j} \boldsymbol{\psi}_{k,j}(\mathbf{x}_i). \end{aligned}$$

The minimizer of (8) can be computed analytically as

$$\hat{\boldsymbol{\theta}}_{k,j} = \underset{\boldsymbol{\theta}_{k,j}}{\operatorname{argmin}} \tilde{J}_{k,j}(\boldsymbol{\theta}_{k,j}) = (-1)^k (\mathbf{G}_{k,j} + \lambda_{k,j} \mathbf{I})^{-1} \mathbf{h}_{k,j},$$

where \mathbf{I} denotes the identity matrix. Finally, a density-derivative-ratio estimator is given by

$$\hat{r}_{k,j}(\mathbf{x}) = \hat{\boldsymbol{\theta}}_{k,j}^\top \boldsymbol{\psi}_{k,j}(\mathbf{x}).$$

We call this method the *least-squares density-derivative ratios* (LSDDR). Note that when $k = 1$, LSDDR is reduced to an existing log-density derivative estimator [Cox, 1985, Sasaki et al., 2014]. Therefore, LSDDR can be regarded as its generalization to higher-order derivatives.

3.2 Theoretical Analysis

Next, we perform theoretical analysis of LSDDR. Let \mathcal{H} be a reproducing kernel Hilbert space (RKHS) associated with the kernel $K(\cdot, \cdot)$. In our analysis, we assume that the true density-derivative-ratio is contained in \mathcal{H} :

$$r_{k,j}^*(\mathbf{x}) := \frac{\partial^{k,j} p(\mathbf{x})}{p(\mathbf{x})} \in \mathcal{H}.$$

According to Zhou [2008], $\partial^{k,j} K(\mathbf{x}, \cdot)$ belongs to \mathcal{H} under regularity conditions. Thus, the linear-in-parameter model employed for $r_{k,j}$ also belongs to \mathcal{H} . In this analysis, we define a slightly modified version of the LSDDR estimator $\hat{r}_{k,j}$ as the optimal solution of

$$\begin{aligned} \min_{r \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^n \left\{ \frac{1}{2} r(\mathbf{x}_i)^2 - (-1)^k \partial^{k,j} r(\mathbf{x}_i) \right\} \\ \text{subject to } \|r\|_{\mathcal{H}}^2 \leq M_n^2, \end{aligned} \quad (9)$$

where $\|\cdot\|_{\mathcal{H}}$ denotes the norm in \mathcal{H} and M_n is a constant depending on the sample size n . Note that the optimal solution of (9) can be expressed by the kernelized LSDDR estimator with regularization parameter $\lambda_{k,j}$.

In the following theorem, we establish the convergence rate of LSDDR. The accuracy of the estimator is evaluated by the L_2 -norm under the distribution P defined as

$$\|h\|_P^2 = \int |h(\mathbf{x})|^2 dP(\mathbf{x}).$$

Theorem 1 (Convergence rate of LSDDR). *Let us assume that the kernel function $K(\mathbf{x}, \mathbf{x}')$ of \mathcal{H} is smooth and that there exists a constant $c_K > 0$ such that*

$$\begin{aligned} \partial^{j'} \partial^{j'} K(\mathbf{x}, \mathbf{x}) \\ := \frac{\partial^{j'} |}{\partial^{j'_1} x_1 \cdots \partial^{j'_D} x_D} \frac{\partial^{j'} |}{\partial^{j'_1} x'_1 \cdots \partial^{j'_D} x'_D} K(\mathbf{x}, \mathbf{x}') \Big|_{\mathbf{x}'=\mathbf{x}} \leq c_K \end{aligned} \quad (10)$$

holds for any $|j'| \leq k + \ell$, where $|j'|$ is the sum of the elements in $\mathbf{j}' = (j'_1, \dots, j'_D)$ and ℓ is a natural number greater than $D/2$. Suppose that M_n is of the poly-logarithmic order of n . Then,

$$\|\hat{r}_{k,j} - r_{k,j}^*\|_P^2 = O_P(c_n/n^{1/(2+D/\ell)}),$$

where $O_P(\cdot)$ denotes the probabilistic order and c_n is of the poly-logarithmic order of n .

Due to lack of space, we only provide a sketch of proof below. The full proof is given in the supplementary material.

Sketch of the proof of Theorem 1. Let $L(r)$ be the objective function of (9). Suppose that n is sufficiently large so that $\|r_{k,j}^*\|_{\mathcal{H}} \leq M_n$ holds. Then, the inequality $L(\hat{r}_{k,j}) \leq L(r_{k,j}^*)$ leads to

$$\begin{aligned} \frac{1}{2} \|\hat{r} - r_{k,j}^*\|_P^2 &\leq \frac{1}{2} \int \{(\hat{r})^2 - (r_{k,j}^*)^2\} d(P - P_n) \\ &\quad - (-1)^k \int \{\partial^{k,j} \hat{r} - \partial^{k,j} r_{k,j}^*\} d(P - P_n), \end{aligned}$$

where the equality, $\int (r_{k,j}^*)^2 dP = (-1)^k \int \partial^{k,j} r_{k,j}^* dP$, is used. Using Proposition 4 in Cucker and Smale [2002], we find that the convergence rate of $\|\hat{r} - r_{k,j}^*\|_P^2$ is closely related to the complexity of the function sets,

$$\begin{aligned} \mathcal{F} &= \{r^2 - (r_{k,j}^*)^2 \mid \|r\|_{\mathcal{H}} \leq M_n\}, \\ \mathcal{G} &= \{\partial^{k,j} r - \partial^{k,j} r_{k,j}^* \mid \|r\|_{\mathcal{H}} \leq M_n\}. \end{aligned}$$

Here, the complexity of the function set \mathcal{F}' is measured by the covering number $\mathcal{N}_\infty(\mathcal{F}', \varepsilon)$, which is defined as the minimal number $n \in \mathbb{N}$ such that there exist n disks with radius ε covering \mathcal{F}' , i.e.,

$$\begin{aligned} \mathcal{N}_\infty(\mathcal{F}', \varepsilon) &= \min\{n \in \mathbb{N} \mid \exists \mathcal{S} = \{g_1, \dots, g_n\} \subset \mathcal{F}', \\ &\quad \forall f \in \mathcal{F}', \exists g \in \mathcal{S}, \|f - g\|_\infty < \varepsilon\}. \end{aligned}$$

Under the assumption (10) with $|\mathbf{j}'| = 0$, the standard technique shown in Cucker and Smale [2002] can be used to obtain an upper bound of $\mathcal{N}_\infty(\mathcal{F}, \varepsilon)$. In order to obtain an upper bound of $\mathcal{N}_\infty(\mathcal{G}, \varepsilon)$, the formula on the derivative in the RKHS [Zhou, 2008], i.e.,

$$|\partial^{k'} \cdot \mathbf{j}' r(\mathbf{x})| \leq \|r\|_{\mathcal{H}} \sqrt{\partial \mathbf{j}' \partial' \mathbf{j}' K(\mathbf{x}, \mathbf{x})},$$

is used up to the order $k' \leq k + \ell$. Then, the uniform convergence rates of $\int (r^2 - (r_{k,j}^*)^2) d(P - P_n)$ and $\int (\partial^{k'} \cdot \mathbf{j}' r - \partial^{k'} \cdot \mathbf{j}' r_{k,j}^*) d(P - P_n)$ are respectively evaluated by these covering numbers. Eventually, we obtain the upper bound of $\|\hat{r} - r_{k,j}^*\|_P^2$. \square

3.3 Model Selection by Cross-Validation

In practice, the performance of LSDDR depends on the choice of hyper-parameters such as parameters in $\psi_{k,j}^{(i)}(\mathbf{x})$ and the regularization parameter $\lambda_{k,j}$. Such hyper-parameters can be chosen by cross-validation with respect to the squared-loss criterion in a straightforward way as follows:

1. Divide the samples $\mathcal{X} = \{\mathbf{x}_i\}_{i=1}^n$ into T disjoint subsets $\{\mathcal{X}_t\}_{t=1}^T$.
2. Obtain the estimator $\hat{r}_{k,j}^{(t)}(\mathbf{x})$ from $\mathcal{X} \setminus \mathcal{X}_t$, and then compute $\tilde{J}_{k,j}$ from the hold-out samples as

$$\begin{aligned} & \text{CV}(t) \\ &= \frac{1}{|\mathcal{X}_t|} \sum_{\mathbf{x} \in \mathcal{X}_t} \left[\left\{ \hat{r}_{k,j}^{(t)}(\mathbf{x}) \right\}^2 - 2(-1)^k \partial^{k,j} \hat{r}_{k,j}^{(t)}(\mathbf{x}) \right], \end{aligned}$$

where $|\mathcal{X}_t|$ denotes the number of elements in \mathcal{X}_t .

3. Choose the model that minimizes $\text{CV} = \frac{1}{T} \sum_{t=1}^T \text{CV}(t)$.

4 Application in Density Ridge Estimation

In this section, based on LSDDR proposed in Section 3, we develop a novel density ridge estimator called the *least-squares density ridge finder* (LSDRF). For LSDDR, we employ the Gaussian kernel with bandwidth $\sigma_{k,j}$.¹

The algorithm of LSDRF essentially follows the same line as SCMS (see Figure 1), i.e., projected gradient ascent is performed:

$$\mathbf{x} \leftarrow \mathbf{x} + \hat{\mathbf{V}}_{\text{LS}} \hat{\mathbf{V}}_{\text{LS}}^\top \hat{\mathbf{m}}_{\text{LS}}(\mathbf{x}).$$

¹If the sample size n is large, we may only use a subset of data samples as Gaussian centers in LSDDR.

However, compared with the update rule (3) used in SCMS, $\hat{\mathbf{V}}_{\text{KDE}}$ and $\hat{\mathbf{m}}_{\text{MS}}(\mathbf{x})$ are replaced with $\hat{\mathbf{V}}_{\text{LS}}$ and $\hat{\mathbf{m}}_{\text{LS}}(\mathbf{x})$, respectively.²

$\hat{\mathbf{V}}_{\text{LS}}$ is obtained by applying PCA to an estimate of the inverse local-covariance matrix obtained based on LSDDR, not KDE:

$$\hat{\Sigma}_{\text{LS}}^{-1}(\mathbf{x}) = -\hat{\mathbf{H}}_{\text{LS}}(\mathbf{x}) + \hat{\mathbf{g}}_{\text{LS}}(\mathbf{x}) \hat{\mathbf{g}}_{\text{LS}}(\mathbf{x})^\top, \quad (11)$$

where the elements in $\hat{\mathbf{g}}_{\text{LS}}(\mathbf{x})$ and $\hat{\mathbf{H}}_{\text{LS}}(\mathbf{x})$ are the LSDDR solutions $\hat{r}_{1,j}(\mathbf{x})$ and $\hat{r}_{2,j}(\mathbf{x})$, respectively.

$\hat{\mathbf{m}}_{\text{LS}}(\mathbf{x})$ is given by

$$\begin{aligned} [\hat{\mathbf{m}}_{\text{LS}}(\mathbf{x})]_\ell &= \frac{\sum_{i=1}^n \hat{\theta}_{1,\ell}^{(i)}[\mathbf{x}_i]_\ell K(\mathbf{x}, \mathbf{x}_i)}{\sum_{i=1}^n \hat{\theta}_{1,\ell}^{(i)} K(\mathbf{x}, \mathbf{x}_i)} - [\mathbf{x}]_\ell \\ &= \frac{1}{\sum_{i=1}^n \hat{\theta}_{1,\ell}^{(i)} K(\mathbf{x}, \mathbf{x}_i)} [\hat{\mathbf{g}}_{\text{LS}}(\mathbf{x})]_\ell, \end{aligned} \quad (12)$$

where $[\mathbf{x}]_\ell$ denotes the ℓ -th element in \mathbf{x} . This vector comes from another mode-seeking clustering method called *least-squares log-density gradients* (LSLDG) clustering [Sasaki et al., 2014], which was experimentally shown to work much better than MS especially for higher-dimensional data.

5 Experiments

In this section, we experimentally demonstrate the usefulness of LSDRF.

5.1 Illustration on Simulated Data

First, we investigate the performance of LSDRF and compare it with SCMS on a variety of simulated datasets.³ The i -th observation of data was generated according to $x_i^{(j)} = f^{(j)}(t_i) + n_i^{(j)}$, where t_i was taken from some range at regular intervals, $f^{(j)}(\cdot)$ denotes some fixed function, and $n_i^{(j)}$ was the Gaussian noise with mean 0 and standard deviation γ . The bandwidth h of the Gaussian kernel in SCMS was determined by least-squares cross-validation using ten candidates from $10^{-1.5}$ to 10^0 at regular intervals in logarithmic scale. For LSDDR, model selection was performed by five-fold cross-validation using ten candidates from 10^{-1} (or 10^{-4}) to $10^{0.5}$ (or 10^0) for $\sigma_{k,j}$

²To avoid numerical instability, we stop updating a data point \mathbf{x} if $\sum_{i=1}^n \hat{\theta}_{1,\ell}^{(i)} K(\mathbf{x}, \mathbf{x}_i)$ is less than a very small positive constant.

³Most of the datasets are generated using a MATLAB package made by Jakob Verbeek, which is available at http://lear.inrialpes.fr/people/verbeek/code/kseg_soft.tar.gz.

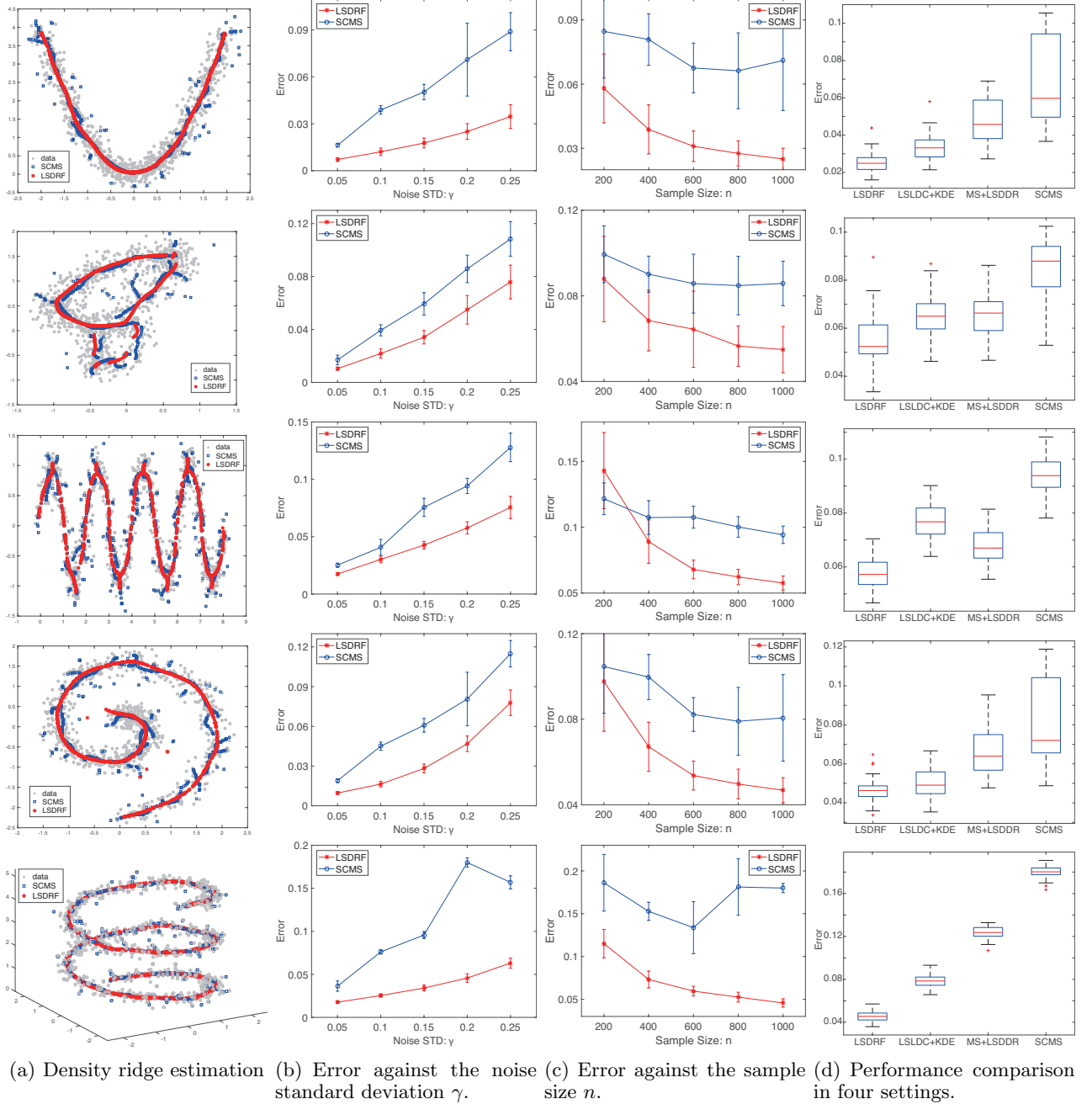


Figure 2: Density ridge estimation on simulated data. In (b) and (c), each point and error bar denote the average and standard deviation of estimation errors over 50 runs, respectively. We set $(n, \gamma) = (1000, 0.15)$ for (a), $n = 1000$ for (b) and (d), and $\gamma = 0.2$ for (c).

(or $\lambda_{k,j}$) at regular intervals in logarithmic scale. The estimation error was measured by

$$\text{Error} = \frac{1}{n} \sum_{i=1}^n \min_l \|\hat{\mathbf{y}}_i - \mathbf{f}(t_l)\|_2,$$

where $\hat{\mathbf{y}}_i$ denotes an estimate of the density ridge obtained from \mathbf{x}_i and $\mathbf{f}(\cdot) = (f^{(1)}(\cdot), f^{(2)}(\cdot), \dots, f^{(D)}(\cdot))^\top$.

Two density ridge estimates from LSDRF and SCMS are visualized in Figures 2(a). LSDRF provides more accurate and smooth ridge estimates than SCMS on all datasets. Figures 2(b) show the noise tolerance property of both methods. As the noise standard deviation γ increases, the performance of both LSDRF and SCMS gets worse, but LSDRF still works better than SCMS in all five cases. Figures 2(c) further indi-

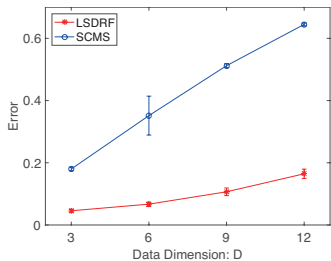


Figure 3: Performance comparison for higher-dimensional data. We set $(n, \gamma) = (1000, 0.2)$.

icates that the estimation error of LSDRF more quickly decreases with respect to the sample size n .

Next, we investigate how LSDRF and SCMS perform in higher-dimensional cases. For this experiment, to create higher-dimensional data, Gaussian variables with mean 0 and standard deviation $\gamma = 0.2$ were appended to the three-dimensional spiral data in Figure 2. The result shown in Figure 3 indicates that LSDRF is more useful than SCMS especially for higher-dimensional data.

LSDRF makes two major modifications from SCMS: (a) $\widehat{\mathbf{m}}_{\text{LS}}$ from a previously proposed mode-seeking clustering method [Sasaki et al., 2014] and (b) $\widehat{\Sigma}_{\text{LS}}$ from LSDDR. To understand how these modifications improved the performance in ridge estimation, we performed experiments with the following additional methods:

- LSDRF: $\widehat{\mathbf{m}}_{\text{LS}}$ and $\widehat{\Sigma}_{\text{LS}}$ are used.
- (LS+KDE): $\widehat{\mathbf{m}}_{\text{LS}}$ and $\widehat{\Sigma}_{\text{KDE}}$ are used.
- (MS+LSDDR): $\widehat{\mathbf{m}}_{\text{MS}}$ and $\widehat{\Sigma}_{\text{LS}}$ are used.
- SCMS: $\widehat{\mathbf{m}}_{\text{MS}}$ and $\widehat{\Sigma}_{\text{KDE}}$ are used.

Figures 2(d) show that (LS+KDE) and (MS+LSDDR) improve SCMS, while LSDRF performs better than (LS+KDE) and (MS+LSDDR). Therefore, both $\widehat{\Sigma}_{\text{LS}}$ and $\widehat{\mathbf{m}}_{\text{LS}}$ surely contribute to improving the performance of SCMS. However, the amount of performance improvement obtained from $\widehat{\mathbf{m}}_{\text{LS}}$ or from $\widehat{\Sigma}_{\text{LS}}$ seems dependent on datasets.

5.2 Density Ridge Visualization on Real-World Datasets

Finally, we apply LSDRF to real-world datasets. As in Pulkkinen [2015], we employed the following two datasets:

- *New Madrid earthquake* dataset: This seismological dataset was downloaded from the Cen-

ter for Earthquake Research and Information.⁴ The dataset contains positional information for earthquakes around the New Madrid seismic zone from 1974 to 2016, providing 11,131 samples. The three regions in Figures 4(a,b,c) were extracted according to (a) $(-90.2, -89.25)$, (b) $(-92.5, -92.15)$ and (c) $(-85.5, -83.5)$ degrees for the latitude range. For the longitude range, (a) $(36, 36.8)$, (b) $(35.2, 35.4)$ and (c) $(34.5, 36.5)$ degrees were selected. The total number of the original data samples and reduced data samples in each region was (a) $(n, n') = (5902, 500)$, (b) $(n, n') = (1548, 300)$ and (c) $(n, n') = (594, 200)$.

- *Shapley galaxy* dataset: This dataset was downloaded from the Center for Astrostatistics at Pennsylvania State University.⁵ The dataset contains information about the three-dimensional sky angles and recession velocity of 4,215 galaxies. As done in Pulkkinen [2015], we transformed the data samples into the three-dimensional Cartesian coordinates based on the fact that the recession velocity is proportional to the radial distance [Drinkwater et al., 2004]. The three regions in Figures 4(a,b,c) were extracted according to a velocity range: (a) $(6000, 20000)$ km/s, (b) $(1500, 6000)$ km/s and (c) $(6000, 10500)$ km/s, respectively. The total number of the original data samples and reduced data samples in each region was (a) $(n, n') = (2849, 500)$, (b) $(n, n') = (595, 200)$ and (c) $(n, n') = (351, 150)$.

In each dataset, we focused on three regions containing prominent features, and standardized data samples in each region by subtracting the mean value and dividing by standard deviation in a dimension-wise manner. For performance comparison, we computed the log-likelihood of density ridge estimates, which is defined by $\mathcal{L} = \frac{1}{n'} \sum_{l=1}^{n'} \log \widehat{p}_{\text{KDE}}(\widehat{\mathbf{y}}_l)$: The kernel centers in \widehat{p}_{KDE} were set at data samples $\{\mathbf{x}_i\}_{i=1}^n$ in each region, while density ridges estimates $\{\widehat{\mathbf{y}}_l\}_{l=1}^{n'}$ were obtained from $n' (< n)$ data samples randomly chosen from $\{\mathbf{x}_i\}_{i=1}^n$. If \mathcal{L} is larger, the performance can be interpreted to be better because ridges are defined on relatively high density areas. Unlike the last experiment, we used the following adaptive-bandwidth Gaussian kernel in LSDDR: The bandwidth parameter of each Gaussian kernel was set at the Euclidean distance to the m -nearest sample from the Gaussian center where m was cross-validated. The regularization parameter λ_k , which is supposed to be common to all j in this section, was also determined by cross-validation as in

⁴<http://www.memphis.edu/cei/seismic/>

⁵http://astrostatistics.psu.edu/datasets/Shapley_galaxy.html

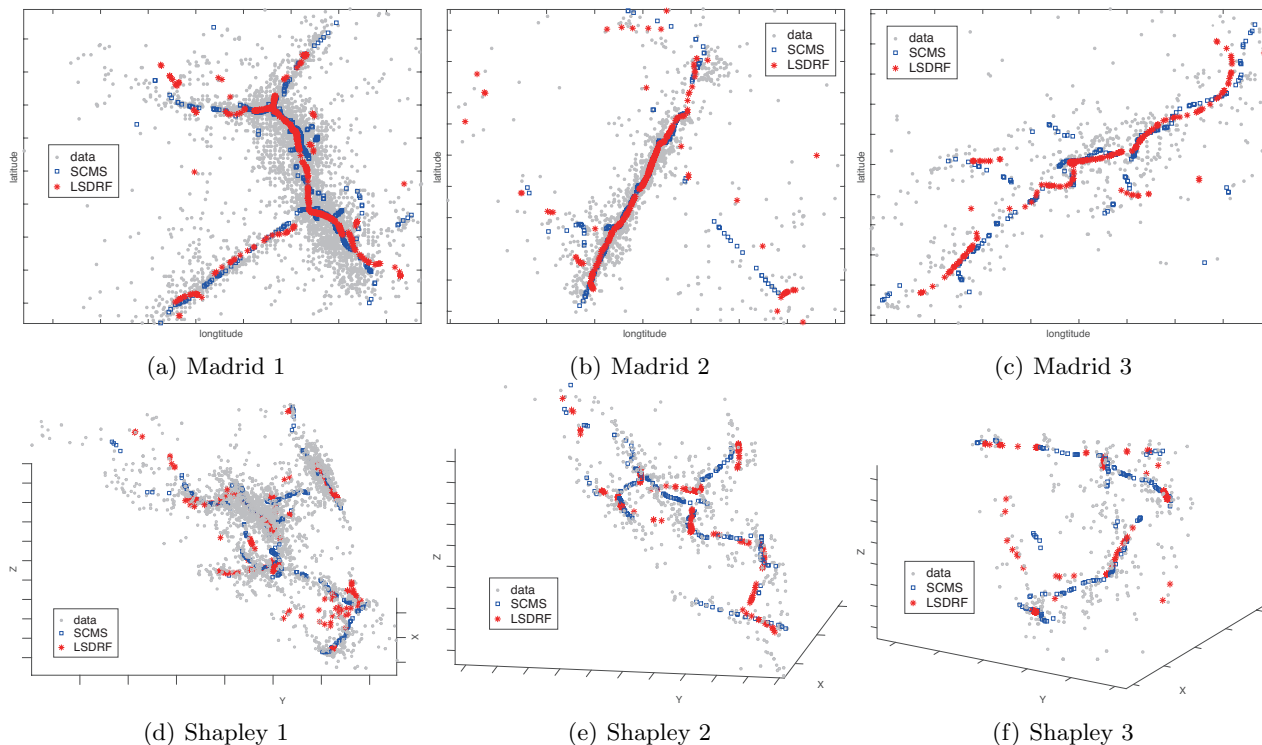


Figure 4: Density ridge estimation to the (a,b,c) New Madrid earthquake and (d,e,f) Shapley galaxy datasets. Three regions were extracted from each dataset according to (a,b,c) a range of latitude and longitude, and (d,e,f) a range of recession velocity.

Section 3.3. The adaptive-bandwidth Gaussian kernel was also used for SCMS, and a similar cross-validation procedure was performed.

Ridges estimated by LSDRF are often smooth and seem to well-match the ridges of the underlying data (Figure 4). Table 1 quantitatively substantiates that LSDRF overall performs better than SCMS.

6 Conclusion

In this paper, we proposed a new algorithm for density ridge estimation. Our main contribution was the *least-squares density-derivative-ratios* (LSDDR) estimator, which avoids density estimation and division by an estimated density. We theoretically established a non-parametric convergence bound for LSDDR and experimentally demonstrated the superior performance of the density ridge estimator constructed based on LSDDR.

Acknowledgements

HS was supported by KAKENHI 15H06103, TK was supported by KAKENHI 16K00044 and MS acknowledges the JST CREST program

Table 1: The average and standard deviation of the log-likelihood of density ridges over 50 runs. A larger value means a better result. Numbers in the parentheses are standard deviations. The best and comparable methods judged by the unpaired t-test at the significance level 5% are described in boldface.

New Madrid earthquake		
	LSDRF	SCMS
Madrid 1	-0.611(0.109)	-0.632(0.056)
Madrid 2	-0.004(0.108)	-0.051(0.086)
Madrid 3	-1.097(0.146)	-1.238(0.086)
Shapley galaxy		
	LSDRF	SCMS
Shapley 1	0.125(0.101)	0.038(0.082)
Shapley 2	-1.285(0.113)	-1.216(0.089)
Shapley 3	-1.252(0.197)	-1.550(0.086)

References

Y.-C. Chen, C. R. Genovese, S. Ho, and L. Wasserman. Optimal ridge detection using coverage risk. In *Advances in Neural Information Processing Systems 28*, pages 316–324. 2015a.

- Y.-C. Chen, C. R. Genovese, and L. Wasserman. Asymptotic theory for density ridges. *The Annals of Statistics*, 43(5):1896–1928, 2015b.
- Y.-C. Chen, S. Ho, P. Freeman, C. Genovese, and L. Wasserman. Cosmic web reconstruction through density ridges: Method and algorithm. *Monthly Notices of the Royal Astronomical Society*, 454(1): 1140–1156, 2016.
- Y. Cheng. Mean shift, mode seeking, and clustering. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 17(8):790–799, 1995.
- D. Comaniciu and P. Meer. Mean shift: A robust approach toward feature space analysis. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 24(5):603–619, 2002.
- D. D. Cox. A penalty method for nonparametric estimation of the logarithmic derivative of a density function. *Annals of the Institute of Statistical Mathematics*, 37(1):271–288, 1985.
- F. Cucker and S. Smale. On the mathematical foundations of learning. *Bulletin of the American Mathematical Society*, 39:1–49, 2002.
- M. J. Drinkwater, Q. A. Parker, D. Proust, E. Slezak, and H. Quintana. The large scale distribution of galaxies in the shapley supercluster. *Publications of the Astronomical Society of Australia*, 21(1):89–96, 2004.
- D. Eberly. *Ridges in Image and Data Analysis*. Springer, 1996.
- J. Einbeck and J. Dwyer. Using principal curves to analyse traffic patterns on freeways. *Transportmetrica*, 7(3):229–246, 2011.
- K. Fukunaga and L. Hostetler. The estimation of the gradient of a density function, with applications in pattern recognition. *IEEE Transactions on Information Theory*, 21(1):32–40, 1975.
- C. R. Genovese, M. Perone-Pacifico, I. Verdinelli, and L. Wasserman. Nonparametric ridge estimation. *The Annals of Statistics*, 42(4):1511–1545, 2014.
- Y. A. Ghassabeh, T. Linder, and G. Takahara. On some convergence properties of the subspace constrained mean shift. *Pattern Recognition*, 46(11): 3140–3147, 2013.
- T. Hastie. *Principal Curves and Surfaces*. PhD thesis, Stanford University, 1984.
- T. Hastie and W. Stuetzle. Principal curves. *Journal of the American Statistical Association*, 84(406):502–516, 1989.
- B. Kégl and A. Krzyżak. Piecewise linear skeletonization using principal curves. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 24(1): 59–74, 2002.
- U. Ozertem and D. Erdogmus. Locally defined principal curves and surfaces. *Journal of Machine Learning Research*, 12:1249–1286, 2011.
- S. Pulkkinen. Ridge-based method for finding curvilinear structures from noisy data. *Computational Statistics & Data Analysis*, 82:89–109, 2015.
- H. Sasaki, A. Hyvärinen, and M. Sugiyama. Clustering via mode seeking by direct estimation of the gradient of a log-density. In *Machine Learning and Knowledge Discovery in Databases Part III- European Conference, ECML/PKDD 2014*, volume 8726, pages 19–34, 2014.
- H. Sasaki, Y. K. Noh, and M. Sugiyama. Direct density-derivative estimation and its application in KL-divergence approximation. In *Proceedings of the 18th International Conference on Artificial Intelligence and Statistics (AISTATS)*, pages 809–818, 2015.
- J. Su, A. Srivastava, and F. Huffer. Detection, classification and estimation of individual shapes in 2D and 3D point clouds. *Computational Statistics & Data Analysis*, 58:227–241, 2013.
- S. You, E. Bas, D. Erdogmus, and J. Kalpathy-Cramer. Principal curved based retinal vessel segmentation towards diagnosis of retinal diseases. In *IEEE International Conference on Healthcare Informatics, Imaging and Systems Biology (HISB)*, pages 331–337. IEEE, 2011.
- D. Zhou. Derivative reproducing properties for kernel methods in learning theory. *Journal of Computational and Applied Mathematics*, 220(1-2):456–463, 2008.