

# *Sparse density estimator with tunable kernels*

Article

Accepted Version

Creative Commons: Attribution-Noncommercial-No Derivative Works 4.0

Hong, X., Chen, S. and Becerra, V. (2016) Sparse density estimator with tunable kernels. *Neurocomputing*, 173 (3). pp. 1976-1982. ISSN 0925-2312 doi: <https://doi.org/10.1016/j.neucom.2015.08.021> Available at <http://centaur.reading.ac.uk/65632/>

It is advisable to refer to the publisher's version if you intend to cite from the work. See [Guidance on citing](#).

To link to this article DOI: <http://dx.doi.org/10.1016/j.neucom.2015.08.021>

Publisher: Elsevier

All outputs in CentAUR are protected by Intellectual Property Rights law, including copyright law. Copyright and IPR is retained by the creators or other copyright holders. Terms and conditions for use of this material are defined in the [End User Agreement](#).

[www.reading.ac.uk/centaur](http://www.reading.ac.uk/centaur)

## **CentAUR**

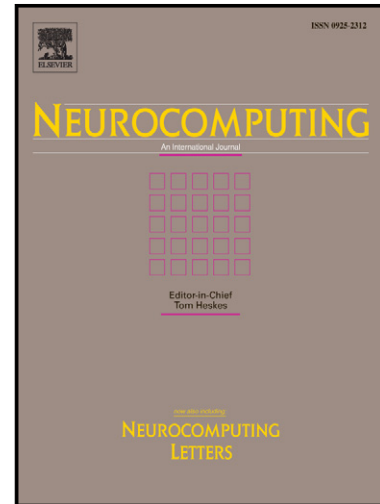
Central Archive at the University of Reading

Reading's research outputs online

# Author's Accepted Manuscript

Sparse Density Estimator with Tunable Kernels

Xia Hong, Sheng Chen, Victor M. Becerra



[www.elsevier.com/locate/neucom](http://www.elsevier.com/locate/neucom)

PII: S0925-2312(15)01173-X  
DOI: <http://dx.doi.org/10.1016/j.neucom.2015.08.021>  
Reference: NEUCOM15949

To appear in: *Neurocomputing*

Received date: 30 March 2015  
Revised date: 29 June 2015  
Accepted date: 3 August 2015

Cite this article as: Xia Hong, Sheng Chen, Victor M. Becerra, Sparse Density Estimator with Tunable Kernels, *Neurocomputing*, <http://dx.doi.org/10.1016/j.neucom.2015.08.021>

This is a PDF file of an unedited manuscript that has been accepted for publication. As a service to our customers we are providing this early version of the manuscript. The manuscript will undergo copyediting, typesetting, and review of the resulting galley proof before it is published in its final citable form. Please note that during the production process errors may be discovered which could affect the content, and all legal disclaimers that apply to the journal pertain.

# Sparse Density Estimator with Tunable Kernels

Xia Hong, Sheng Chen, and Victor M. Becerra

## Abstract

A new sparse kernel density estimator with tunable kernels is introduced within a forward constrained regression framework whereby the nonnegative and summing-to-unity constraints of the mixing weights can easily be satisfied. Based on the minimum integrated square error criterion, a recursive algorithm is developed to select significant kernels one at time, and the kernel width of the selected kernel is then tuned using the gradient descent algorithm. Numerical examples are employed to demonstrate that the proposed approach is effective in constructing very sparse kernel density estimators with competitive accuracy to existing kernel density estimators.

## Index Terms

Probability density function, kernel density estimator, sparse modeling, minimum integrated square error

## I. INTRODUCTION

The probability density function (PDF) estimation, e.g., the Parzen window (PW) and finite mixture model, is of fundamental importance to many data analysis and pattern recognition applications [1]–[8]. There is a considerable interest into research on sparse PDF estimation which can be summarized into two categories. The first category is based on constrained optimization. For example, the support vector machine (SVM) density estimation was researched [9], [10], in which the density estimation problem is formulated as a supervised learning mode whilst the mean absolute deviation between the empirical cumulative distribution function (CDF) calculated from the training data and the CDF based on the PDF estimator also calculated from the training data are minimized. This yields the sparsity inducing property, i.e., at the optimality,

X. Hong and V.M. Becerra are with School of Systems Engineering, University of Reading, Reading, RG6 6AY, UK (E-mails: x.hong@reading.ac.uk, v.m.becerra@reading.ac.uk).

S. Chen is with Electronics and Computer Science, University of Southampton, Southampton SO17 1BJ, UK (E-mail: sqc@ecs.soton.ac.uk), and also with King Abdulaziz University, Jeddah 21589, Saudi Arabia.

many kernels' weights are driven to zero. The desirable property of sparsity inducing also happens in the interesting approach of reduced set density estimator (RSDE) [11], which is based on the minimization of the integrated square error (ISE) between the estimator and the true density evaluated on the training data [2], [11], [12], and two efficient optimization algorithms were introduced. Alternatively, by exploiting the first and second order Riemannian geometry of the multinomial manifold, the Riemannian trust-region algorithm [13] was recently applied to find the set of sparse mixing coefficients based on the minimum ISE (MISE), referred to as the RTR-MISE algorithm [14].

The second category of sparse kernel density estimators construct the PDF estimator in a forward regression manner. A regression-based PDF estimation method was introduced [15], in which the empirical CDF is constructed and used as the desired response. In order to automatically determine the model structure with the improved model generalization, the regression-based idea of [15] and the approach of [16] were extended to yield an orthogonal forward regression (OFR) based sparse density estimation algorithm [17] which is capable of automatically constructing a very sparse kernel density estimate, with comparable performance to that of the PW estimate. A simpler and viable alternative approach was proposed to use kernels directly as regressors by adopting the PW estimate as the target response [18]. A sparse kernel density estimator [19] was introduced based on the MISE and the forward constrained regression (FCR) [20] to select significant kernels one at time, which has very low computational cost and is referred to as the FCR-MISE algorithm.

With the exception of [4], in all the above-mentioned sparse kernel estimators, including those based on the MISE approach [11], [14], [19], the PDF kernels involve a single and *fixed* kernel bandwidth parameter that needs to be empirically predetermined. By contrast, this paper introduces a new sparse kernel density estimator with *tunable* kernels also based on the MISE. Specifically, a new recursive algorithm is developed to select significant kernels one at time, followed by tuning the kernel width of the selected kernel using the gradient descent algorithm. This means that there is no need to determine the bandwidth parameters empirically outside the algorithm loop. Numerical examples are employed to demonstrate that the proposed approach can construct very sparse kernel density estimates with competitive accuracy, compared to the existing kernel density estimators.

## II. FORWARD CONSTRUCTION OF TUNABLE SPARSE KERNEL DENSITY ESTIMATOR

Given the finite data set  $D_N = \{\mathbf{x}_j\}_{j=1}^N$  consisting of  $N$  data samples, where the data  $\mathbf{x}_j \in \mathbb{R}^m$  follows an unknown PDF  $p(\mathbf{x})$ , the problem under study is to find a sparse approximation of  $p(\mathbf{x})$  by forward construction based on the subset  $D_M = \{\mathbf{x}'_1, \mathbf{x}'_2, \dots, \mathbf{x}'_M\}$  of  $M$  data samples selected from  $D_N$ . For example, if  $\mathbf{x}_6$  from  $D_N$  is selected to form the first kernel, it is denoted as  $\mathbf{x}'_1$  in  $D_M$ . A general kernel based density estimate of  $p(\mathbf{x})$  is given by

$$\hat{p}^{(M)}(\mathbf{x}; \boldsymbol{\beta}_M, \boldsymbol{\sigma}_M) = \sum_{i=1}^M \beta_i K_{\sigma_i}(\mathbf{x}, \mathbf{x}'_i) \quad (1)$$

subject to

$$\beta_i \geq 0, \quad \text{and} \quad \boldsymbol{\beta}_M^T \mathbf{1}_M = 1, \quad (2)$$

where  $K_{\sigma_i}(\mathbf{x}, \mathbf{x}'_i)$  is the Gaussian kernel with the kernel center vector  $\mathbf{x}'_i$  and an *adjustable* kernel width  $\sigma_i$  given by

$$K_{\sigma_i}(\mathbf{x}, \mathbf{x}'_i) = \frac{1}{(2\pi\sigma_i^2)^{m/2}} \exp\left(-\frac{\|\mathbf{x} - \mathbf{x}'_i\|^2}{2\sigma_i^2}\right), \quad (3)$$

and  $\beta_i$  is the  $i$ th kernel weight, while  $\boldsymbol{\sigma}_M = [\sigma_1 \ \sigma_2 \ \dots \ \sigma_M]^T$ ,  $\boldsymbol{\beta}_M = [\beta_1 \ \beta_2 \ \dots \ \beta_M]^T$ , and  $\mathbf{1}_M$  is the  $M$ -dimensional vector whose elements are all equal to one.

We form the kernel density estimator (1) from the subset  $D_M$  in a forward construction manner. Specifically given the initial condition  $\sigma_i = \sigma_0, \forall i$ , and starting from an empty model set, our proposed sparse kernel density estimation algorithm selects the kernel functions  $K_{\sigma_0}(\mathbf{x}, \mathbf{x}'_i)$  into the model set one at a time from  $D_N$ . At each forward step, the associated kernel width  $\sigma_i$  is then optimized to obtain  $K_{\sigma_i}(\mathbf{x}, \mathbf{x}'_i)$ .

Let the superscript  $(l)$  denote the  $l$ th forward selection step. At the  $l$ th forward selection step, further denote the intermediate kernel density estimator  $\hat{p}^{(l)}(\mathbf{x}; \boldsymbol{\beta}_l^{(l)}, \boldsymbol{\sigma}_l)$  as  $\hat{y}^{(l)}(\mathbf{x})$ , where  $\boldsymbol{\sigma}_l^{(l)} = [\sigma_1 \ \sigma_2 \ \dots \ \sigma_l]^T$  and  $\boldsymbol{\beta}_l^{(l)} = [\beta_1^{(l)} \ \beta_2^{(l)} \ \dots \ \beta_l^{(l)}]^T$ , with  $\beta_i^{(l)}, 1 \leq i \leq l$ , as the kernels weights at the  $l$ th forward selection step, i.e.,

$$\hat{y}^{(l)}(\mathbf{x}) = \sum_{i=1}^l \beta_i^{(l)} K_{\sigma_i}(\mathbf{x}, \mathbf{x}'_i). \quad (4)$$

The proposed algorithm integrates the FCR procedure [20] described below:

(i) At the first step, the PDF estimator is simply

$$\hat{y}^{(1)}(\mathbf{x}) = K_{\sigma_1}(\mathbf{x}, \mathbf{x}'_1), \quad (5)$$

where  $K_{\sigma_1}(\mathbf{x}, \mathbf{x}'_1)$  is obtained by adjusting the kernel width from  $\sigma_0$  to  $\sigma_1$  based on the selected kernel center  $\mathbf{x}'_1$ . Clearly  $\beta_1^{(1)} = 1$ .

- (ii) At the  $l$ th step, where  $l \geq 2$ , the PDF estimator is constructed by adding the  $l$ th kernel  $K_{\sigma_l}(\mathbf{x}, \mathbf{x}'_l)$  to  $\hat{y}^{(l-1)}(\mathbf{x})$  according to

$$\hat{y}^{(l)}(\mathbf{x}) = \lambda_l \hat{y}^{(l-1)}(\mathbf{x}) + (1 - \lambda_l) K_{\sigma_l}(\mathbf{x}, \mathbf{x}'_l), \quad (6)$$

where  $K_{\sigma_l}(\mathbf{x}, \mathbf{x}'_l)$  is obtained by adjusting the kernel width from  $\sigma_0$  to  $\sigma_l$  based on the selected kernel center  $\mathbf{x}'_l$ , while  $0 \leq \lambda_l \leq 1, \forall l$ , and  $\lambda_1 = 0$ .

It can be straightforwardly verified that the model constructed using the FCR procedure satisfies the convex constraint conditions of (2), namely,  $\beta_i^{(l)} \geq 0, 1 \leq i \leq l$ , and  $\sum_{i=1}^l \beta_i^{(l)} = 1, \forall l \geq 1$ , see [20]. Moreover, given  $\lambda_l$  and  $\beta_{l-1}^{(l-1)}, \beta_l^{(l)}$  can be recursively computed via

$$\beta_l^{(l)} = \begin{bmatrix} \lambda_l \beta_{l-1}^{(l-1)} \\ 1 - \lambda_l \end{bmatrix}, \quad (7)$$

where  $l > 1$  and  $\beta_1^{(1)} = \beta_1^{(1)} = 1$ .

It can be seen that the key issues at each forward selection step  $l$  are 1) how to initially select the kernel center vector  $\mathbf{x}'_l$  with the kernel width  $\sigma_l = \sigma_0$ , followed by adjusting the kernel width  $\sigma_l$  for the selected kernel; and 2) how to compute  $\lambda_l$  and hence the kernel weight  $\beta_l^{(l)}$ .

### III. JOINT KERNEL SELECTION AND KERNEL WIDTH OPTIMIZATION BASED ON THE MISE

We now introduce our new algorithm integrating the kernel term selection, the kernel width optimization and the kernel weight calculation based on MISE [2], [11], [12] and the FCR framework described in the previous section. In particular, we detail the joint kernel selection, the tunable kernel width optimization and kernel weight estimation at the  $l$ th forward selection stage. Specifically, based on the ISE criterion, we formulate initially the kernel weight estimation problem for a given kernel per forward selection step, and then the kernel width optimization using the gradient descent algorithm for the selected kernel. Joint kernel selection together with the kernel width/weights optimization are finally presented.

#### A. Kernel weight estimation

At the  $l$ th forward selection stage,  $K_{\sigma_i}(\mathbf{x}, \mathbf{x}'_i)$  are given for  $1 \leq i \leq l-1$ , and we consider the problem of determining  $\lambda_l$  and  $\sigma_l$  for a fixed  $\mathbf{x}'_l$  based on the ISE which is the global accuracy

measure for density estimate and is given by [11]

$$\begin{aligned}
\text{ISE}(\boldsymbol{\beta}_l^{(l)}, \boldsymbol{\sigma}_l) &= \int \left( p(\mathbf{x}) - \sum_{i=1}^l \beta_i^{(l)} K_{\sigma_i}(\mathbf{x}, \mathbf{x}'_i) \right)^2 d\mathbf{x} \\
&= \int p^2(\mathbf{x}) d\mathbf{x} + \int \left( \sum_{i=1}^l \beta_i^{(l)} K_{\sigma_i}(\mathbf{x}, \mathbf{x}'_i) \right)^2 d\mathbf{x} \\
&\quad - 2E \left[ \sum_{i=1}^l \beta_i^{(l)} K_{\sigma_i}(\mathbf{x}, \mathbf{x}'_i) \right] \\
&= \int p^2(\mathbf{x}) d\mathbf{x} + \sum_{i=1}^l \sum_{j=1}^l \beta_i^{(l)} \beta_j^{(l)} \int K_{\sigma_i}(\mathbf{x}, \mathbf{x}'_i) K_{\sigma_j}(\mathbf{x}, \mathbf{x}'_j) d\mathbf{x} \\
&\quad - 2 \sum_{i=1}^l \beta_i^{(l)} E \left[ K_{\sigma_i}(\mathbf{x}, \mathbf{x}'_i) \right] \\
&= \int p^2(\mathbf{x}) d\mathbf{x} + Q^{(l)}(\lambda_l, \sigma_l), \tag{8}
\end{aligned}$$

in which  $E[\bullet]$  denotes the expectation with respect to the true density  $p(\mathbf{x})$ . Since the unknown term  $\int p^2(\mathbf{x}) d\mathbf{x}$  is independent of  $\beta_l^{(l)}$ , it can be dropped from the objective function. We write the argument directly as  $\{\lambda_l, \sigma_l\}$  for the last term  $Q^{(l)}(\lambda_l, \sigma_l)$ , which becomes our objective function. We point out that since our algorithm is based on the FCR framework, only these two parameters need to be estimated at the  $l$ th selection stage. Note that  $\beta_l^{(l)}$  depends on  $\lambda_l$  and  $\beta_{l-1}^{(l-1)}$ , i.e., the sequence  $\{\lambda_1, \lambda_2, \dots, \lambda_{l-1}\}$ , that have already been obtained from the previous forward selection steps (see (7)). Similarly  $\{\sigma_1, \sigma_2, \dots, \sigma_{l-1}\}$  are also obtained from the previous forward selection steps.

Using the following unbiased estimator of  $E[K_{\sigma_i}(\mathbf{x}, \mathbf{x}'_i)]$

$$E[K_{\sigma_i}(\mathbf{x}, \mathbf{x}'_i)] \approx \frac{1}{N} \sum_{k=1}^N K_{\sigma_i}(\mathbf{x}_k, \mathbf{x}'_i) \tag{9}$$

as well as noting the result of  $\int K_{\sigma_i}(\mathbf{x}, \mathbf{x}'_i) K_{\sigma_j}(\mathbf{x}, \mathbf{x}'_j) d\mathbf{x}$  given in Appendix yield

$$\begin{aligned}
Q^{(l)}(\lambda_l, \sigma_l) &\triangleq \sum_{i=1}^l \sum_{j=1}^l \beta_i^{(l)} \beta_j^{(l)} K_{\sigma_{i,j}}(\mathbf{x}'_i, \mathbf{x}'_j) \\
&\quad - \frac{2}{N} \sum_{i=1}^l \beta_i^{(l)} \sum_{k=1}^N K_{\sigma_i}(\mathbf{x}_k, \mathbf{x}'_i), \tag{10}
\end{aligned}$$



where  $\sigma_{i,j} = \sqrt{\sigma_i^2 + \sigma_j^2}$ . Using matrix expression, we easily obtain the recursive form of  $Q^{(l)}(\lambda_l, \sigma_l)$  which is given by

$$Q^{(l)}(\lambda_l, \sigma_l) = \mu^{(l)} - 2\nu^{(l)} \quad (11)$$

where

$$\begin{cases} \mu^{(l)} = (\boldsymbol{\beta}_l^{(l)})^T \mathbf{C}_l^{(l)} \boldsymbol{\beta}_l^{(l)}, \\ \nu^{(l)} = (\boldsymbol{\beta}_l^{(l)})^T \mathbf{p}_l^{(l)}, \end{cases} \quad (12)$$

in which  $\mathbf{p}_l^{(l)}$  and  $\mathbf{C}_l^{(l)}$  can be computed recursively as

$$\mathbf{p}_l^{(l)} = \left[ (\mathbf{p}_{l-1}^{(l-1)})^T \frac{1}{N} \sum_{k=1}^N K_{\sigma_l}(\mathbf{x}_k, \mathbf{x}'_l) \right]^T, \quad (13)$$

$$\mathbf{C}_l^{(l)} = \begin{bmatrix} \mathbf{C}_{l-1}^{(l-1)} & \mathbf{b}_{l-1}^{(l)} \\ (\mathbf{b}_{l-1}^{(l)})^T & \gamma_l \end{bmatrix}, \quad (14)$$

and

$$\begin{cases} \gamma_l = 1/(4\pi\sigma_l^2)^{m/2}, \\ \mathbf{b}_{l-1}^{(l)} = [K_{\sigma_{1,l}}(\mathbf{x}'_1, \mathbf{x}'_l) \cdots K_{\sigma_{l-1,l}}(\mathbf{x}'_{l-1}, \mathbf{x}'_l)]^T. \end{cases} \quad (15)$$

This recursion is initialized at the first step ( $l = 1$ ) as

$$\mathbf{C}_1^{(1)} = K_{\sqrt{2}\sigma_1}(\mathbf{x}'_1, \mathbf{x}'_1) = \gamma_1 \quad (16)$$

and

$$\mathbf{p}_1^{(1)} = \frac{1}{N} \sum_{k=1}^N K_{\sigma_1}(\mathbf{x}_k, \mathbf{x}'_1). \quad (17)$$

By substituting (7) and (12)-(14) into (11), we have

$$\begin{aligned} Q^{(l)}(\lambda_l, \sigma_l) &= \\ & \begin{bmatrix} \lambda_l \boldsymbol{\beta}_{l-1}^{(l-1)} \\ 1 - \lambda_l \end{bmatrix}^T \begin{bmatrix} \mathbf{C}_{l-1}^{(l-1)} & \mathbf{b}_{l-1}^{(l)} \\ (\mathbf{b}_{l-1}^{(l)})^T & \gamma_l \end{bmatrix} \begin{bmatrix} \lambda_l \boldsymbol{\beta}_{l-1}^{(l-1)} \\ 1 - \lambda_l \end{bmatrix} \\ & - 2[\lambda_l (\boldsymbol{\beta}_{l-1}^{(l-1)})^T \quad 1 - \lambda_l] \begin{bmatrix} \mathbf{p}_{l-1}^{(l-1)} \\ \frac{1}{N} \sum_{k=1}^N K_{\sigma_l}(\mathbf{x}_k, \mathbf{x}'_l) \end{bmatrix} \\ & = \lambda_l^2 \mu^{(l-1)} + (1 - \lambda_l)^2 \gamma_l + 2\lambda_l(1 - \lambda_l) (\mathbf{b}_{l-1}^{(l)})^T \boldsymbol{\beta}_{l-1}^{(l-1)} \\ & - 2\lambda_l \nu^{(l-1)} - \frac{2(1 - \lambda_l)}{N} \sum_{k=1}^N K_{\sigma_l}(\mathbf{x}_k, \mathbf{x}'_l). \end{aligned} \quad (18)$$

For  $l > 1$ ,  $Q^{(l)}(\lambda_l, \sigma_l)$  is a quadratic function with respect to  $\lambda_l$ . Hence there exists a unique minimum of  $Q^{(l)}(\lambda_l, \sigma_l)$  for a given  $\sigma_l$ , which can be found by setting  $\frac{\partial}{\partial \lambda_l} Q^{(l)}(\lambda_l, \sigma_l) = 0$ , followed by the constraint satisfaction operation. This yields the closed-form solution for  $\lambda_l$  for the given  $\sigma_l$  as

$$\lambda_l = \min \{ \max \{ u_l, 0 \}, 1 \}, \quad (19)$$

with

$$u_l = \frac{\gamma_l - (\mathbf{b}_{l-1}^{(l)})^T \boldsymbol{\beta}_{l-1}^{(l-1)} + \nu^{(l-1)} - \frac{1}{N} \sum_{k=1}^N K_{\sigma_l}(\mathbf{x}_k, \mathbf{x}'_l)}{\mu^{(l-1)} + \gamma_l - 2(\mathbf{b}_{l-1}^{(l)})^T \boldsymbol{\beta}_{l-1}^{(l-1)}}. \quad (20)$$

It is easy to verify that the constraint satisfaction operator

$$\min \{ \max \{ u, 0 \}, 1 \} = \begin{cases} 1, & u > 1, \\ 0, & u < 0, \\ u, & 0 \leq u \leq 1. \end{cases} \quad (21)$$

Therefore,  $0 \leq \lambda_l \leq 1$  is guaranteed. By plugging  $\lambda_l$  back to (18), we obtain the MISE value  $Q^{(l)}(\lambda_l, \sigma_l)$  for this given kernel. The computational cost of parameter estimation for a kernel with fixed width is in the order of  $O(l)$ , which is extremely low, owing to the recursive computation and the closed-form solution for the parameter  $\lambda_l$  when  $\sigma_l$  is fixed.

### B. Kernel width optimization with MISE criterion

We now consider the problem of optimizing  $K_{\sigma_l}(\mathbf{x}, \mathbf{x}'_l)$  by adjusting  $\sigma_l$ , also based on the MISE, when  $\lambda_l$  is fixed. Express (18) as

$$Q^{(l)}(\lambda, \sigma_l) = \lambda_l^2 \mu^{(l-1)} - 2\lambda_l \nu^{(l-1)} + S^{(l)}(\lambda_l, \sigma_l), \quad (22)$$

where

$$S^{(l)}(\lambda_l, \sigma_l) = 2\lambda_l(1 - \lambda_l) \sum_{i=1}^{l-1} \beta_i^{(l-1)} K_{\sigma_{i,l}}(\mathbf{x}'_i, \mathbf{x}'_l) + (1 - \lambda_l)^2 \gamma_l - \frac{2(1 - \lambda_l)}{N} \sum_{k=1}^N K_{\sigma_l}(\mathbf{x}_k, \mathbf{x}'_l) \quad (23)$$

which excludes all the components independent of  $\sigma_l$ . The gradient descent algorithm for minimizing  $S^{(l)}(\lambda_l, \sigma_l)$  and hence  $\text{ISE}^{(l)} = \text{ISE}(\lambda_l, \sigma_l)$  for the selected  $\mathbf{x}'_l$  and the fixed  $\lambda_l$  is given as follows.

Starting with  $\sigma_l^{\text{old}} = \sigma_0$ , repeat the following iterations for a sufficiently large number of times Iter, e.g., Iter = 20

$$\begin{cases} \sigma_l^{\text{new}} = \sigma_l^{\text{old}} - \eta \frac{\partial S^{(l)}(\lambda_l, \sigma_l^{\text{old}})}{\partial \sigma_l}, \\ \sigma_l^{\text{old}} \leftarrow \max\{\sigma_l^{\text{new}}, \sigma_{\min}\}, \end{cases} \quad (24)$$

where  $\eta > 0$  is a small positive learning rate,  $\sigma_{\min}$  is a small positive value representing the lower bound of the kernel width parameter, and the gradient is given by

$$\begin{aligned} \frac{\partial S^{(l)}(\lambda_l, \sigma_l)}{\partial \sigma_l} &= 2\lambda_l(1 - \lambda_l) \sum_{i=1}^{l-1} \beta_i^{(l-1)} \frac{\partial K_{\sigma_{i,l}}(\mathbf{x}'_i, \mathbf{x}'_l)}{\partial \sigma_l} \\ &\quad - \frac{m(1 - \lambda_l)^2 \gamma_l}{\sigma_l} - \frac{2(1 - \lambda_l)}{N} \sum_{k=1}^N \frac{\partial K_{\sigma_l}(\mathbf{x}_k, \mathbf{x}'_l)}{\partial \sigma_l} \\ &= 2\lambda_l(1 - \lambda_l) \sum_{i=1}^{l-1} \beta_i^{(l-1)} K_{\sigma_{i,l}}(\mathbf{x}'_i, \mathbf{x}'_l) \left( -\frac{m\sigma_l}{\sigma_i^2 + \sigma_l^2} + \right. \\ &\quad \left. \frac{\|\mathbf{x}'_i - \mathbf{x}'_l\|^2 \sigma_l}{(\sigma_i^2 + \sigma_l^2)^2} \right) - \frac{m(1 - \lambda_l)^2 \gamma_l}{\sigma_l} \\ &\quad - \frac{2(1 - \lambda_l)}{N} \sum_{k=1}^N K_{\sigma_l}(\mathbf{x}_k, \mathbf{x}'_l) \left( -\frac{m}{\sigma_l} + \frac{\|\mathbf{x}_k - \mathbf{x}'_l\|^2}{\sigma_l^3} \right). \end{aligned} \quad (25)$$

### C. Joint kernel selection and parameter estimation algorithm

At the  $l$ th forward selection stage, a data sample is to be selected from the remaining  $(N - l + 1)$  candidate data samples based on the fixed kernel width  $\sigma_0$ , while the associated kernel width  $\sigma_l$  is optimized, and the  $l$  kernel weights are adjusted. More specifically, we initially review the contribution of each candidate data sample according to its associated MISE value, based on the fixed kernel width  $\sigma_0$ , and decide which is to be added to the model. The data point producing the smallest MISE value amongst all the candidate data samples is selected as  $\mathbf{x}'_l$ . With the kernel weights being fixed, we then adjust the kernel width  $\sigma_l$  using the gradient descent algorithm described in Section III-B. Finally, the optimal kernel weights are recalculated for the given  $\sigma_l$  as described in Section III-A.

First define  $\mathbf{X}_N^{(l-1)} \in \mathbb{R}^{m \times N}$  as

$$\mathbf{X}_N^{(l-1)} = [\mathbf{x}'_1 \cdots \mathbf{x}'_{l-1} \quad \mathbf{x}_l^{(l-1)} \cdots \mathbf{x}_N^{(l-1)}], \quad (26)$$

and  $\mathbf{q}_N^{(l-1)} \in \mathbb{R}^{1 \times N}$  as

$$\mathbf{q}_N^{(l-1)} = \left[ \frac{1}{N} \sum_{k=1}^N K_{\sigma_0}(\mathbf{x}_k, \mathbf{x}'_1) \cdots \frac{1}{N} \sum_{k=1}^N K_{\sigma_0}(\mathbf{x}_k, \mathbf{x}'_{l-1}) \right. \\ \left. \frac{1}{N} \sum_{k=1}^N K_{\sigma_0}(\mathbf{x}_k, \mathbf{x}_l^{(l-1)}) \cdots \frac{1}{N} \sum_{k=1}^N K_{\sigma_0}(\mathbf{x}_k, \mathbf{x}_N^{(l-1)}) \right], \quad (27)$$

with

$$\mathbf{X}_N^{(0)} = [\mathbf{x}_1^{(0)} \ \mathbf{x}_2^{(0)} \ \cdots \ \mathbf{x}_N^{(0)}] = [\mathbf{x}_1 \ \mathbf{x}_2 \ \cdots \ \mathbf{x}_N], \quad (28)$$

$$\mathbf{q}_N^{(0)} = \left[ \frac{1}{N} \sum_{k=1}^N K_{\sigma_0}(\mathbf{x}_k, \mathbf{x}_1) \ \frac{1}{N} \sum_{k=1}^N K_{\sigma_0}(\mathbf{x}_k, \mathbf{x}_2) \cdots \right. \\ \left. \frac{1}{N} \sum_{k=1}^N K_{\sigma_0}(\mathbf{x}_k, \mathbf{x}_N) \right]. \quad (29)$$

If the  $j_l$ th column, where  $l \leq j_l \leq N$ , and the  $l$ th column of  $\mathbf{X}_N^{(l-1)}$  are interchanged,  $\mathbf{X}_N^{(l-1)}$  becomes  $\mathbf{X}_N^{(l)}$ . Similarly, if the  $j_l$ th column and the  $l$ th column of  $\mathbf{q}_N^{(l-1)}$  are interchanged,  $\mathbf{q}_N^{(l-1)}$  becomes  $\mathbf{q}_N^{(l)}$ . Further define the  $j$ th element of  $\mathbf{q}_N^{(l-1)}$  as  $q^{(l-1)}(j) = \frac{1}{N} \sum_{k=1}^N K_{\sigma_0}(\mathbf{x}_k, \mathbf{x}_j^{(l-1)})$  for  $l \leq j \leq N$ . We are now ready to present our proposed algorithm.

*Initialization:* At the 1st stage of the selection procedure, set  $\beta_1^{(1)} = \beta_1^{(1)} = 1$  and  $\lambda_1 = 0$ .

Step 1). For  $1 \leq j \leq N$ , based on  $\sigma_0$ , compute

$$Q^{(1,j)} = \gamma - 2\mathbf{p}_1^{(1,j)}, \quad (30)$$

where  $\gamma = \frac{1}{(4\pi\sigma_0^2)^{m/2}}$  and  $\mathbf{p}_1^{(1,j)} = q^{(0)}(j)$ .

Step 2). Find

$$Q^{(1,j_1)} = \min \{Q^{(1,j)}, 1 \leq j \leq N\}. \quad (31)$$

Then the  $j_1$ th column and the first column of  $\mathbf{X}_N^{(0)}$  are interchanged to yield  $\mathbf{X}_N^{(1)}$ , and the  $j_1$ th column and the first column of  $\mathbf{q}_N^{(0)}$  are interchanged to yield  $\mathbf{q}_N^{(1)}$ . This effectively selects the

first kernel.

Step 3). Apply (24) to find  $\sigma_1$ .

Step 4). Calculate  $\mu^{(1)} = \mathbf{C}_1^{(1)}$  and  $\nu^{(1)} = \mathbf{p}_1^{(1)}$  using (16) and (17). Update  $Q^{(1)} = \mu^{(1)} - 2\nu^{(1)}$ .

*The  $l$ th stage of the selection procedure, where  $l \geq 2$ :*

Step 1). For  $l \leq j \leq N$ , set  $\sigma_j = \sigma_0$ , compute

$$\begin{aligned} \mathbf{b}_{l-1}^{(l,j)} &= [K_{\sigma_{1,j}}(\mathbf{x}'_1, \mathbf{x}_j^{(l-1)}) \cdots K_{\sigma_{l-1,j}}(\mathbf{x}'_{l-1}, \mathbf{x}_j^{(l-1)})]^\top, \\ d^{(l,j)} &= (\mathbf{b}_{l-1}^{(l,j)})^\top \boldsymbol{\beta}_{l-1}^{(l-1)}, \\ \lambda_l^{(j)} &= \min \left\{ \max \left\{ \frac{\gamma - d^{(l,j)} + \nu^{(l-1)} - q^{(l-1)}(j)}{\mu^{(l-1)} + \gamma - 2d^{(l,j)}}, 0 \right\}, 1 \right\} \end{aligned}$$

and

$$\begin{aligned} Q^{(l,j)}(\lambda_l^{(j)}) &= (\lambda_l^{(j)})^2 \mu^{(l-1)} + (1 - \lambda_l^{(j)})^2 \gamma + \\ & 2\lambda_l^{(j)}(1 - \lambda_l^{(j)})d^{(l,j)} - 2\lambda_l^{(j)}\nu^{(l-1)} - 2(1 - \lambda_l^{(j)})q^{(l-1)}(j). \end{aligned}$$

Step 2): Find

$$Q^{(l,j_i)} = \min \{Q^{(l,j)}, l \leq j \leq N\}. \quad (32)$$

Then the  $j_l$ th column and the  $l$ th column of  $\mathbf{X}_N^{(l-1)}$  are interchanged to yield  $\mathbf{X}_N^{(l)}$ . Also the  $j_l$ th column and the  $l$ th column of  $\mathbf{q}_N^{(l-1)}$  are interchanged to yield  $\mathbf{q}_N^{(l)}$ . This effectively selects the  $l$ th kernel.

Step 3): With  $\lambda_l = \lambda_l^{(j_i)}$ , calculate  $\boldsymbol{\beta}_l^{(l)}$  using (7). Then apply (24) to find  $\sigma_l$ .

Step 4). Update  $\mathbf{p}_l^{(l)}$ ,  $\mathbf{C}_l^{(l)}$ . Recalculate  $\lambda_l$  using (19) and (20). Recalculate  $\boldsymbol{\beta}_l^{(l)}$  using (7). Update  $\mu^{(l)}$ ,  $\nu^{(l)}$  and  $Q^{(l)}(\lambda_l, \sigma_l)$  using (11)-(14).

*Termination:* The selection procedure is terminated at the  $(M + 1)$ th stage when the following condition is satisfied

$$|Q^{(M+1)} - Q^{(M)}| \leq \delta Q,$$

where  $\delta Q$  is a predetermined very small positive number, and this produces a subset model with the  $M$  significant kernels.

TABLE I  
COMPUTATIONAL COST OF THE PROPOSED ALGORITHM IN COMPARISON TO THE FCR-MISE ALGORITHM AT THE  $l$ TH FORWARD STEP.

Method	Kernel selection	Kernel width tuning	kernel weight re-estimation
FCR-MISE	$(N - l + 1)\mathcal{O}(l)$	none	none
The proposed	$(N - l + 1)\mathcal{O}(l)$	Iter $\times \mathcal{O}(N + l)$	$\mathcal{O}(N)$

#### D. Remarks:

*Remark 1:* The reason that the optimisation of (18) with respect to  $\lambda_l$  and  $\sigma_l$  is carried out separately is that the optimal value  $\lambda_l$  can be expressed in closed form for fixed  $\sigma_l$ , thus significantly reducing computational costs. Alternatively both of them can be optimized using gradient descent algorithm simultaneously. Since the relationship with respect to  $\sigma_l$  is not quadratic, the results will not be the same, each only achieving a local minimum. However the property that (18) is quadratic in  $\lambda_l$  cannot be exploited for computational advantage.

*Remark 2:* In FCR-MISE algorithm [19], each kernel has a common fixed width, and appropriate kernel value can be determined empirically through trial and error based on cross-validation. More specifically, a suitable kernel width value can be found using a line search based on the cross-validation performance. In the proposed algorithm, the kernel width is given as a  $\sigma_M$ , where each element in  $\sigma_M$  is optimized from an initial  $\sigma_0$  which needs to be preset. Unlike the fixed kernel width in FCR-MISE algorithm [19], the choice of  $\sigma_0$  is more relaxed, since there is a wide range of suitable values.

#### E. Computational cost

The proposed algorithm is an extension to the low cost FCR-MISE algorithm [19], with the difference that each kernel is tuned after it has been selected. The FCR-MISE algorithm [19] has a significant advantage in that it offers a much lower complexity in constructing PDF estimate than other existing sparse estimators with  $\mathcal{O}(N^2)$  complexity. Table I compares the computational cost of the proposed algorithm with that of the FCR-MISE algorithm at the  $l$ th forward step. Overall the computational cost is increased at each forward stage, compared to the FCR-MISE algorithm. Since the tuning of the kernel is only applied to the selected kernel, the extra cost is small. In contrast to our proposed algorithm which automatically tunes each kernel width, however, there exists extra computational cost for any estimator based on a pre-set fixed single

kernel width, such as the FCR-MISE algorithm, since this kernel width has to be empirically tuned outside the algorithm loop. Moreover, the total computational cost of an algorithm is dependent on the model size  $M$  of the final selected model. Since  $M$  is usually much smaller than  $N$ , the total computational cost is approximately linear with respect to the model size  $M$ . Since our proposed algorithm can produce a much smaller model, its total computational cost can actually be lower than that of the FCR-MISE algorithm.

#### IV. SIMULATION STUDY

Two numerical examples are provided. In each example, we randomly draw a data set of  $N$  points from a known distribution  $p(\mathbf{x})$  to construct the PDF estimate  $\hat{p}^{(M)}(\mathbf{x}_k; \boldsymbol{\beta}_M, \boldsymbol{\sigma}_M)$  based on the proposed approach. A separate test data set of  $N_{\text{test}} = 10000$  sample points was used for evaluation according to the  $L_1$  norm

$$L_1 = \frac{1}{N_{\text{test}}} \sum_{k=1}^{N_{\text{test}}} |p(\mathbf{x}_k) - \hat{p}^{(M)}(\mathbf{x}_k; \boldsymbol{\beta}_M, \boldsymbol{\sigma}_M)|. \quad (33)$$

The experiment was repeated for 100 different random runs.

*Example 1:* The density to be estimated for this 2-dimensional (2-D) example was given by the mixture of two densities, a Gaussian and a Laplacian, as defined by

$$p(\mathbf{x}) = \frac{1}{4\pi} \exp\left(-\frac{(x_1 - 2)^2}{2}\right) \exp\left(-\frac{(x_2 - 2)^2}{2}\right) + \frac{0.35}{8} \exp(-0.7|x_1 + 2|) \exp(-0.5|x_2 + 2|). \quad (34)$$

The estimation data set had  $N = 500$  samples.

*Example 2:* The density to be estimated for this 6-D example was the mixture of three Gaussians defined by

$$p(\mathbf{x}) = \frac{1}{3} \sum_{i=1}^3 \frac{1}{(2\pi)^3 \sqrt{\det(\boldsymbol{\Gamma}_i)}} \times \exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_i)^T \boldsymbol{\Gamma}_i^{-1} (\mathbf{x} - \boldsymbol{\mu}_i)\right), \quad (35)$$

with  $\boldsymbol{\mu}_1 = [1 \ 1 \ 1 \ 1 \ 1 \ 1]^T$ ,  $\boldsymbol{\mu}_2 = [-1 \ -1 \ -1 \ -1 \ -1 \ -1]^T$ ,  $\boldsymbol{\mu}_3 = [0 \ 0 \ 0 \ 0 \ 0 \ 0]^T$ ,  $\boldsymbol{\Gamma}_1 = \text{diag}\{1, 2, 1, 2, 1, 2\}$ ,  $\boldsymbol{\Gamma}_2 = \text{diag}\{2, 1, 2, 1, 2, 1\}$ , and  $\boldsymbol{\Gamma}_3 = \text{diag}\{2, 1, 2, 1, 2, 1\}$ . The estimation data set had  $N = 600$  samples.

TABLE II  
PERFORMANCE COMPARISON OF KERNEL DENSITY ESTIMATORS.

(a) Example 1

Method	$L_1$ test error (mean $\pm$ STD)	Kernel number (mean $\pm$ STD)
PW	$(4.18 \pm 0.8) \times 10^{-3}$	$500 \pm 0$
SDC [17]	$(3.83 \pm 0.8) \times 10^{-3}$	$11.9 \pm 2.6$
SKD [18]	$(3.84 \pm 0.8) \times 10^{-3}$	$15.3 \pm 3.9$
RSDE-MNQP [11]	$(4.24 \pm 0.8) \times 10^{-3}$	$129.4 \pm 35.7$
FCR-MISE [19]	$(3.33 \pm 0.8) \times 10^{-3}$	$25.1 \pm 2.7$
RTR-MISE [14]	$(3.13 \pm 0.7) \times 10^{-3}$	$36.7 \pm 11.3$
The proposed	$(3.57 \pm 0.7) \times 10^{-3}$	$7.6 \pm 1.4$

(b) Example 2

Method	$L_1$ test error (mean $\pm$ STD)	Kernel number (mean $\pm$ STD)
PW	$(3.18 \pm 0.13) \times 10^{-5}$	$600 \pm 0$
SDC [17]	$(4.48 \pm 1.2) \times 10^{-5}$	$14.9 \pm 2.1$
SKD [18]	$(3.11 \pm 0.5) \times 10^{-5}$	$9.4 \pm 1.9$
RSDE-MNQP [11]	$(3.67 \pm 0.7) \times 10^{-5}$	$29.4 \pm 10.1$
FCR-MISE [19]	$(2.82 \pm 0.1) \times 10^{-5}$	$19.4 \pm 0.9$
RTR-MISE [14]	$(2.53 \pm 0.1) \times 10^{-5}$	$81.2 \pm 20$
The proposed	$(2.64 \pm 0.2) \times 10^{-5}$	$2.9 \pm 0.2$

Six methods were used for comparison: (a) the well known PW estimate; (b) the sparse density construction (SDC) algorithm [17]; (c) the sparse kernel density construction (SKD) algorithm [18]; (d) the reduced set density estimator with multiplicative nonnegative quadratic programming (RSDE-MNQP) [11]; (e) the FCR-MISE algorithm [19]; and (f) the RTR-MISE algorithm [14].

We briefly explain these six algorithms. Both the SDC algorithm [17] and the SKD algorithm [18] are regression-based PDF estimation methods that construct sparse PDF forwardly. For the SDC algorithm, the empirical CDF is constructed and used as the desired response, but for the SKD algorithm the PW estimate is constructed and used as the desired response. The RSDE-MNQP [11], the FCR-MISE [19] and the RTR-MISE [14] are all based on the MISE, but employ different optimization algorithms. Specifically, the RSDE-MNQP algorithm uses the MNQP algorithm, the FCR-MISE algorithm formulates the density estimation in a forward constrained regression manner by selecting one kernel at a time forwardly, and the RTR-MISE algorithm is based on the Riemannian trust-region algorithm [13]. We also point out that the MISE cost function is used in PW estimate using grid search for an optimal kernel width. However, the single kernel width for the other five algorithms needs to be preset empirically.



The algorithmic parameters of the proposed approach were set to  $\sigma_{\min} = 0.1$  and  $\sigma_{\max} = 1$  for *Example 1* and *Example 2*, respectively,  $\text{Iter} = 20$  and  $\eta = 0.02$  for the both examples, while  $\delta Q$  was set to  $10^{-4}$  and  $10^{-5}$ , respectively for the two example. The results obtained by the seven kernel density estimators are listed in Table II (a) and (b), respectively, for the two examples, where the results of the SDC, SKD, FCR-MISE and RTR-MISE are quoted from [14], [17]–[19], respectively. The results of Table II clearly show that our proposed algorithm can construct much sparser kernel density estimates than the five state-of-the-art benchmark sparse kernel density estimators compared, with competitive accuracy. Compared to the low cost FCR-MISE algorithm, the proposed algorithm increases the computational complexity per forward step of  $\text{Iter} \times \mathcal{O}(N + l)$  due to the tunable kernel calculation. However it is clear that the resultant models are much sparser leading to fewer forward regression steps for computational cost reduction. Note that the computational costs of [19] have already been shown to be better than the other algorithms.

## V. CONCLUSIONS

We have introduced a new sparse kernel density estimator with tunable kernels based on the idea of forward constrained regression by adding one kernel at a time based on the minimum ISE criterion. Our main contribution has been to develop a new recursive algorithm which selects a significant kernel at each forward construction stage, and then optimizes the kernel width of the selected kernels based on the gradient descent algorithm. The significant advantages of the proposed method are that it is able to obtain very sparse PDF estimates due to the individually tunable kernel width parameters, and it requires no empirically predetermined parameters outside the algorithm. Numerical examples have been employed to demonstrate that the proposed approach can construct *very sparse* kernel density estimators with competitive accuracy to the existing state-of-the-art sparse kernel density estimators.

## APPENDIX

INTEGRATING  $\int K_{\sigma_i}(\mathbf{x}, \mathbf{x}'_i) K_{\sigma_j}(\mathbf{x}, \mathbf{x}'_j) d\mathbf{x}$

With the notations  $\mathbf{x} = [x_1 \ x_2 \ \cdots \ x_m]^T$  and  $\mathbf{x}'_i = [x'_{i,1} \ x'_{i,2} \ \cdots \ x'_{i,m}]^T$  for  $1 \leq i \leq l$ , we have

$$\begin{aligned} \int K_{\sigma_i}(\mathbf{x}, \mathbf{x}'_i) K_{\sigma_j}(\mathbf{x}, \mathbf{x}'_j) d\mathbf{x} &= \frac{1}{(2\pi\sigma_i\sigma_j)^m} \\ &\times \prod_{k=1}^m \int \exp\left(-\frac{(x_k - x'_{i,k})^2}{2\sigma_i^2} - \frac{(x_k - x'_{j,k})^2}{2\sigma_j^2}\right) dx_k \end{aligned} \quad (36)$$

in which

$$\begin{aligned} &\int \exp\left(-\frac{(x_k - x'_{i,k})^2}{2\sigma_i^2} - \frac{(x_k - x'_{j,k})^2}{2\sigma_j^2}\right) dx_k \\ &\quad (\sigma_i^2 + \sigma_j^2)x_k^2 - 2(x'_{i,k}\sigma_j^2 + x'_{j,k}\sigma_i^2)x_k \\ &= \int \exp\left(-\frac{(\sigma_i^2 + \sigma_j^2)x_k^2 - 2(x'_{i,k}\sigma_j^2 + x'_{j,k}\sigma_i^2)x_k + (x'_{i,k})^2\sigma_j^2 + (x'_{j,k})^2\sigma_i^2}{2\sigma_i^2\sigma_j^2}\right) dx_k \\ &= \exp\left(-\frac{\frac{(x'_{i,k})^2\sigma_j^2 + (x'_{j,k})^2\sigma_i^2}{\sigma_i^2 + \sigma_j^2} - \left(\frac{x'_{i,k}\sigma_j^2 + x'_{j,k}\sigma_i^2}{\sigma_i^2 + \sigma_j^2}\right)^2}{2\sigma_i^2\sigma_j^2/(\sigma_i^2 + \sigma_j^2)}\right) \\ &\times \int \exp\left(-\frac{\left(x_k - \frac{x'_{i,k}\sigma_j^2 + x'_{j,k}\sigma_i^2}{\sigma_i^2 + \sigma_j^2}\right)^2}{2\sigma_i^2\sigma_j^2/(\sigma_i^2 + \sigma_j^2)}\right) dx_k \\ &= \exp\left(-\frac{(x'_{i,k} - x'_{j,k})^2}{2(\sigma_i^2 + \sigma_j^2)}\right) \\ &\times \int \exp\left(-\frac{\left(x_k - \frac{x'_{i,k}\sigma_j^2 + x'_{j,k}\sigma_i^2}{\sigma_i^2 + \sigma_j^2}\right)^2}{2\sigma_i^2\sigma_j^2/(\sigma_i^2 + \sigma_j^2)}\right) dx_k. \end{aligned} \quad (37)$$

Noting  $\int \frac{1}{\sqrt{2\pi s}} \exp\left(-\frac{(x-\mu)^2}{2s}\right) dx = 1$ , we have

$$\begin{aligned} &\int \exp\left(-\frac{(x_k - x'_{i,k})^2}{2\sigma_i^2} - \frac{(x_k - x'_{j,k})^2}{2\sigma_j^2}\right) dx_k \\ &= \sqrt{2\pi\sigma_i^2\sigma_j^2/(\sigma_i^2 + \sigma_j^2)} \exp\left(-\frac{(x'_{i,k} - x'_{j,k})^2}{2(\sigma_i^2 + \sigma_j^2)}\right) \end{aligned} \quad (38)$$

so that

$$\begin{aligned} &\int K_{\sigma_i}(\mathbf{x}, \mathbf{x}'_i) K_{\sigma_j}(\mathbf{x}, \mathbf{x}'_j) d\mathbf{x} \\ &= \frac{1}{(2\pi\sigma_{i,j}^2)^{m/2}} \exp\left(-\frac{\|\mathbf{x}'_i - \mathbf{x}'_j\|^2}{2\sigma_{i,j}^2}\right) \\ &= K_{\sigma_{i,j}}(\mathbf{x}'_i, \mathbf{x}'_j) \end{aligned} \quad (39)$$

with  $\sigma_{i,j} = \sqrt{\sigma_i^2 + \sigma_j^2}$ .

## REFERENCES

- [1] G. McLachlan and D. Peel, *Finite Mixture Models*. Wiley: New York, 2000.
- [2] B. W. Silverman, *Density Estimation for Statistics and Data Analysis*. Chapman and Hall: London, 1986.
- [3] R. O. Duda and P. E. Hart, *Pattern Classification and Scene Analysis*. Wiley: New York, 1973.
- [4] S. Chen, X. Hong, and C. J. Harris, "Particle swarm optimization aided orthogonal forward regression for unified data modelling," *IEEE Trans. Evolutionary Computation*, vol. 14, no. 4, pp. 477–499, Aug. 2010.
- [5] L. Rutkowski, "Adaptive probabilistic neural networks for pattern classification in time-varying environment," *IEEE Trans. Neural Networks*, vol. 15, no. 4, pp. 811–827, July 2004.
- [6] H. Yin and N. W. Allinson, "Self-organizing mixture networks for probability density estimation," *IEEE Trans. Neural Networks*, vol. 12, no. 2, pp. 405–411, March 2001.
- [7] A. P. Dempster, N. M. Laird, and D. B. Rubin, "Maximum likelihood from incomplete data via the EM algorithm," *J. Royal Statistical Society B*, vol. 39, no. 1, pp. 1–38, 1977.
- [8] E. Parzen, "On estimation of a probability density function and mode," *Annals of Mathematical Statistics*, vol. 33, no. 3, pp. 1066–1076, Sept. 1962.
- [9] J. Weston, A. Gammerman, M. O. Stitson, V. Vapnik, V. Vovk, and C. Watkins, "Support vector density estimation," in *Advances in Kernel Methods – Support Vector Learning*, B. Schölkopf, C. Burges, and A. J. Smola, Eds. MIT Press: Cambridge, MA, 1999, pp. 293–306.
- [10] V. Vapnik and S. Mukherjee, "Support vector method for multivariate density estimation," in *Advances in Neural Information Processing Systems*, S. Solla, T. Leen, and K. R. Müller, Eds. MIT Press: Cambridge, MA, 2000, pp. 659–665.
- [11] M. Girolami and C. He, "Probability density estimation from optimally condensed data samples," *IEEE Trans. Pattern Analysis and Machine Intelligence*, vol. 25, no. 10, pp. 1253–1264, Oct. 2003.
- [12] S. W. Scott, "Parametric statistical modeling by minimum integrated square error," *Technometrics*, vol. 43, no. 3, pp. 274–285, Aug. 2001.
- [13] N. Boumal, B. Mishra, P.-A. Absil, and R. Sepulchre, "Manopt, a Matlab toolbox for optimization on manifolds," *J. Machine Learning Research*, vol. 15, pp. 1455–1450, April 2014.
- [14] X. Hong, J. B. Gao, S. Chen, and T. Zia, "Sparse density estimation on the multinomial manifold," *IEEE Trans. Neural Networks and Learning Systems*, To appear, 2015.
- [15] A. Choudhury, *Fast Machine Learning Algorithms for Large Data*, Ph.D. dissertation, School of Engineering Sciences, University of Southampton, 2002.
- [16] S. Chen, X. Hong, C. J. Harris, and P. M. Sharkey, "Sparse modeling using forward regression with PRESS statistic and regularization," *IEEE Trans. Systems, Man and Cybernetics, Part B*, vol. 34, no. 2, pp. 898–911, April 2004.
- [17] S. Chen, X. Hong, and C. J. Harris, "Sparse kernel density construction using orthogonal forward regression with leave-one-out test score and local regularization," *IEEE Trans. Systems, Man, and Cybernetics, Part B*, vol. 34, no. 4, pp. 1708–1717, Aug. 2004.
- [18] S. Chen, X. Hong, and C. J. Harris, "An orthogonal forward regression techniques for sparse kernel density estimation," *Neurocomputing*, vol. 71, nos. 4-6, pp. 931–943, Jan. 2008.
- [19] X. Hong, S. Chen, A. Qataweh, K. Daqrouq, M. Sheikh, and A. Morfeq, "Sparse probability density function estimation using the minimum integrated square error," *Neurocomputing*, vol. 115, pp. 122–129, 2013.

- [20] X. Hong and C. J. Harris, "A mixture of experts network structure construction algorithm for modelling and control," *Applied Intelligence*, vol. 16, no. 1, pp. 59–69, 2002.

Accepted manuscript