

# Direct conditional probability density estimation with sparse feature selection

Motoki Shiga · Voot Tangkaratt · Masashi Sugiyama

Received: 13 February 2014 / Accepted: 15 October 2014 / Published online: 15 January 2015  
© The Author(s) 2015

**Abstract** Regression is a fundamental problem in statistical data analysis, which aims at estimating the conditional mean of output given input. However, regression is not informative enough if the conditional probability density is multi-modal, asymmetric, and heteroscedastic. To overcome this limitation, various estimators of conditional densities themselves have been developed, and a kernel-based approach called *least-squares conditional density estimation* (LS-CDE) was demonstrated to be promising. However, LS-CDE still suffers from large estimation error if input contains many irrelevant features. In this paper, we therefore propose an extension of LS-CDE called *sparse additive CDE* (SA-CDE), which allows automatic feature selection in CDE. SA-CDE applies kernel LS-CDE to each input feature in an additive manner and penalizes the whole solution by a group-sparse regularizer. We also give a subgradient-based optimization method for SA-CDE training that scales well to high-dimensional large data sets. Through experiments with benchmark and humanoid robot transition datasets, we demonstrate the usefulness of SA-CDE in noisy CDE problems.

**Keywords** Conditional density estimation · Feature selection · Sparse structured norm

---

Editors: Toon Calders, Rosa Meo, Floriana Esposito and Eyke Hüllermeier.

---

M. Shiga (✉)  
Gifu University, 1-1 Yanagido, Gifu-city, Gifu 501-1193, Japan  
e-mail: shiga\_m@gifu-u.ac.jp

V. Tangkaratt · M. Sugiyama  
Tokyo Institute of Technology, 2-12-1 O-okayama, Meguro-ku, Tokyo 152-8552, Japan  
e-mail: voot@sg.cs.titech.ac.jp

M. Sugiyama  
e-mail: sugi@cs.titech.ac.jp

## 1 Introduction

Estimating the statistical dependency between input  $\mathbf{x}$  and output  $\mathbf{y}$  plays a crucial role in various real-world applications. For example, in robot transition estimation which is highly useful in *model-based reinforcement learning* (Sutton and Barto 1998), input  $\mathbf{x}$  corresponds to the pair of the current state of a robot and an action the robot takes, and output  $\mathbf{y}$  corresponds to the destination state after taking the action. Another application is disease diagnosis, in which input  $\mathbf{x}$  corresponds to measurements of biomarkers and/or clinical images and output  $\mathbf{y}$  corresponds to the presence (or the progression level) of a disease. Thus, accurately estimating the statistical dependency is an important and fundamental problem in statistical data analysis. The most basic approach to this problem is regression, which estimates the conditional *mean* of output  $\mathbf{y}$  given input  $\mathbf{x}$ . Regression gives the optimal estimation of output  $\mathbf{y}$  for additive Gaussian output noise. However, if the conditional probability density of output  $\mathbf{y}$  given input  $\mathbf{x}$ , denoted by  $p(\mathbf{y}|\mathbf{x})$ , possesses more complex structure such as multimodality, asymmetry, and heteroscedasticity, estimating the conditional mean by regression is not necessarily informative.

To overcome the limitation of regression, estimation of conditional densities from paired samples  $\{(\mathbf{x}^{(n)}, \mathbf{y}^{(n)})\}_{n=1}^N$  has been investigated. The most naive approach to estimating  $p(\mathbf{y}|\mathbf{x} = \tilde{\mathbf{x}})$ , the conditional density of output  $\mathbf{y}$  at test input point  $\mathbf{x} = \tilde{\mathbf{x}}$ , is to use the *kernel density estimator* (KDE) (Silverman 1986) with samples such that  $\|\mathbf{x}^{(n)} - \tilde{\mathbf{x}}\|_2^2 \leq \epsilon$ . However, this naive method does not work well in high-dimensional problems. Slightly more sophisticated variants have been proposed that use weighted KDE (Fan et al. 1996; Wolff et al. 1999), but they still share the same weakness.

The *mixture density network* (MDN) (Bishop 2006) uses a mixture of parametric densities for modeling the conditional density, and the parameters are estimated by a neural network as functions of input  $\mathbf{x}$ . MDN was demonstrated to work well, but its training is time-consuming and only a local optimal solution may be found due to the non-convexity of neural network training. A similar method based on a mixture of Gaussian processes was developed (Tresp 2001), which can be trained in a computationally more efficient way by the expectation-maximization algorithm (Dempster et al. 1977). However, due to the non-convexity of the optimization problem, it is difficult to find the global optimal solution.

*Kernel quantile regression* (KQR) (Takeuchi et al. 2006; Li et al. 2007) gives non-parametric percentile estimates of conditional distributions through convex optimization. KQR can be used for estimating the entire conditional cumulative distribution by solving KQR for all percentiles. It was shown that the regularization path tracking technique (Hastie et al. 2004) can be employed for efficiently computing the entire conditional cumulative distribution (Takeuchi et al. 2009). However, KQR is applicable only to one-dimensional output, which limits the range of applications significantly.

*Least-squares conditional density estimation* (LS-CDE) allows estimation of multiple-input-multiple-output conditional densities by directly learning a conditional density model with least-squares estimation (Sugiyama et al. 2010). For linear-in-parameter models such as a linear combination of Gaussian kernels, LS-CDE is formulated as a convex optimization problem and its solution can be obtained efficiently and analytically just by solving a system of linear equations. Furthermore, kernel LS-CDE was proved to achieve the optimal non-parametric convergence rate to the true conditional density in the mini-max sense, meaning that no method can be better than LS-CDE asymptotically. Through extensive experiments, LS-CDE was demonstrated to compare favorably with competing approaches.

However, LS-CDE still suffers from large estimation error when many irrelevant features exist in input  $\mathbf{x}$ . Such irrelevant features are conceivable in many real-world problems. For

example, in gene expression analysis for diseased cells, only a small subset of biomarker genes (input) affects the disease status (output). A standard way to cope with high input dimensionality is to select relevant features with forward selection or backward elimination (Guyon and Elisseeff 2003), but this often leads to a local optimal set of features.

In this paper, we propose extending LS-CDE to allow simultaneous feature selection during conditional density estimation. More specifically, we apply kernel LS-CDE to each input feature in an additive manner and penalize the whole solution by a group-sparse regularizer (Yuan and Lin 2006). Our subgradient-based optimization solver allows computationally efficient selection of relevant features that are even non-linearly correlated with output  $\mathbf{y}$ . Numerical experiments on noisy conditional density estimation demonstrate that our proposed method, which we call *sparse additive CDE* (SA-CDE), compares favorably with baseline approaches in estimation accuracy and computational efficiency.

The remainder of this paper is structured as follows. In Sect. 2, we formulate the problem of conditional density estimation and describe our proposed SA-CDE method. We experimentally evaluate the performance of SA-CDE in Sect. 3, and we summarize our contribution in Sect. 4.

## 2 Conditional density estimation with sparse feature selection

In this section, we formulate the problem of conditional density estimation and describe our proposed SA-CDE method.

### 2.1 Problem formulation

Let

$$\mathbf{x} = (x_1, \dots, x_{D_x})^T \in \mathbb{R}^{D_x}$$

be an input vector and  $\mathbf{y} \in \mathbb{R}^{D_y}$  be an output vector, where  $\mathbb{R}$  is the set of all real numbers,  $D_x$  is the dimension of the input vector, and  $D_y$  is the dimension of the output vector. We are given i.i.d. input-output paired samples of size  $N$  following the joint probability distribution with density  $p(\mathbf{x}, \mathbf{y})$ :

$$(\mathbf{x}^{(1)}, \mathbf{y}^{(1)}), \dots, (\mathbf{x}^{(N)}, \mathbf{y}^{(N)}). \quad (1)$$

We assume that only some of the input features are relevant to output  $\mathbf{y}$ . Then, such relevant features are sufficient for predicting output  $\mathbf{y}$  (Li 1991; Cook and Ni 2005):

$$p(\mathbf{y}|\mathbf{x}) = p(\mathbf{y}|\mathbf{x}^+), \quad (2)$$

where  $\mathbf{x}^+$  is the sub-vector of  $\mathbf{x}$  that consists only of relevant features.

Our goal is to estimate the conditional density  $p(\mathbf{y}|\mathbf{x})$  from the training samples (1) via sufficient feature selection (2).

### 2.2 Sparse additive conditional density estimation

We use an additive conditional density model:

$$\hat{p}(\mathbf{y}|\mathbf{x}) := \sum_{d=1}^{D_x} \hat{r}_d(\mathbf{y}, x_d), \quad (3)$$

where  $\hat{r}_d(\mathbf{y}, x_d)$  is an unnormalized estimator with the  $d$ th input feature  $x_d$ . We use a linear combination of  $B$  basis functions as  $\hat{r}_d(\mathbf{y}, x)$ :

$$\begin{aligned}\hat{r}_d(\mathbf{y}, x) &:= \sum_{b=1}^B \alpha_{d,b} \left\{ \eta_b(\mathbf{y}) \cdot \varphi_{d,b}(x) \right\} \\ &= \boldsymbol{\alpha}_d^T \left\{ \boldsymbol{\eta}(\mathbf{y}) \circ \boldsymbol{\varphi}_d(x) \right\}, \quad d = 1, \dots, D_x,\end{aligned}\quad (4)$$

where

$$\boldsymbol{\alpha}_d = (\alpha_{d,1}, \dots, \alpha_{d,B})^T \in \mathbb{R}_+^B$$

is a parameter vector for the  $d$ th input  $x_d$ ,

$$\boldsymbol{\eta}(\mathbf{y}) = (\eta_1(\mathbf{y}), \dots, \eta_B(\mathbf{y}))^T$$

is a vector of basis functions for output  $\mathbf{y}$ ,

$$\boldsymbol{\varphi}_d(x) = (\varphi_{d,1}(x), \dots, \varphi_{d,B}(x))^T$$

is a vector of basis functions for the  $d$ th input  $x_d$ , and  $\circ$  denotes the Hadamard (or element-wise) product. We use the following Gaussian kernels as the basis functions:

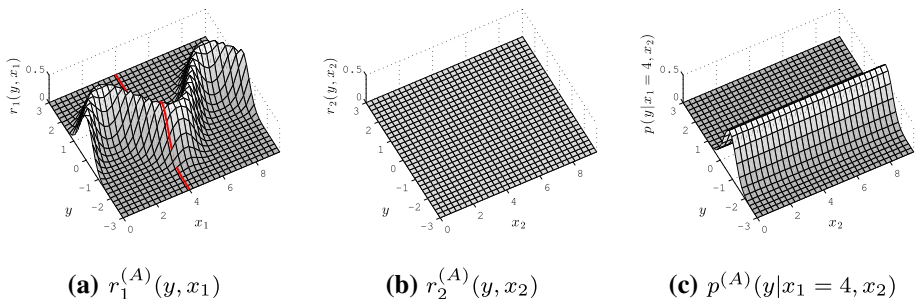
$$\eta_b(\mathbf{y}) := \exp\left(-\frac{\|\mathbf{y} - \mathbf{v}_b\|^2}{2\sigma^2}\right), \quad b = 1, \dots, B, \quad (5)$$

$$\varphi_{d,b}(x_d) := \exp\left(-\frac{(x_d - \mu_{d,b})^2}{2\sigma^2}\right), \quad d = 1, \dots, D_x, \quad b = 1, \dots, B, \quad (6)$$

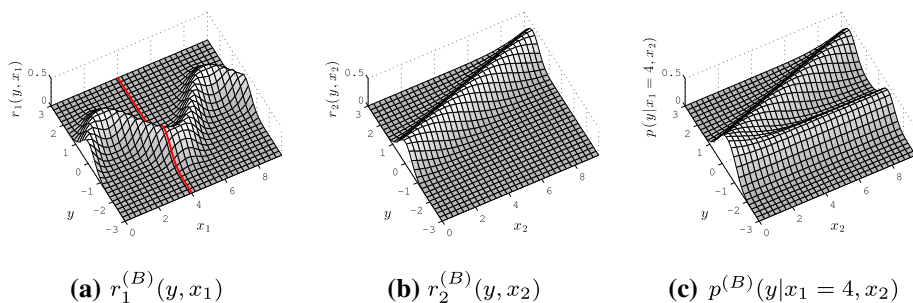
where  $\mathbf{v}_b$  and  $\mu_{d,b}$  are the Gaussian centers and  $\sigma$  is the Gaussian width. The Gaussian centers are fixed at the points chosen randomly from training samples. On the other hand, the Gaussian width will be optimized by cross-validation (see Sect. 2.5). For simplicity, we assume that input  $\mathbf{x}$  and output  $\mathbf{y}$  of the training samples (1) are both normalized element-wise to have the unit variance in advance and use the common Gaussian width  $\sigma$ . Note that the magnitude of parameter  $\boldsymbol{\alpha}_d$  may reflect the relevance between the  $d$ th feature and output  $\mathbf{y}$ .

Figures 1 and 2 illustrate two examples of the additive function model (3). The first example shown in Fig. 1 is

$$r_1^{(A)}(\mathbf{y}, x_1) = \mathcal{N}\left(y | \sin(x_1), \frac{1}{4}\right) \quad \text{and} \quad r_2^{(A)}(\mathbf{y}, x_2) = 0,$$



**Fig. 1** Conditional density function in which output  $\mathbf{y}$  is relevant to only input  $x_1$



**Fig. 2** Conditional density function in which output  $y$  is relevant to both input  $x_1$  and  $x_2$

while the second example shown in Fig. 2 is

$$r_1^{(B)}(y, x_1) = \frac{1}{2} \mathcal{N}\left(y | \sin(x_1), \frac{1}{4}\right) \quad \text{and} \quad r_2^{(B)}(y, x_2) = \frac{1}{2} \mathcal{N}\left(y | \frac{1}{4}x_2, \frac{1}{4}\right).$$

Note that  $r_2(y, x_2)$  is only different in these two examples. In the density function  $p^{(A)}(y|x_1, x_2)$  shown in Fig. 1c,  $x_2$  does not affect output  $y$ . On the other hand, in density function  $p^{(B)}(y|x_1, x_2)$  shown in 2c, output  $y$  varies depending on  $x_2$ . This illustrates that the magnitude of  $r_d^{(B)}(y, x_d)$ , which is controlled by  $\alpha_d$ , may reflect the relevance between the  $d$ th feature and output  $y$ .

Our optimization parameters to be learned in model (3) are

$$\alpha := (\alpha_1^T, \dots, \alpha_{D_x}^T)^T \in \mathbb{R}^{B \cdot D_x}.$$

We learn  $\alpha$  to minimize the squared error between true conditional density  $p(y|x)$  and our estimator  $\hat{p}(y|x)$

$$J_0(\alpha) := \frac{1}{2} \int \int \left( \hat{p}(y|x) - p(y|x) \right)^2 p(x) dy dx. \quad (7)$$

Substituting our model (3) and  $p(y, x) = p(y|x)p(x)$  into (7), we have

$$\begin{aligned} J_0(\alpha) = & \frac{1}{2} \sum_{d_1, d_2=1}^{D_x} \int \int \int \hat{r}_{d_1}(y, x_{d_1}) \hat{r}_{d_2}(y, x_{d_2}) p(x_{d_1}, x_{d_2}) dy dx_{d_1} dx_{d_2} \\ & - \sum_{d=1}^{D_x} \int \int \hat{r}_d(y, x_d) p(y, x_d) dy dx_d + \text{Const.} \end{aligned} \quad (8)$$

Using empirical approximation and ignoring the constant term, we can approximate the loss function as follows (the detailed derivation is described in “Appendix 1”):

$$\begin{aligned} \hat{J}_0(\alpha) = & \frac{1}{2N} \sum_{d_1, d_2=1}^{D_x} \sum_{n=1}^N \int \hat{r}_{d_1}(y, x_{d_1}^{(n)}) \hat{r}_{d_2}(y, x_{d_2}^{(n)}) dy \\ & - \frac{1}{N} \sum_{d=1}^{D_x} \sum_{n=1}^N \hat{r}_d(y^{(n)}, x_d^{(n)}) \end{aligned} \quad (9)$$

$$= \frac{1}{2} \alpha^T H \alpha - h^T \alpha, \quad (10)$$

where  $H \in \mathbb{R}^{B \cdot D_x \times B \cdot D_x}$  and  $h \in \mathbb{R}^{B \cdot D_x}$  are given by

$$\mathbf{H} := \begin{pmatrix} \mathbf{H}_{1,1} & \cdots & \mathbf{H}_{1,D_x} \\ \vdots & \ddots & \vdots \\ \mathbf{H}_{D_x,1} & \cdots & \mathbf{H}_{D_x,D_x} \end{pmatrix}. \quad (11)$$

$$\begin{aligned} [\mathbf{H}_{d_1,d_2}]_{b_1,b_2} &:= \frac{(\sqrt{\pi}\sigma)^{D_y}}{N} \exp\left(-\frac{\|\mathbf{v}_{b_1} - \mathbf{v}_{b_2}\|^2}{4\sigma^2}\right) \\ &\quad \times \sum_{n=1}^N \exp\left(-\frac{(x_{d_1}^{(n)} - \mu_{d_1,b_1})^2}{2\sigma^2} - \frac{(x_{d_2}^{(n)} - \mu_{d_2,b_2})^2}{2\sigma^2}\right), \\ &\quad d_1, d_2 = 1, \dots, D_x, \quad b_1, b_2 = 1, \dots, B, \end{aligned} \quad (12)$$

$$\mathbf{h} := (\mathbf{h}_1^T, \dots, \mathbf{h}_{D_x}^T)^T, \quad (13)$$

$$\begin{aligned} h_{d,b} &:= \frac{1}{N} \sum_{n=1}^N \exp\left(-\frac{\|\mathbf{y}^{(n)} - \mathbf{v}_b\|^2}{2\sigma^2} - \frac{(x_d^{(n)} - \mu_{d,b})^2}{2\sigma^2}\right), \\ &\quad d = 1, \dots, D_x, \quad b = 1, \dots, B. \end{aligned} \quad (14)$$

To perform feature selection in our additive CDE model, we introduce an  $(\ell_1, \ell_2)$ -mixed norm:

$$\Omega(\boldsymbol{\alpha}) := \sum_{d=1}^{D_x} \|\boldsymbol{\alpha}_d\|_2. \quad (15)$$

In the mixed norm  $\Omega(\boldsymbol{\alpha})$ , the parameter vector  $\boldsymbol{\alpha}$  is grouped in the sub-vectors  $\boldsymbol{\alpha}_1, \dots, \boldsymbol{\alpha}_{D_x}$ . Minimizing a loss function penalized with  $\Omega(\boldsymbol{\alpha})$  tends to produce a group-wise sparse solution (Yuan and Lin 2006), which means that the penalized optimization can be useful for selecting a relevant subset of input variables:

$$\min_{\boldsymbol{\alpha} \geq 0} J(\boldsymbol{\alpha}) := \hat{J}_0(\boldsymbol{\alpha}) + \lambda \Omega(\boldsymbol{\alpha}), \quad (16)$$

where  $\lambda \geq 0$  is the regularization parameter. Note that empirical squared error  $\hat{J}_0(\boldsymbol{\alpha})$  is differentiable and convex because the Hessian matrix  $\mathbf{H}$  is positive definite.  $\Omega(\boldsymbol{\alpha})$  is also convex but non-differentiable. Overall, (16) is a convex optimization problem and we develop a fast optimization algorithm below.

### 2.3 Optimization algorithm

We use a *proximal method* (Sra et al. 2012; Beck and Teboulle 2009) to solve the optimization problem (16). More specifically, we consider a linear approximation to function  $\hat{J}_0$  at the current solution  $\boldsymbol{\alpha}^{(t)}$ , penalized by a proximal term to keep the update confined in the neighborhood:

$$\min_{\boldsymbol{\alpha} \geq 0} \left\{ \hat{J}_0(\boldsymbol{\alpha}^{(t)}) + \nabla \hat{J}_0(\boldsymbol{\alpha}^{(t)})^T (\boldsymbol{\alpha} - \boldsymbol{\alpha}^{(t)}) + \lambda \cdot \Omega(\boldsymbol{\alpha}) + \frac{L}{2} \|\boldsymbol{\alpha} - \boldsymbol{\alpha}^{(t)}\|_2^2 \right\}. \quad (17)$$

Here,  $L$  is the Lipchitz constant,<sup>1</sup> which is given by the maximum eigenvalue of  $\mathbf{H}$  in the current setup. We can describe our update rule analytically as follows (the detailed derivation is described in “Appendix 2”):

<sup>1</sup> The Lipchitz constant  $L$  for  $f(x)$  satisfies  $\|\nabla f(x) - \nabla f(y)\|_2 \leq L\|x - y\|_2$  for arbitrary  $x$  and  $y$ .

**Algorithm 1** SA-CDE

---

```

1: Initialization:  $t \leftarrow 0$ ,  $\alpha^{(0)} = \mathbf{0}$ , and  $\epsilon > 0$ .
2: Compute  $\mathbf{H}$  and  $\mathbf{h}$  by (11) and (13), respectively.
3: Set  $L$  to the maximum eigenvalue of  $\mathbf{H}$ .
4: repeat
5:    $\mathbf{u} \leftarrow \left[ \alpha^{(t)} - \frac{1}{L} (\mathbf{H} \alpha^{(t)} - \mathbf{h}) \right]_+$ 
6:    $\alpha_d^{(t+1)} \leftarrow \left[ 1 - \frac{\lambda}{L \cdot \|\mathbf{u}_d\|_2} \right]_+ \cdot \mathbf{u}_d$ , for  $d = 1$  to  $D_x$ .
7: until  $\|\alpha^{(t+1)} - \alpha^{(t)}\| < \epsilon$ 
8: Normalize the optimized conditional density by (20).

```

---

$$\alpha_d^{(t+1)} \leftarrow \left[ 1 - \frac{\lambda}{L \cdot \|\mathbf{u}_d\|_2} \right]_+ \cdot \mathbf{u}_d, \quad d = 1, \dots, D_x, \quad (18)$$

where  $\mathbf{u}_d$  is the sub-vector of  $\mathbf{u}$  associated with the  $d$ th feature:

$$\mathbf{u} := \left[ \alpha^{(t)} - \frac{1}{L} \nabla \hat{J}_0(\alpha^{(t)}) \right]_+. \quad (19)$$

The operator  $[\cdot]_+$  rounds up negative values to zero. This update rule follows the standard proximal method, meaning that it has the global convergence rate of  $O(\frac{1}{t})$ , where  $t$  is the number of update iterations.

## 2.4 Post processing

Because we did not explicitly include the normalization constraint, the optimized conditional density estimator  $\hat{r}(\mathbf{y}, \mathbf{x}) = \sum_{d=1}^{D_x} \hat{r}_d(\mathbf{y}, x_d)$  may not be integrated to one with respect to  $\mathbf{y}$ . Here, we renormalize the estimator after optimization as

$$\hat{p}(\mathbf{y}|\mathbf{x}) = \frac{\hat{r}(\mathbf{y}, \mathbf{x})}{\int \hat{r}(\mathbf{y}', \mathbf{x}) d\mathbf{y}'}. \quad (20)$$

The denominator can be analytically calculated as

$$\int \hat{r}(\mathbf{y}, \mathbf{x}) d\mathbf{y} = \left( \sqrt{2\pi} \sigma \right)^{D_y} \sum_{d=1}^{D_x} \alpha_d^T \boldsymbol{\varphi}_d(x_d, \boldsymbol{\mu}_d). \quad (21)$$

Algorithm 2.4 summarizes our algorithm, which we call *sparse additive CDE* (SA-CDE).

## 2.5 Cross-validation for model selection

Performance of SA-CDE depends on the choice of model parameters such as the Gaussian width  $\sigma$  and the regularization parameter  $\lambda$ . Cross-validation (CV) is available to systematically choose these model parameters. Throughout this paper, we use fivefold CV: we first divide the samples into five subsets, then learn the parameter using four subsets, and evaluate the test error using the held-out subset. This procedure is iterated five times with different training-test choice and the error is averaged.

We use the negative log-likelihood (NLL) as our metric for evaluating the test error:

$$\text{NLL} = -\frac{1}{|T|} \sum_{n \in T} \log \hat{p}(\mathbf{y}^{(n)} | \mathbf{x}^{(n)}), \quad (22)$$

where  $\mathcal{T}$  is the set of indices of test samples. The smaller the value of NLL is, the better the performance of the conditional density estimator is. Thus, we chose the model parameters that minimize the averaged NLL by fivefold CV.

### 3 Numerical experiments

In this section, we experimentally evaluate the performance of our proposed method, SA-CDE. Throughout the experiments, the number of basis functions is fixed to  $B = \min(100, N)$ . The model parameters  $\sigma$  and  $\lambda$  were chosen from the twenty values between  $10^{-2}$  and 2 at the equal interval in the logarithmic scale by fivefold CV. We use NLL (22) for the performance measure of conditional density estimation. NLL is computed from test samples, which are not used for learning parameters and hyper-parameters. All experiments were implemented by Matlab 2013b and an HP DL360p Gen8 E5 v2 server with two CPUs of Xeon E5-2650 v2 2.60 GHz (8 Core) and the main memory of 96 GB.

#### 3.1 Compared methods

We compare SA-CDE with the following methods:

- **Sparse additive feature selection LSCDE (SA-LSCDE):** SA-LSCDE is a variation of the proposed SA-CDE, which first runs SA-CDE for feature selection and then estimates the conditionally density by LS-CDE with only selected features.
- **$\epsilon$ -neighbor kernel density estimation (eKDE):** eKDE estimates a conditional density by standard kernel density estimation using neighborhood samples in the domain of input  $\mathbf{x}$ , denoted by  $\mathcal{I}_{\mathbf{x}, \epsilon} := \{\mathbf{x}^{(i)} : \|\mathbf{x}^{(i)} - \mathbf{x}\|_2^2 \leq \epsilon\}$  for threshold  $\epsilon$ . In the case of Gaussian kernels, eKDE is given as

$$\hat{p}(\mathbf{y}|\mathbf{x}) = \frac{1}{|\mathcal{I}_{\mathbf{x}, \epsilon}|} \sum_{i \in \mathcal{I}_{\mathbf{x}, \epsilon}} \mathcal{N}(\mathbf{y}, \mathbf{y}^{(i)}, \sigma^2 \mathbf{I}_{D_y}), \quad (23)$$

where  $\mathcal{N}(\mathbf{y}, \boldsymbol{\mu}, \boldsymbol{\Sigma})$  denotes the Gaussian density function with respect to  $\mathbf{y}$  with mean  $\boldsymbol{\mu}$  and covariance matrix  $\boldsymbol{\Sigma}$ , and  $\mathbf{I}_{D_x}$  is the identity matrix of size  $D_x$ . In experiments, threshold  $\epsilon$  and bandwidth  $\sigma$  were chosen based on fivefold CV with respect to NLL, where the candidate values of  $\epsilon$  are the twenty values between  $10^{-2}$  and 5 at the equal interval in the logarithmic scale.

- **Least-squares conditional density estimation (LS-CDE):** The original LS-CDE method. This corresponds to a multi-dimensional non-sparse version of SA-CDE where, instead of the group-sparse penalty and an additive model, an  $\ell_2$ -penalty  $\lambda \|\boldsymbol{\alpha}\|_2^2$  and a multi-dimensional linear-in-parameter model,

$$\begin{aligned} \hat{p}(\mathbf{y}|\mathbf{x}) &:= \hat{r}(\mathbf{y}, \mathbf{x}) \\ &= \sum_{b=1}^B \alpha_b \left\{ \eta_b(\mathbf{y}) \cdot \varphi_b(\mathbf{x}) \right\}, \end{aligned} \quad (24)$$

is used. We use the Gaussian kernels for both  $\eta_b(\cdot)$  and  $\varphi_b(\cdot)$ , where the bandwidth  $\sigma$  and the regularization parameter  $\lambda$  are chosen based on fivefold CV with respect to NLL.



**Table 1** Computational complexities of our methods and existing CDEs

Method	SA-CDE	SA-LSCDE	LS-CDE	e-KDE	NW-CDE
Time	$O(N^3 D_x^3)$	$O(N^3 D_x^3)$	$O(N^3)$	$O(N^2 D_x)$	$O(N^2 D_x)$
Space	$O(N^2 D_x^2)$	$O(N^2 D_x^2)$	$O(N^2)$	$O(N^2)$	$O(N^2)$

**Table 2** Computational complexities of our methods and existing CDEs with forward feature selection

Method	SA-CDE	SA-LSCDE	FW-LSCDE	FW-eKDE	FW-NWCDE
Time	$O(N^3 D_x^3)$	$O(N^3 D_x^3)$	$O(D_x! N^3)$	$O(D_x! N^2)$	$O(D_x! N^2)$
Space	$O(N^2 D_x^2)$	$O(N^2 D_x^2)$	$O(N^2)$	$O(N^2)$	$O(N^2)$

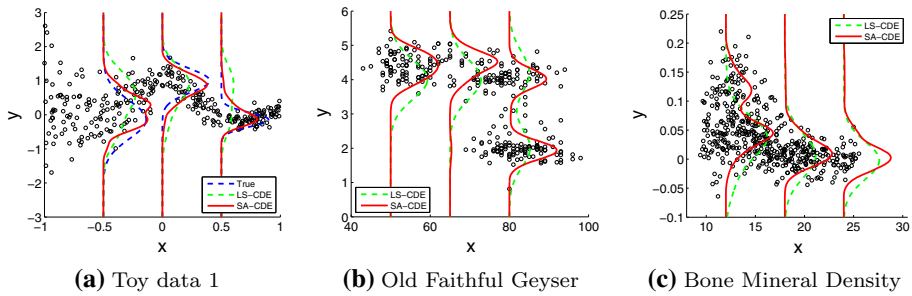
- **Nadaraya-Watson CDE (NW-CDE):** This corresponds to a simple version of LS-CDE, which fixes weights of basis functions to  $\frac{1}{B}$ :

$$\begin{aligned}\hat{p}(y|\mathbf{x}) &:= \frac{\hat{p}(y, \mathbf{x})}{\hat{p}(\mathbf{x})} \\ &= \frac{\sum_{b=1}^B \eta_b(y) \cdot \varphi_b(\mathbf{x})}{\sum_{b=1}^B \varphi_b(\mathbf{x})}.\end{aligned}\quad (25)$$

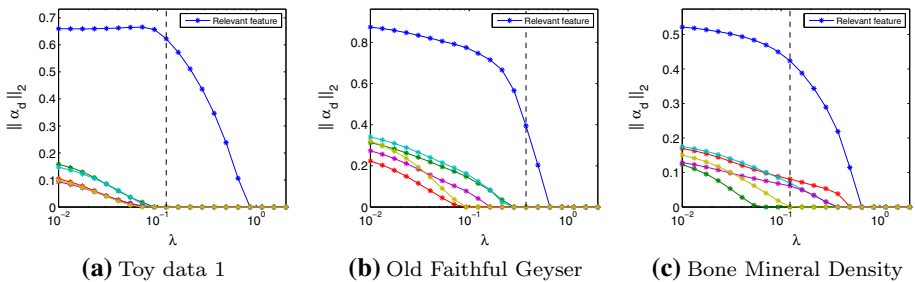
We use the Gaussian kernels for both  $\eta_b(\cdot)$  and  $\varphi_b(\cdot)$ , where the bandwidth  $\sigma$  is chosen based on leave-one-out CV for the exact likelihood formulated in Holmes et al. (2007). To directly employ the method in Holmes et al. (2007), we only use  $B$  samples in this CV procedure.

- **Forward feature selection + eKDE (FW-eKDE):** Forward feature selection is performed based on fivefold CV with respect to NLL. That is, the most useful feature that maximally reduces the cross-validated NLL by eKDE is selected one by one until the cross-validated NLL no longer decreases.
- **Forward feature selection + LS-CDE (FW-LSCDE):** Similarly, forward feature selection is performed for LS-CDE.
- **Forward feature selection + NW-CDE (FW-NWCDE):** Similarly, forward feature selection is performed for NW-CDE.

Tables 1 and 2 summarize the time and space complexities of our proposed method (SA-CDE and SA-LSCDE) and compared methods for a fixed hyper-parameter, under the assumption of  $N \gg D_x \gg D_y$  and  $O(N) = O(B)$ . All time complexities for  $T$  candidates of hyper-parameters are  $T$  times larger than those in Tables 1 and 2, while all space complexities does not depend on  $T$ . The time complexities of e-KDE and NW-CDE are the smallest, while CDEs with forward feature selections require larger computational costs. Especially, FW-LSCDE is worst in terms of the computational complexity. Thanks to the single-shot procedure of feature selection, SA-CDE and SA-LSCDE are computationally much more efficient than FW-LSCDE. The space complexities of all methods except for SA-CDE and SA-LSCDE are  $O(N^2)$ , while those of SA-CDE and SA-LSCDE are  $O(N^2 D_x^2)$ , which increases with  $D_x$ . Overall, the time complexities of SA-CDE and SA-LSCDE are much smaller than CDEs with forward selection procedures in return for increasing the space complexity.



**Fig. 3** Training samples and estimated conditional densities



**Fig. 4** Regularized path of SA-CDE. The blue line denotes the path of the relevant feature,  $x_1$ , while lines with other colors denote paths of irrelevant features  $x_2, \dots, x_6$ . The vertical black dashed line indicates the value of  $\lambda$  selected by fivefold CV

### 3.2 Illustrative examples

We first illustrate the behavior of our proposed method, SA-CDE, using toy and benchmark datasets having one relevant feature  $x_1$  and five irrelevant features  $x_2, \dots, x_6$ ,

- **Toy data 1:**  $x_1$  is independently generated following the uniform distribution on  $[-1, 1]$ , while each of  $x_2, \dots, x_6$  is generated by  $x_1 + \epsilon_c$  where  $\epsilon_c$  is a noise variable following the normal distribution with mean 0 and standard deviation  $3\hat{\sigma}$ , and  $\hat{\sigma}$  is the standard deviation of  $x_1$ . Output  $y$  is generated as a function of  $x_1$  as

$$y|x_1 \sim \text{sinc}\left(\frac{3}{4}\pi x_1\right) + \frac{1}{8} \exp(1 - x_1) \cdot \varepsilon, \quad (26)$$

where  $\varepsilon$  is standard normal noise. We generate  $N = 300$  samples for estimating the conditional density.

- **Old Faithful Geyser:** A benchmark dataset with  $D_x = D_y = 1$  that consists of durations of  $N = 299$  eruptions of the Old Faithful Geyser (Weisberg 1985). We add five irrelevant features  $x_2, \dots, x_6$  in a similar manner to Toy data 1.
- **Bone Mineral Density:** A benchmark dataset with  $D_x = D_y = 1$  that consists of relative spinal bone mineral density measurements on  $N = 485$  North American adolescents (Hastie et al. 2001). We add five irrelevant features  $x_2, \dots, x_6$  in a similar manner to Toy data 1.

Figures 3 and 4 show estimation results for these three datasets. Figure 3 shows estimated conditional densities: the black circles denote training samples, and the red solid and green

dashed lines denote the estimates obtained by the proposed SA-CDE and the plain LS-CDE, respectively. Figure 3a also contains the true conditional density drawn by the blue dashed line. Figure 4 shows the regularization paths of SA-CDE, i.e. the magnitude of each learned parameter  $\|\alpha_d\|_2$  as a function of the regularization parameter (Hastie et al. 2004). The blue line denotes the path of the relevant feature,  $x_1$ , while lines with other colors denote the paths of irrelevant features  $x_2, \dots, x_6$ . The vertical black dashed line indicates the value of  $\lambda$  selected by fivefold CV.

The regularization paths in Fig. 4 show that, in (a) Toy data 1 and (b) Old Faithful Geyser, the parameters corresponding to the irrelevant feature are zero and that corresponding to the relevant feature is non-zero for the cross-validated solution, which means that SA-CDE optimally performs feature selection. In Fig. 4c Bone Mineral Density, some of irrelevant features are non-zero because features with the skewed distribution are strongly correlated with the relevant feature despite additive Gaussian noise. Thus these features may still contain some information on the output value.

The estimation results in Fig. 3a show that, SA-CDE gives more accurate estimates than the plain LS-CDE. In Fig. 3b, c, SA-CDE tends to give sharper conditional density estimates than the plain LS-CDE. This is because relatively large Gaussian kernel widths are chosen in LS-CDE to incorporate irrelevant noisy features. This indicates that LS-CDE with many irrelevant features tend to produce too flat conditional densities which are not informative, while SA-CDE can avoid this problem by automatically eliminating irrelevant features.

### 3.3 Comparison of performance and computation time for different numbers of samples

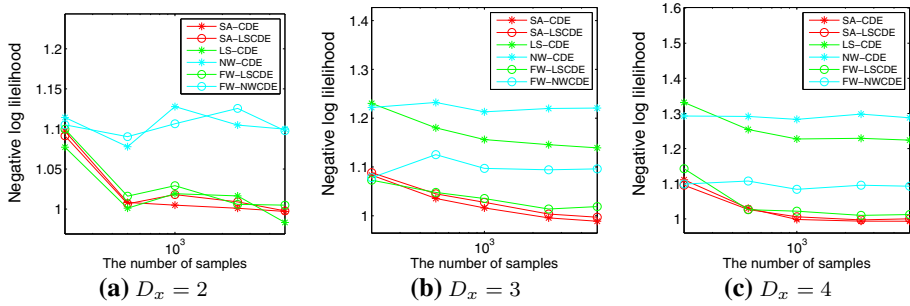
We compare the estimation performance and computation time of our proposed methods with existing CDEs for different numbers of training samples,  $N = 200, 500, 1,000, 2,500, 5,000$ . Performance evaluation measures are NLL and MSE (Mean Squared Error) in (7), both of which are computed from  $10^4$  test samples. The integral in MSE (7) is computed based on the test sample average. Datasets in this experiments are

- **Toy data 1:** The generation procedure is the same as the one in the previous section, in which both dimensions of relevant feature (input) and output are one.
- **Toy data 2:** Each irrelevant  $x_1, x_2, x_4, x_5, \dots, x_{D_x-2}, x_{D_x-1}$  is independently generated following the uniform distribution on  $[-1, 1]$ . Relevant features are generated by  $x_{3d} = x_{3d-2} + x_{3d-1}$ , and the  $d$ th dimension of output  $y$  is generated as a function of  $x_{3d}$ ,  $d = 1, 2, \dots, D_x/3$  as

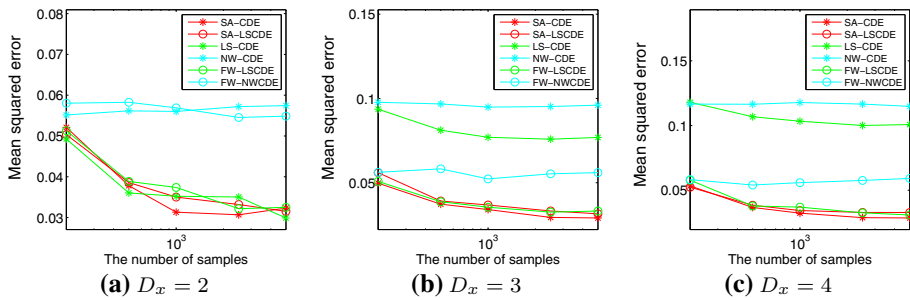
$$y_d | x_{3d} \sim \text{sinc}\left(\frac{3}{4}\pi x_{3d}\right) + \frac{1}{8} \exp(1 - x_{3d}) \cdot \varepsilon, \quad (27)$$

where  $\varepsilon$  is standard normal noise. This dataset has multi-dimensional relevant features and outputs.

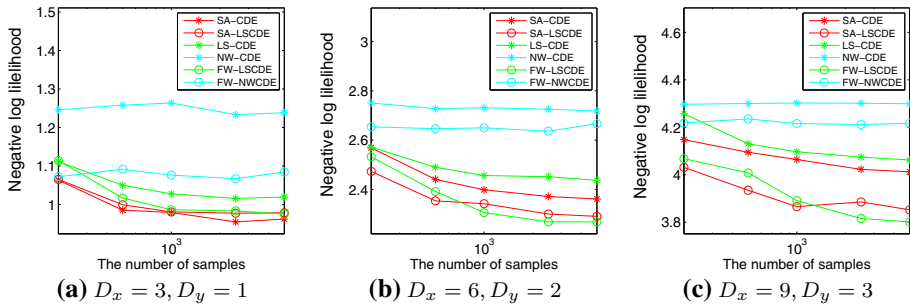
Figures 5 and 6 show NLL and MSE on Toy data 1, and Figs. 7 and 8 NLL and MSE on Toy data 2. These figures show that CDEs with feature selection except for FW-NWCDE, i.e. SA-CDE, SA-LSCDE, and FW-LSCDE, decrease both NLL and MSE with increasing the number of samples. However FW-NWCDE cannot decrease estimation errors even when the number of samples is increased. This is because NW-CDE does not optimize weight values of basis functions, resulting in poor performance. For Toy data 1 which has only single relevant input and output, the performance of SA-CDE and SA-LSCDE is almost the same. However, for Toy data 2 which has multiple relevant inputs and outputs, SA-LSCDE is much better than SA-CDE. This is because SA-LSCDE can represent more complex conditional densities than SA-CDE which is limited to additive models.



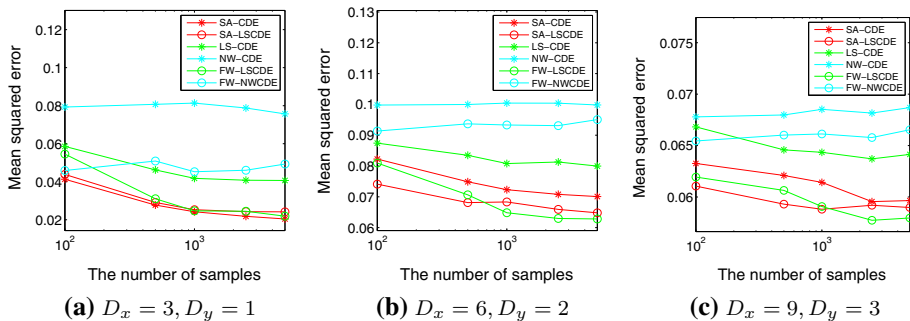
**Fig. 5** Negative log likelihood on Toy data 1 ( $D_y = 1$ )



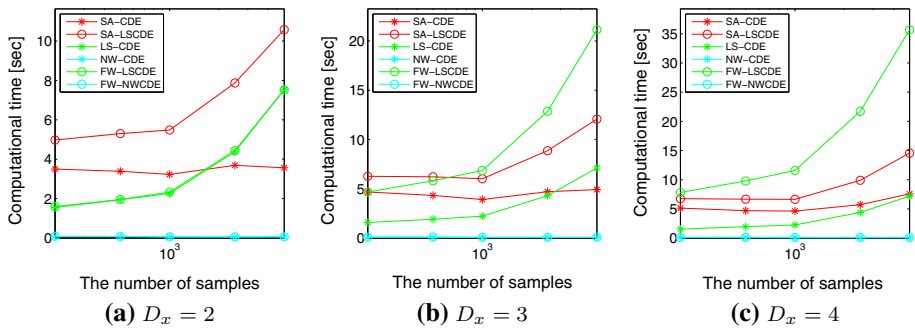
**Fig. 6** Mean squared error on Toy data 1 ( $D_y = 1$ )



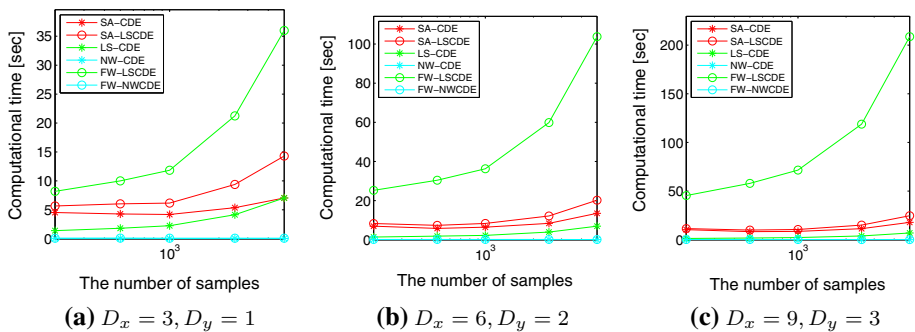
**Fig. 7** Negative log likelihood on Toy data 2



**Fig. 8** Mean squared error on Toy data 2



**Fig. 9** Computational time on Toy data 1 ( $D_y = 1$ )

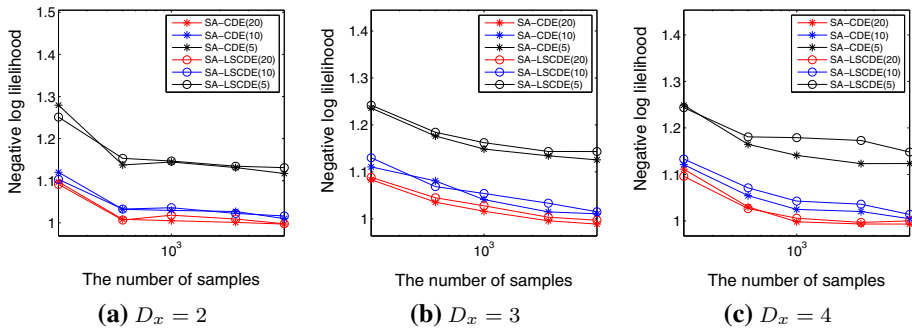


**Fig. 10** Computational time on Toy data 2

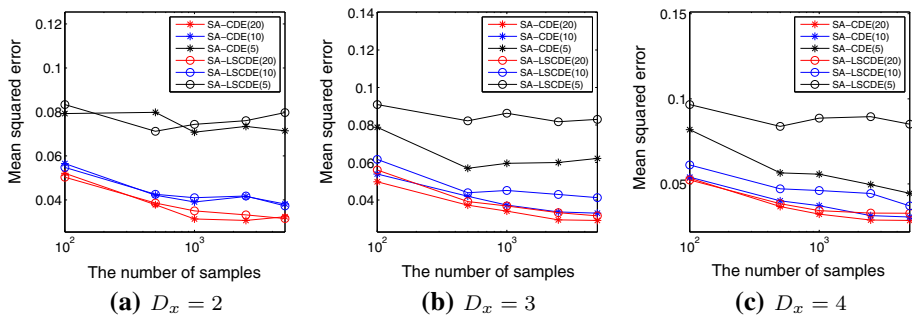
The computation time on Toy data 1 and Toy data 2, which is actual run-time including both fivefold CV to optimize hyper-parameters and conditional density estimation for test data, are shown in Figs. 9 and 10, respectively. These figures show that the computation time of FW-LSCDE is much longer than others when the number of input features is larger than 2 because of the time consuming procedure of forward feature selection. This weakness becomes more critical when the number of training samples is increasing. These results show that our proposed methods, SA-CDE and SA-LSCDE, are much faster than the existing feature selection methods, and our methods tend to improve the performance by feature selection.

### 3.4 Hyper-parameter selection

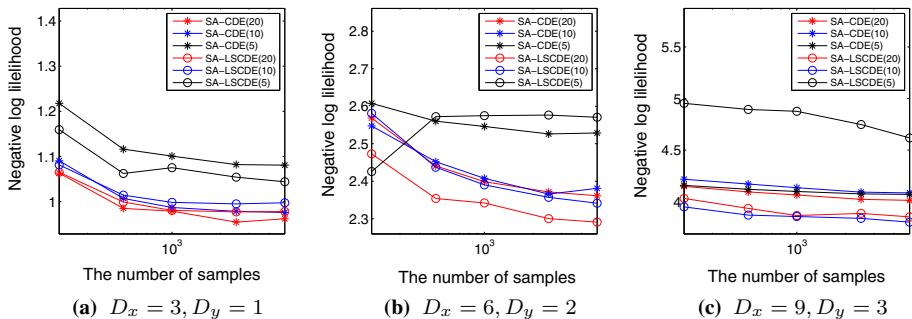
Our proposed methods (SA-CDE and SA-LSCDE) need to optimize hyper-parameters based on CV, which is a time consuming procedure. Thus it is a bottleneck of their computational time. Here, we compare the performance and computation time when the number of hyper-parameter candidates is changed. We keep choosing the hyper-parameters  $\sigma$  and  $\lambda$  between  $10^{-2}$  and 2 at the equal interval in the logarithmic scale, and we only change the number of hyper-parameter candidates:  $n_g = 5, 10, 20$ .  $n_g = 20$  is the default setting for all other experiments. Datasets and evaluation measures we use in this experiments are the same as Sect. 3.3. These experimental results are shown in Figs. 11, 12, 13, 14, 15 and 16. All performance results in Figs. 11, 12, 13 and 14 show that SA-CDE and SA-LSCDE with  $n_g = 20$  are the best. On the other hand, those with  $n_g = 5$ , which cannot improve



**Fig. 11** Negative log likelihood on Toy data 1 ( $D_y = 1$ )

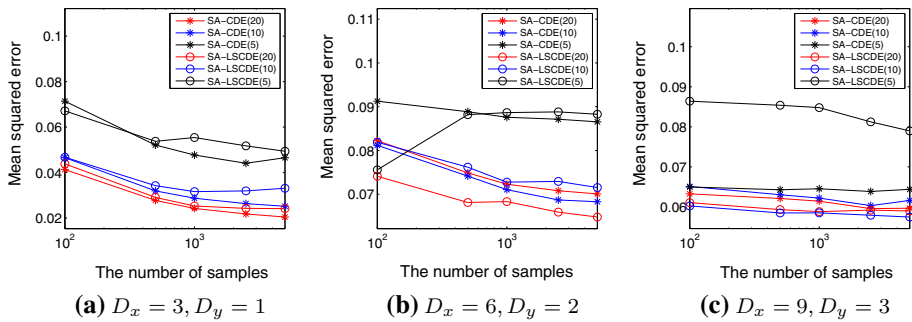


**Fig. 12** Mean squared error on Toy data 1 ( $D_y = 1$ )

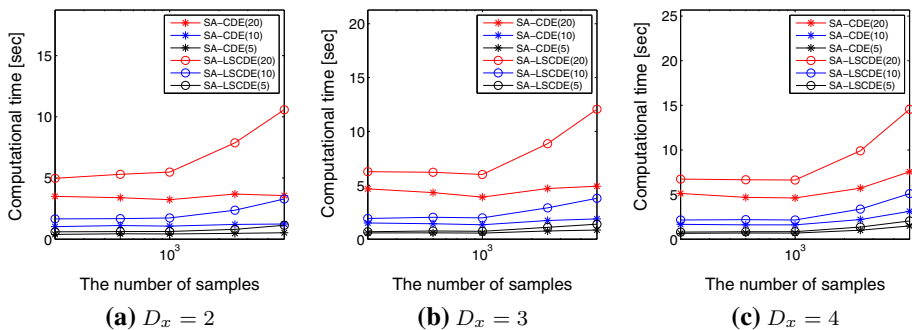


**Fig. 13** Negative log likelihood on Toy data 2

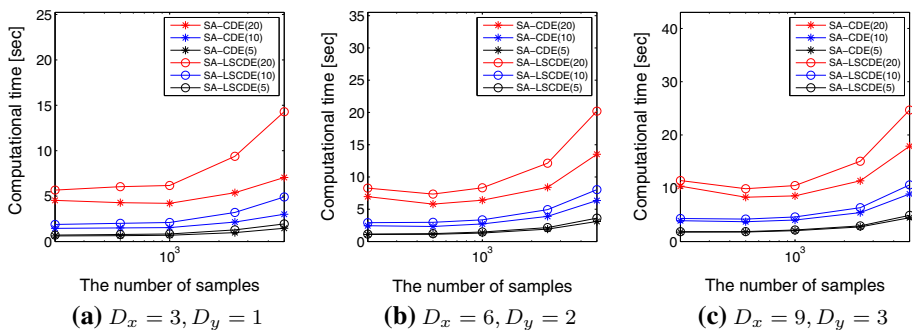
the performance even when the number of samples increases, are the worst. On the other hand, in terms of the computation time, Figs. 15 and 16 show that the  $n_g = 5$  is the best, while  $n_g = 20$  is the worst. This demonstrates the tradeoff between the performance and computation time on optimizing hyper-parameters using CV. We note that the computation time of SA-CDE and SA-LSCDE with  $n_g = 20$  is the worst, but they are still much faster than FW-LSCDE. Thus we keep using  $n_g = 20$  in all numerical experiments of later sections.



**Fig. 14** Mean squared error on Toy data 2



**Fig. 15** Computational time on Toy data 1 ( $D_y = 1$ )

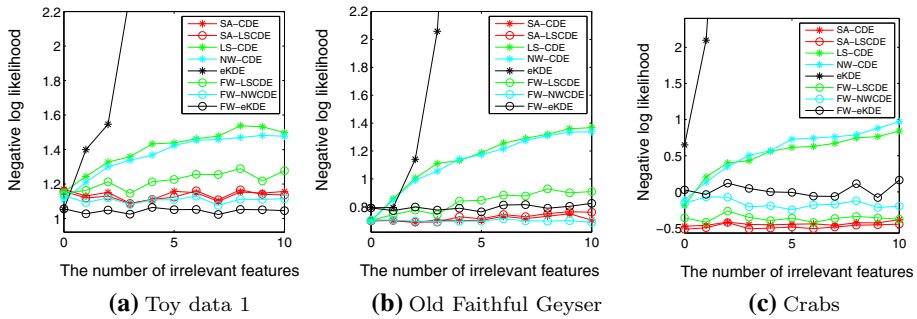


**Fig. 16** Computational time on Toy data 2

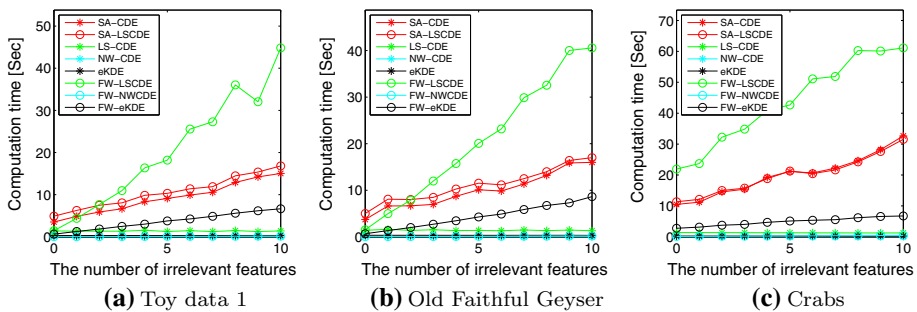
### 3.5 Performance comparison for different numbers of irrelevant features

Next, we compare the performance of SA-CDE with LS-CDE, eKDE, FW-LSCDE, and FW-eKDE for different numbers of irrelevant features. We use three datasets: (a) the same Toy data 1 ( $N = 300$  and  $D_x = 1$ ) used in the previous experiments, (b) the same *Old Faithful Geyser* benchmark dataset ( $N = 299$  and  $D_x = 1$ ), and (c) the *crabs* benchmark dataset ( $N = 200$  and  $D_x = 6$ ) taken from the R package.<sup>2</sup>

<sup>2</sup> <http://www.r-project.org/>.



**Fig. 17** Negative log likelihood with increasing noise dimensions



**Fig. 18** Computation time

For each dataset, we add  $m$  ( $= 0, 1, \dots, 10$ ) irrelevant features by copying  $x_1$  and adding Gaussian noise in a similar manner to the previous experiments. We randomly choose a half of samples as training samples to estimate conditional densities and use the rest as test samples to compute the test NLL. This procedure is repeated 100 times and the averaged test NLL is computed. The experimental results are summarized in Fig. 17.

In all three cases, the NLL values of LS-CDE, NW-CDE, and eKDE (without feature selection) grow as the number of irrelevant features increases. On the other hand, the NLL values of SA-CDE, SA-LSCDE, FW-LSCDE, FW-NWCDE, and FW-eKDE do not grow that much when the number of irrelevant features increases. This clearly demonstrates an advantage of performing feature selection.

The total computation time of each method, i.e. run-time including fivefold CV to optimize hyper-parameters and conditional density estimation for test data, is plotted in Fig. 18. This shows that all methods require more computation costs as the number of features increases. The computation time of FW-eKDE, FW-LSCDE, and FW-NWCDE is much longer than the plain eKDE, LS-CDE, and NW-CDE, implying that forward feature selection is highly time-consuming. Indeed, forward feature selection involves repetitious conditional density estimation to find the best feature to add, which is computationally highly demanding. On the other hand, the computation time of SA-CDE and SA-LSCDE grows less sharply than FW-LSCDE, thanks to the single-shot procedure of feature selection and conditional density estimation.



**Table 3** NLL for benchmark datasets with five dimensional irrelevant features

Name	$ \mathcal{F} $	$N$	SA-CDE	SA-LSCDE	LS-CDE	eKDE	NW-CDE
caution	2	100	$1.34 \pm 0.6$	<b><math>1.24 \pm 0.4</math></b>	$1.38 \pm 0.3$	$24.25 \pm 3.4$	$1.36 \pm 0.2$
CobarOre	2	38	$1.71 \pm 0.5$	$1.70 \pm 0.4$	<b><math>1.62 \pm 0.2</math></b>	$31.81 \pm 2.9$	<b><math>1.62 \pm 0.2</math></b>
snowgeese	2	45	<b><math>1.80 \pm 2.0</math></b>	<b><math>1.76 \pm 1.8</math></b>	$1.85 \pm 1.3$	$22.04 \pm 6.1$	<b><math>1.59 \pm 1.0</math></b>
topo	2	52	$1.17 \pm 0.3$	$1.14 \pm 0.3$	$1.22 \pm 0.2$	$29.30 \pm 3.0$	$1.21 \pm 0.1$
sniffer	4	125	$0.70 \pm 0.6$	<b><math>0.60 \pm 0.7</math></b>	$0.85 \pm 0.2$	$16.91 \pm 3.2$	$0.83 \pm 0.2$
crabs	6	200	$-0.44 \pm 0.1$	<b><math>-0.47 \pm 0.3</math></b>	$0.53 \pm 0.1$	$26.03 \pm 3.1$	$0.58 \pm 0.1$
UN3	6	125	<b><math>1.27 \pm 0.2</math></b>	<b><math>1.35 \pm 0.4</math></b>	$1.57 \pm 0.6$	$33.36 \pm 1.6$	$1.54 \pm 0.6$
birthwt	7	189	<b><math>1.49 \pm 0.2</math></b>	<b><math>1.52 \pm 0.1</math></b>	<b><math>1.51 \pm 0.1</math></b>	$31.77 \pm 1.6$	$1.67 \pm 0.2$
cpus	7	209	<b><math>0.36 \pm 0.6</math></b>	$0.80 \pm 0.7$	$1.19 \pm 0.5$	$22.29 \pm 3.4$	$1.17 \pm 0.6$
gilgais	8	365	<b><math>0.70 \pm 0.2</math></b>	$0.89 \pm 0.2$	$1.16 \pm 0.2$	$27.77 \pm 2.2$	$1.11 \pm 0.2$
BigMac	9	69	$1.33 \pm 0.8$	$1.37 \pm 0.7$	$1.42 \pm 0.7$	$35.79 \pm 0.5$	$1.34 \pm 0.5$
highway	11	39	<b><math>1.38 \pm 0.7</math></b>	$1.60 \pm 0.7$	$1.71 \pm 0.8$	$36.04 \pm 0.0$	$1.74 \pm 0.7$
Time			1.00	0.00	0.06	0.02	0.00
Name	FW-LSCDE	FW-eKDE	FW-NWCDE				
caution	$1.33 \pm 0.6$	$1.35 \pm 0.6$	$1.30 \pm 0.5$				
CobarOre	$1.95 \pm 0.6$	$2.45 \pm 1.9$	<b><math>1.65 \pm 0.4</math></b>				
snowgeese	$2.09 \pm 1.9$	$3.03 \pm 2.4$	<b><math>1.82 \pm 1.8</math></b>				
topo	$1.19 \pm 0.4$	$1.73 \pm 1.2$	<b><math>1.07 \pm 0.2</math></b>				
sniffer	$0.74 \pm 0.8$	$0.96 \pm 0.8$	$0.96 \pm 1.1$				
crabs	$-0.37 \pm 0.3$	$0.08 \pm 0.6$	$-0.12 \pm 0.8$				
UN3	<b><math>1.27 \pm 0.3</math></b>	$1.60 \pm 0.6$	<b><math>1.34 \pm 0.3</math></b>				
birthwt	$1.67 \pm 0.2$	$1.75 \pm 0.5$	$3.85 \pm 2.1$				
cpus	$0.70 \pm 0.8$	$1.00 \pm 0.9$	$0.76 \pm 0.9$				
gilgais	$0.76 \pm 0.2$	$0.97 \pm 0.3$	$1.20 \pm 0.3$				
BigMac	$1.45 \pm 0.9$	$2.54 \pm 1.7$	<b><math>1.23 \pm 0.8</math></b>				
highway	$2.06 \pm 1.0$	$3.17 \pm 1.9$	$2.18 \pm 1.8$				
Time	2.73	0.54	0.01				

### 3.6 Benchmark datasets

We further compare the performance of the proposed SA-CDE with other methods on twelve benchmark datasets accompanied with the R package. All of these datasets have one-dimensional output, i.e.,  $D_y = 1$ . The number of features  $|\mathcal{F}|$  and the number of samples  $N$  are listed in Table 3. For all datasets, five irrelevant features are added which are the copy of the relevant variable or a linear combination of two relevant variables contaminated with Gaussian noise. The type of noise and relevant variables are chosen at random. Gaussian noise is generated in a similar manner to the previous experiments. In this experiment, we randomly choose a half of samples as training data for estimating conditional densities, and the rest is used as test data for computing the test NLL. The experimental results are summarized in Table 3. The values described in the table are averaged NLL values and standard deviations over one hundred runs with different random seeds. The bold letter means the best NLL and comparable results that could not be rejected by the two-sided paired

$t$ -test at significance level 5 %. The bottom row shows the averaged normalized computation time.

Table 3 shows that the performance of our methods (SA-CDE and SA-LSCDE) is best on nine datasets. For high-dimensional datasets, especially when  $|\mathcal{F}|$  is seven or more, SA-CDE tends to outperform other methods with statistical significance. For low-dimensional datasets with large  $N$ , the performance of SA-LSCDE outperforms SA-CDE because of their expressive power of functions. For low-dimensional datasets with small  $N$ , FW-NWCDE performs the best because all other methods optimizing weights of basis functions cause overfitting. LSCDE, eKDE, NW-CDE, and FW-eKDE are computationally much more efficient than SA-CDE and SA-LSCDE, but these methods tend to perform poorly for high-dimensional relevant features with noisy dimensions.

### 3.7 Humanoid robot transition dataset

Finally, we evaluate the performance of the proposed method on humanoid robot transition estimation with multiple inputs and multiple outputs. The dataset was generated from a simulator of the upper-body part of the humanoid robot *CB-i* (Cheng et al. 2007). The robot has 9 controllable joints: shoulder pitch, shoulder roll and elbow pitch of the right arm, shoulder pitch, shoulder roll and elbow pitch of the left arm, waist yaw, torso roll, and torso pitch joints.

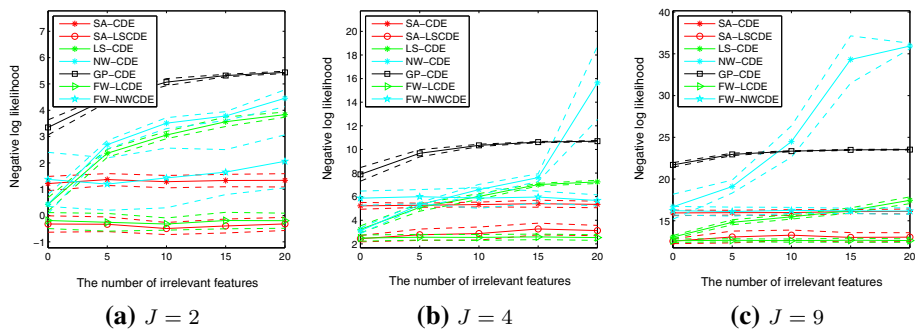
Posture of the robot is described by 18-dimensional real-valued state vector  $s$ , which corresponds to the angle and angular velocity of each joint in radians and radians per seconds, respectively. We can control the robot by sending the action command  $a$  to the system. The action command  $a$  is a 9-dimensional real-valued vector, which corresponds to the target angle of each joint. When the robot is currently at state  $s$  and receive action  $a$ , the physical control system of the simulator calculates the amount of torques to be applied to each joint. These torques are calculated by the *Proportional-Derivative* (PD) controller as

$$\tau_i = K_{p_i}(a_i - s_i) - K_{d_i}\dot{s}_i, \quad (28)$$

where  $s_i$ ,  $\dot{s}_i$ , and  $a_i$  denote the current angle, the current angular velocity, and the received target angle of the  $i$ th joint, respectively.  $K_{p_i}$  and  $K_{d_i}$  denote the position and velocity gains for the  $i$ th joint, respectively. We set  $K_{p_i} = 200$  and  $K_{d_i} = 10$  for the elbow pitch joints and  $K_{p_i} = 2000$  and  $K_{d_i} = 100$  for the other joints. After the torques are applied to the joints, the physical system updates the state of the robot to  $s'$ . We simulate a noisy control system by perturbing action vectors with independent bi-modal Gaussian noise. More specifically, for each action element, we add Gaussian noise with mean 0 and standard deviation 0.052 with probability 0.6, and Gaussian noise with mean  $-0.087$  and standard deviation 0.052 with probability 0.4.

To generate transition samples, we first generated the initial posture of the robot  $s^{(1)}$  at random and then simulated a trajectory with 100 steps, i.e.  $s^{(2)}, \dots, s^{(100)}$ . For each step, we additionally generated  $m$  irrelevant input features  $z^{(n)} \in \mathbb{R}^m$  by copying a relevant variable or by linearly combining two relevant variables contaminated with Gaussian noise in a similar manner to the previous experiments. By iterating these procedures, we obtained the transition samples  $\{(s^{(n)}, a^{(n)}, z^{(n)}, s'^{(n)})\}_{n=1}^{10000}$ .

Our goal is to learn the system dynamics as state transition probability  $p(s'|s, a, z)$  from these samples. Thus, as the conditional density estimation problem, the state-action pair  $(s^T, a^T, z^T)^T$  is regarded as input variable  $x$ , while the next state  $s'$  is regarded as output variable  $y$ . Note that an accurate estimate of the state transition probability is highly useful in *model-based reinforcement learning* (Sutton and Barto 1998).



**Fig. 19** Negative log likelihood on humanoid robot data

From the transition samples, we randomly picked up 5000 samples as training data and used the other 5000 samples as test data to calculate NLL. We compare our proposed method SA-CDE with LS-CDE, NW-CDE, FW-LSCDE, and FW-NWCDE, as well as parametric conditional density estimation by the *Gaussian process regression* (GP-CDE) (Rasmussen and Williams 2005). In this experiment, the candidate values of regularization parameter  $\lambda$  are the twenty values between  $10^{-3}$  and  $10^{-1}$  at the equal interval in the logarithmic scale, while the candidate values of other parameters are the same as the previous setting. We consider three datasets with  $J = 2, 4, 9$  joints, and change the number of irrelevant features as  $m = 0, 5, 10, 15, 20$ . Thus, the input dimensionality is  $3J + m$ , while the output dimensionality is  $2J$ . For each  $J$  and  $m$ , we evaluated the performance of conditional density estimation methods by averaged NLL and averaged computational time over 20 runs.

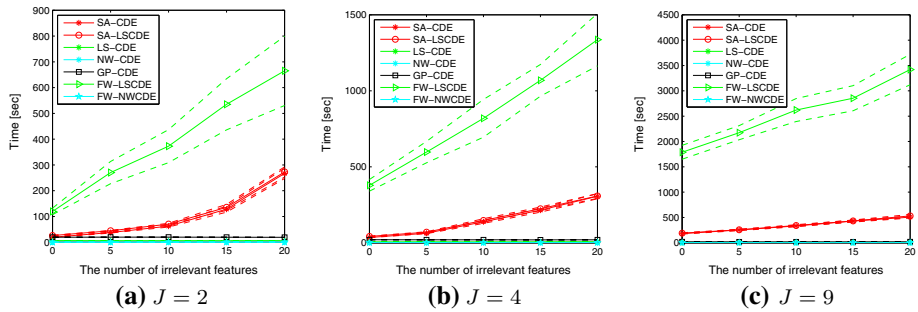
Figure 19 shows the experimental results, where the solid lines denote the averaged values and the dashed lines denote the averaged values with one standard deviation. From the plots, we can confirm that the NLL values of LS-CDE, NW-CDE, and GP-CDE grow sharply as the number of irrelevant features  $m$  is increased. On the other hand, the NLL values of SA-CDE, FW-LSCDE, FW-NWCDE, and SA-LSCDE do not increase even if the number of irrelevant features is increased. Among them, SA-CDE cannot outperform FW-LSCDE and SA-LSCDE, because the additive-model assumption of SA-CDE caused large estimation bias. However, feature selection by SA-CDE itself performs well and SA-LSCDE performs comparably to FW-LSCDE.

Figure 20 plots the computation time. LS-CDE, NW-CDE, and GP-CDE are very fast because no feature selection process is involved. NW-CDE and FW-NWCDE are also very fast because no optimization of weight parameters is involved. Among SA-CDE, FW-LSCDE, and SA-LSCDE, SA-CDE and SA-LSCDE are much faster than FW-LSCDE.

Overall, in this challenging task of robot transition estimation, SA-LSCDE, the combination of SA-CDE and LS-CDE, was shown to be the most promising approach.

## 4 Conclusions

We proposed a direct estimator of conditional probability densities that is equipped with feature selection. Our feature selection strategy is based on the  $\ell_1/\ell_2$  mixed-norm, which tends to produce a group-sparse solution. An optimization algorithm based on a proximal method was presented that is guaranteed to possess fast convergence. The numerical experi-



**Fig. 20** Computation time on humanoid robot data

ments on benchmark and robot transition datasets demonstrated that the proposed method is promising.

SA-CDE assumes the additive structure for feature selection. However, this causes linear increase of the time and space complexities, resulting in high computation costs for datasets with a large number of features. Improving the scalability issue is future work.

**Acknowledgments** Motoki Shiga was supported by JSPS KAKENHI 25870322. Masashi Sugiyama was supported by JSPS KAKENHI 23120004 and AOARD. Authors thank Dr. Ichiro Takeuchi, Nagoya Institute of Technology, for kindly providing his source codes.

## Appendix 1: Derivation of Eq. (10)

Substituting (4) into (9), we can transform the term associated with  $d_1$ ,  $d_2$ , and  $n$  in the above equation as

$$\begin{aligned}
 & \int \hat{r}_{d_1}(\mathbf{y}, x_{d_1}^{(n)}) \hat{r}_{d_2}(\mathbf{y}, x_{d_2}^{(n)}) d\mathbf{y} \\
 &= \int \prod_{d=d_1, d_2} \left\{ \sum_{b=1}^B \alpha_{d,b} \cdot \eta(\mathbf{y}, \mathbf{v}_b) \cdot \varphi_d(x_d^{(n)}, \mu_{d,b}) \right\} d\mathbf{y} \\
 &= \sum_{b_1=1}^B \sum_{b_2=1}^B \left\{ \int \eta(\mathbf{y}, \mathbf{v}_{b_1}) \cdot \eta(\mathbf{y}, \mathbf{v}_{b_2}) d\mathbf{y} \times \varphi_{d_1}(x_{d_1}^{(n)}, \mu_{d_1,b_1}) \cdot \varphi_{d_2}(x_{d_2}^{(n)}, \mu_{d_2,b_2}) \right\}. \quad (29)
 \end{aligned}$$

The integral with respect to  $\mathbf{y}$  in the above equation can be analytically computed as

$$\begin{aligned}
 \int \eta(\mathbf{y}, \mathbf{v}_{b_1}) \cdot \eta(\mathbf{y}, \mathbf{v}_{b_2}) d\mathbf{y} &= \int \exp\left(-\frac{\|\mathbf{y} - \mathbf{v}_{b_1}\|^2 + \|\mathbf{y} - \mathbf{v}_{b_2}\|^2}{2\sigma^2}\right) d\mathbf{y} \\
 &= \left(\sqrt{\pi}\sigma\right)^{D_y} \exp\left(-\frac{\|\mathbf{v}_{b_1} - \mathbf{v}_{b_2}\|_2^2}{4\sigma^2}\right). \quad (30)
 \end{aligned}$$

Using (5), (6) and (30), we can derive (10).

## Appendix 2: Derivation of update rule (18)

By some simple algebra, the optimization problem (17) can be transformed into

$$\min_{\alpha \geq 0} \left\{ \frac{1}{2} \|\alpha - \mathbf{u}\|_2^2 + \frac{\lambda}{L} \Omega(\alpha) \right\} \quad (31)$$

where  $\mathbf{u}$  is defined as (19). Because the optimization problem (31) is strongly convex, update of  $\alpha$  by (31) can be regarded as an operator that maps  $\mathbf{u}$  to a unique value. This is called a proximal operator associated with  $\frac{\lambda}{L} \Omega$ . Denoting the proximal operator by  $\text{Prox}_{\frac{\lambda}{L} \Omega}(\cdot)$ , we can write our update rule as

$$\begin{aligned} \alpha^{(t+1)} &\leftarrow \text{Prox}_{\frac{\lambda}{L} \Omega}(\mathbf{u}) \\ &:= \arg \min_{\alpha \geq 0} \left\{ \frac{1}{2} \|\alpha - \mathbf{u}\|_2^2 + \frac{\lambda}{L} \Omega(\alpha) \right\}. \end{aligned} \quad (32)$$

Next we describe the solution of this proximal operation by transforming the operator into the dual form. By using the Fenchel conjugate of function  $\hat{J}_0$  and the dual norm of  $\Omega$ , the dual problem of (31) can be obtained as

$$\max_{\mathbf{v}} \left\{ -\sup_{\alpha} \left[ \alpha^T(-\mathbf{v}) - \frac{1}{2} \|\alpha - \mathbf{u}\|_2^2 \right] \right\} \text{ subject to } \|\mathbf{v}_d\|_2 \leq \frac{\lambda}{L}, \quad d = 1, \dots, D_x, \quad (33)$$

where  $\mathbf{v}$  is the dual of  $\alpha$  and  $\mathbf{v}_d$  is the sub-vector associated with the  $d$ th feature. The objective function can be transformed into

$$-\sup_{\alpha} \left[ \alpha^T(-\mathbf{v}) - \frac{1}{2} \|\alpha - \mathbf{u}\|_2^2 \right] = \inf_{\alpha} \left[ \frac{1}{2} \|\alpha - \mathbf{u}\|_2^2 + \alpha^T \mathbf{v} \right]. \quad (34)$$

Since the above function is quadratic with respect to  $\alpha$ , its minimum is achieved at the point where the derivative is zero: from this, we have  $\alpha = \mathbf{u} - \mathbf{v}$ . Substituting this back into (34), we have

$$\begin{aligned} \inf_{\alpha} \left[ \frac{1}{2} \|\alpha - \mathbf{u}\|_2^2 + \alpha^T \mathbf{v} \right] &= \frac{1}{2} \|\mathbf{u} - \mathbf{v} - \mathbf{u}\|_2^2 + (\mathbf{u} - \mathbf{v})^T \mathbf{v} \\ &= -\frac{1}{2} (\|\mathbf{v} - \mathbf{u}\|_2^2 - \|\mathbf{u}\|_2^2). \end{aligned} \quad (35)$$

Further substituting (35) into (33), we can express the dual of the proximal operator for optimizing dual vector  $\mathbf{v}$ , denoted by  $\text{Proj}_{\frac{\lambda}{L} \Omega}$ , as

$$\begin{aligned} \mathbf{v}^{(t+1)} &\leftarrow \text{Proj}_{\frac{\lambda}{L} \Omega}(\mathbf{u}) \\ &:= \arg \min_{\mathbf{v}} \left[ \frac{1}{2} \|\mathbf{v} - \mathbf{u}\|_2^2 \right] \text{ subject to } \|\mathbf{v}_d\|_2 \leq \frac{\lambda}{L}, \quad d = 1, \dots, D_x. \end{aligned} \quad (36)$$

Given that the solution for each sub-vector  $\mathbf{v}_d$  can be obtained separately, it is easy to confirm that the solution of the above proximal operation is given by

$$\left[ \text{Proj}_{\frac{\lambda}{L} \Omega}(\mathbf{u}) \right]_d = \left( 1 - \left[ 1 - \frac{\lambda}{L \|\mathbf{u}_d\|_2} \right]_+ \right) \mathbf{u}_d, \quad d = 1, \dots, D_x. \quad (37)$$

Since  $\alpha = \mathbf{u} - \mathbf{v}$ , we have

$$\text{Prox}_{\frac{\lambda}{L} \Omega}(\mathbf{u}) = \mathbf{u} - \text{Proj}_{\frac{\lambda}{L} \Omega}(\mathbf{u}). \quad (38)$$

Using the above equation, we can describe our update rule analytically as (18).

## References

- Beck, A., & Teboulle, M. (2009). A fast iterative shrinkage-thresholding algorithm for linear inverse problems. *SIAM Journal on Imaging Sciences*, 2(1), 183–202.
- Bishop, C. M. (2006). *Pattern recognition and machine learning*. New York, NY: Springer.
- Cheng, C., Hyon, S. H., Morimoto, J., Ude, A., Hale, J. G., Colvin, G., et al. (2007). Cb: A humanoid research platform for exploring neuroscience. *Advanced Robotics*, 21(10), 1097–1114.
- Cook, R. D., & Ni, L. (2005). Sufficient dimension reduction via inverse regression. *Journal of the American Statistical Association*, 100(470), 410–428.
- Dempster, A. P., Laird, N. M., & Rubin, D. B. (1977). Maximum likelihood from incomplete data via the EM algorithm. *Journal of the Royal Statistical Society, Series B*, 39(1), 1–38.
- Fan, J., Yao, Q., & Tong, H. (1996). Estimation of conditional densities and sensitivity measures in nonlinear dynamical systems. *Biometrika*, 83(1), 189–206.
- Guyon, I., & Elisseeff, A. (2003). An introduction to variable and feature selection. *Journal of Machine Learning Research*, 3(Mar.), 1157–1182.
- Hastie, T., Rosset, S., Tibshirani, R., & Zhu, J. (2004). The entire regularization path for the support vector machine. *Journal of Machine Learning Research*, 5, 1391–1415.
- Hastie, T., Tibshirani, R., & Friedman, J. (2001). *The elements of statistical learning: Data mining, inference, and prediction*. New York, NY: Springer.
- Holmes, M. P., Gray, A. G., & Isbell, C. L. (2007). Fast nonparametric conditional density estimation. In *Proceedings of the twenty-third conference annual conference on uncertainty in artificial intelligence* (pp 175–182).
- Li, K. (1991). Sliced inverse regression for dimension reduction. *Journal of the American Statistical Association*, 86(414), 316–342.
- Li, Y., Liu, Y., & Zhu, J. (2007). Quantile regression in reproducing kernel Hilbert spaces. *Journal of the American Statistical Association*, 102(477), 255–268.
- Rasmussen, C. E., & Williams, C. K. I. (2005). *Gaussian processes for machine learning (adaptive computation and machine learning)*. Cambridge: MIT Press.
- Silverman, B. W. (1986). *Density estimation for statistics and data analysis*. London: Chapman and Hall.
- Sra, S., Nowozin, S., & Wright, S. (2012). *Optimization for machine learning. neural information processing series*. Cambridge: MIT Press.
- Sugiyama, M., Takeuchi, I., Suzuki, T., Kanamori, T., Hachiya, H., & Okanohara, D. (2010). Least-squares conditional density estimation. *IEICE Transactions on Information and Systems*, E93-D(3), 583–594.
- Sutton, R. S., & Barto, A. G. (1998). *Introduction to reinforcement learning* (1st ed.). Cambridge, MA: MIT Press.
- Takeuchi, I., Le, Q. V., Sears, T. D., & Smola, A. J. (2006). Nonparametric quantile estimation. *Journal of Machine Learning Research*, 7, 1231–1264.
- Takeuchi, I., Nomura, K., & Kanamori, T. (2009). Nonparametric conditional density estimation using piecewise-linear solution path of kernel quantile regression. *Neural Computation*, 21(2), 533–559.
- Tresp, V. (2001). Mixtures of Gaussian processes. In T. K. Leen, T. G. Dietterich, & V. Tresp (Eds.), *Advances in Neural information processing systems, vol. 13* (pp. 654–660). Cambridge, MA: MIT Press.
- Weisberg, S. (1985). *Applied linear regression*. New York, NY: Wiley.
- Wolff, R. C. L., Yao, Q., & Hall, P. (1999). Methods for estimating a conditional distribution function. *Journal of the American Statistical Association*, 94(445), 154–163.
- Yuan, M., & Lin, Y. (2006). Model selection and estimation in regression with grouped variables. *Journal of the Royal Statistical Society, Series B*, 68(1), 49–67.