# Multi-Kernel Correntropy for Robust Learning

Badong Chen, *Senior Member, IEEE,* Yuqing Xie, Xin Wang, *Student Member, IEEE,*
Zejian yuan, *Member, IEEE,* Pengju Ren, *Member, IEEE,* and Jing Qin, *Member, IEEE*

*Abstract*—As a novel similarity measure that is defined as the expectation of a kernel function between two random variables, correntropy has been successfully applied in robust machine learning and signal processing to combat large outliers. The kernel function in correntropy is usually a zero-mean Gaussian kernel. In a recent work, the concept of mixture correntropy (MC) was proposed to improve the learning performance, where the kernel function is a mixture Gaussian kernel, namely a linear combination of several zero-mean Gaussian kernels with different widths. In both correntropy and mixture correntropy, the center of the kernel function is, however, always located at zero. In the present work, to further improve the learning performance, we propose the concept of multi-kernel correntropy (MKC), in which each component of the mixture Gaussian kernel can be centered at a different location. The properties of the MKC are investigated and an efficient approach is proposed to determine the free parameters in MKC. Experimental results show that the learning algorithms under the maximum multi-kernel correntropy criterion (MMKCC) can outperform those under the original maximum correntropy criterion (MCC) and the maximum mixture correntropy criterion (MMCC).

*Index Terms*—Correntropy, mixture correntropy, multi-kernel correntropy, robust learning, outliers.

## I. INTRODUCTION

A Key problem in supervised machine learning is how to define an objective function to measure the similarity between model output and a target variable. The mean square error (MSE) is one of the most popular similarity measures, which is computationally simple and easy to use as a performance index in many signal processing and machine learning applications. The MSE is, however, vulnerable to non-Gaussian noises, such as impulsive noises or outliers, because the solution that minimizes the squared difference (the error in $L_2$ norm) can deviate far from the optimal solution in the presence of large outliers. To address this problem, many non-MSE similarity measures were proposed in the literature, such as the mean absolute error (MAE)[1, 2], mean p-power error (MPE)[3], M-estimation cost [4] and logarithmic cost [5]. In particular in recent years, the correntropy as a local similarity measure in kernel space has found many successful applications in robust regression [6, 7], classification [8–12], PCA [13], feature extraction [14, 15], adaptive filtering [16–22] and so on. Correntropy defines a non-homogeneous metric

(Correntropy Induced Metric, CIM) that behaves like different norms (from $L_2$ to $L_0$) depending on the actual distance between samples, which can be used as an outlier-robust error measure in robust signal processing or a sparsity penalty term in sparse signal processing [23].

The original correntropy is defined as the expectation of a kernel function between two random variables, where the kernel function is usually a zero-mean Gaussian kernel [23]. The learning methods under maximum correntropy criterion (MCC) may, however, perform poorly when the kernel function in correntropy is limited to a single Gaussian kernel. It is likely that the combination of several kernel functions can perform much better. The mixture correntropy (MC) was thus proposed in a recent work to improve the learning performance, in which the kernel function is implemented by a linear combination of several zero-mean Gaussian kernels with different widths [24]. Similar ideas can be found in multiple kernel learning (MKL) methods [25], such as the Multiple Kernel Support Vector Machine (MKSVM) [26], Multiple Kernel Modification of Ho-Kashyap algorithm with Squared approximation of the misclassification errors (MultiK-MHKS) [27] and Multikernel Adaptive Filtering (MKAF)[28], where a combination of several kernels is used instead of a single kernel. However, there is still a shortcoming in the mixture correntropy that only allows the combination of zero-mean Gaussian kernels, which may perform poorly under some complex non-Gaussian noises such as those from multimodal distributions. To further improve the learning performance, in the present work, we propose a novel concept of multi-kernel correntropy (MKC), where each component of the mixture Gaussian kernel can be centered at a different location (not limited to zero-mean). Some important properties of the MKC are also studied. The MKC involves more free parameters than the MC, so a challenging issue is how to determine the free parameters in a practical application. To address this issue, we propose an efficient approach in this paper to optimize the free parameters in MKC by minimizing a distance between the mixture Gaussian function and the error's probability density function (PDF). Experimental results have confirmed the satisfactory performance of the learning methods under maximum multi-kernel correntropy criterion (MMKCC). Due to its excellent flexibility and robustness, the proposed MKC has great potential to be applied in many fields involving complex noise disturbances, such as biomedical engineering, remote sensing, autonomous systems and many others.

The rest of the paper is organized as follows. In section II, we define the MKC and present several properties. In section III, we propose an effective method to optimize the free parameters in MKC. Experimental results are then presented in section IV and finally, conclusion is given in section V.

## II. MULTI-KERNEL CORRENTROPY

### A. Definitions

Given two random variables $X \in \mathbf{R}$ and $Y \in \mathbf{R}$ with joint PDF $p_{XY}(x, y)$, correntropy is defined by [23]

$$V(X,Y) = \mathbf{E}[\kappa(X,Y)] = \iint \kappa(x,y)p_{XY}(x,y)dxdy \quad (1)$$

where $\kappa(.,.)$ is usually a radial kernel, and $\mathbf{E}[.]$ denotes the expectation operator. If the kernel function $\kappa(.,.)$ satisfies Mercer's condition, correntropy can be expressed as a correlation measure in a functional Hilbert space $\mathcal{F}$:

$$V(X,Y) = \mathbf{E}\left[\langle \varphi(X), \varphi(Y) \rangle_{\mathcal{F}}\right] \quad (2)$$

where $\varphi(.)$ is a nonlinear mapping induced by the kernel to transform the variables from the original space to the functional space $\mathcal{F}$, and $\langle ., . \rangle_{\mathcal{F}}$ stands for the inner product in $\mathcal{F}$. Without explicit mention, the kernel function in correntropy is the well-known Gaussian kernel:

$$\kappa(X,Y) = \kappa_{\sigma}(e) = \frac{1}{\sqrt{2\pi}\sigma}\exp\left(-\frac{e^2}{2\sigma^2}\right) \quad (3)$$

where $e = X - Y$ is the error between $X$ and $Y$, and $\sigma$ is the kernel bandwidth ($\sigma > 0$). It is easy to understand that correntropy measures how similar two random variables are in a local region of the error space controlled by the kernel bandwidth. Correntropy can easily be estimated from finite samples:

$$\hat{V}_{\sigma}(X,Y) = \frac{1}{N}\sum_{i=1}^{N}\kappa_{\sigma}(x_i - y_i) \quad (4)$$

where $\{x_i, y_i\}_{i=1}^{N}$ are $N$ samples of the random variables $X$ and $Y$. In particular, the function $CIM(\tilde{X}, \tilde{Y}) = \sqrt{\kappa_{\sigma}(0) - \hat{V}_{\sigma}(X,Y)}$ defines a metric, namely the correntropy induced metric (CIM) in the sample space, where $\tilde{X} = [x_1, \cdots, x_N]^T, \tilde{Y} = [y_1, \cdots, y_N]^T$. The CIM behaves like an $L_2$ norm distance if samples are close and like an $L_1$ norm distance as samples get further apart and eventually will approach the $L_0$ norm as samples far apart. This property elucidates the robustness of correntropy for outlier rejection. Under the maximum correntropy criterion (MCC), the detrimental effect of outliers can effectively be eliminated by maximizing the correntropy between the model output and target response [29].

The kernel function in correntropy is usually limited to a zero-mean Gaussian kernel and this may seriously restricts its performance when used as a cost function in machine learning. To improve the learning performance, the mixture correntropy (MC) was proposed in a recent paper [24] by using a linear combination of several zero-mean Gaussian kernels (with different bandwidths) as the kernel function. The mixture correntropy with $m$ sub-kernels is

$$V_{\lambda,\sigma}(X,Y) = \sum_{i=1}^{m}\lambda_i V_{\sigma_i}(X,Y)$$
$$= \mathbf{E}\left[\sum_{i=1}^{m}\lambda_i\kappa_{\sigma_i}(X - Y)\right] \quad (5)$$
$$= \iint \left(\sum_{i=1}^{m}\lambda_i\kappa_{\sigma_i}(x - y)\right)p_{XY}(x,y)dxdy$$
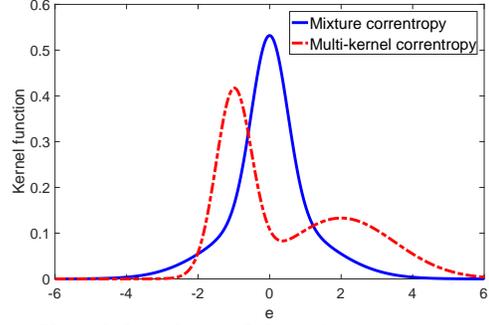


Fig. 1: Kernel functions of the mixture correntropy and multi-kernel correntropy

where $\boldsymbol{\lambda} = [\lambda_1, \lambda_2, \cdots, \lambda_m]^T$ is the mixture coefficient vector, and $\boldsymbol{\sigma} = [\sigma_1, \sigma_2, \cdots, \sigma_m]^T$ is the bandwidth vector. Usually, the mixture coefficient vector satisfies $\sum_{i=1}^{m}\lambda_i = 1$ with $\lambda_i \geq 0(i = 1, \cdots, m)$. In [24], for simplicity, only the case of $m = 2$ is considered. There is still a limitation in the mixture correntropy, that is, all the sub-kernels are centered at zero. To solve this limitation and further enhance the learning performance, in the present paper, we propose a more general definition of correntropy, namely, the multi-kernel correntropy (MKC), in which the sub-kernels can be centered at different locations (not limited to zero-mean). Specifically, the MKC between random variables $X$ and $Y$ is defined by

$$V_{\boldsymbol{\lambda},\boldsymbol{c},\boldsymbol{\sigma}}(X,Y) = \mathbf{E}\left[\sum_{i=1}^{m}\lambda_i\kappa_{\sigma_i}(X - Y - c_i)\right]$$
$$= \iint \left(\sum_{i=1}^{m}\lambda_i\kappa_{\sigma_i}(x - y - c_i)\right)p_{XY}(x,y)dxdy \quad (6)$$

where $\boldsymbol{c} = [c_1, c_2, \cdots, c_m]^T \in \mathbf{R}^m$ is the center vector.

*Remark*: The kernel function in the above MKC is a multi-Gaussian function that usually does not satisfy Mercer's condition. This is not a problem, however, because for a similarity measure the Mercer's condition is not necessary.

The MKC $V_{\boldsymbol{\lambda},\boldsymbol{c},\boldsymbol{\sigma}}(X,Y)$ will reduce to the MC $V_{\boldsymbol{\lambda},\boldsymbol{\sigma}}(X,Y)$ when $\boldsymbol{c} = [0, \cdots, 0]^T$. Fig. 1 shows the kernel functions of the mixture correntropy($m = 2, \lambda_1 = 0.5, \lambda_2 = 0.5, \sigma_1 = 0.5, \sigma_2 = 1.5$) and multi-kernel correntropy($m = 2, \lambda_1 = 0.5, \lambda_2 = 0.5, \sigma_1 = 0.5, \sigma_2 = 1.5, c_1 = -1.0, c_2 = 2.0$). Compared with the MC, the MKC is much more general and flexible and can adapt to more complicated error distribution, such as skewed, multi-peak, discrete-valued distribution, and hence it may achieve much better performance with proper setting of the centers when used as a cost function in machine learning. However, the MKC contains $3m$ free parameters, which have to be determined in practical applications. We will develop an efficient method in section IV to determine these free parameters.

### B. Properties

In the following, we present several properties of the MKC. The first and second properties are very straightforward and will not be proved here.

*Property 1*: The MKC $V_{\boldsymbol{\lambda},\boldsymbol{c},\boldsymbol{\sigma}}(X,Y)$ is positive and bounded: $0 < V_{\boldsymbol{\lambda},\boldsymbol{c},\boldsymbol{\sigma}}(X,Y) \leq \sum_{i=1}^{m} \frac{\lambda_i}{\sqrt{2\pi}\sigma_i}$

*Property 2*: The MKC $V_{\boldsymbol{\lambda},\boldsymbol{c},\boldsymbol{\sigma}}(X,Y)$ involves all the even moments of the error $e = X - Y$ about the centers $\{c_i\}$, that is,

$$V_{\boldsymbol{\lambda},\boldsymbol{c},\boldsymbol{\sigma}}(X,Y) = \sum_{i=1}^{m}\left( \frac{\lambda_i}{\sqrt{2\pi}\sigma_i} \sum_{n=0}^{\infty} \frac{(-1)^n}{2^n n!}\mathbf{E}\left[ \frac{(e-c_i)^{2n}}{\sigma_i^{2n}} \right] \right) \quad (7)$$

*Remark*: As $\{\sigma_i\}$ increases, the high-order moments will decay fast, and the second-order moments will tend to dominate the value.

*Property 3*: As $min\{\sigma_i\}$ is large enough, it holds that

$$V_{\boldsymbol{\lambda},\boldsymbol{c},\boldsymbol{\sigma}}(X,Y) \approx \sum_{i=1}^{m} \frac{\lambda_i}{\sqrt{2\pi}\sigma_i}\left( 1 - \frac{1}{2\sigma_i^2}\mathbf{E}\left[ (e-c_i)^2 \right] \right) \quad (8)$$

*Proof*: Since $\exp(x) \approx 1 + x$ for $x$ small enough, as $min\{\sigma_i\}$ is large enough, we have

$$\begin{aligned} V_{\boldsymbol{\lambda},\boldsymbol{c},\boldsymbol{\sigma}}(X,Y) &= \mathbf{E}\left[ \sum_{i=1}^{m} \lambda_i \kappa_{\sigma_i}(X - Y - c_i) \right] \\ &\approx \mathbf{E}\left[ \sum_{i=1}^{m} \frac{\lambda_i}{\sqrt{2\pi}\sigma_i}\left( 1 - \frac{(e-c_i)^2}{2\sigma_i^2} \right) \right] \quad (9) \\ &= \sum_{i=1}^{m} \frac{\lambda_i}{\sqrt{2\pi}\sigma_i}\left( 1 - \frac{1}{2\sigma_i^2}\mathbf{E}\left[ (e-c_i)^2 \right] \right) \end{aligned}$$

which completes the proof.

*Remark*: According to Property 3, when $min\{\sigma_i\}$ is very large, maximizing the MKC will be equivalent to minimizing a weighted sum of the error's second-order moments about the centers $\{c_i\}$.

*Property 4*: Let $p_e(.)$ be the PDF of the error variable $e = X - Y$. It holds that

$$\lim_{max\{\sigma_i\}\to 0+} V_{\boldsymbol{\lambda},\boldsymbol{c},\boldsymbol{\sigma}}(X,Y) = \sum_{i=1}^{m} \lambda_i p_e(c_i) \quad (10)$$

*Proof*: When $max\{\sigma_i\}$ shrinks to zero, the Gaussian kernel function $\kappa_{\sigma_i}(.)$ will approach the Dirac delta function $\delta(.)$. Thus we have

$$\begin{aligned} \lim_{max\{\sigma_i\}\to 0+} V_{\boldsymbol{\lambda},\boldsymbol{c},\boldsymbol{\sigma}}(X,Y) &= \lim_{max\{\sigma_i\}\to 0+} \mathbf{E}\left[ \sum_{i=1}^{m} \lambda_i \kappa_{\sigma_i}(e - c_i) \right] \\ &= \lim_{max\{\sigma_i\}\to 0+} \sum_{i=1}^{m} \lambda_i \int \kappa_{\sigma_i}(\varepsilon - c_i)p_e(\varepsilon)d\varepsilon \\ &= \sum_{i=1}^{m} \lambda_i \int \delta(\varepsilon - c_i)p_e(\varepsilon)d\varepsilon \\ &= \sum_{i=1}^{m} \lambda_i p_e(c_i) \end{aligned}$$
$$(11)$$

which completes the proof.

*Remark*: According to Property 4, when $max\{\sigma_i\}$ is very small, the MKC will approach a weighted sum of the values of $p_e(\varepsilon)$ evaluated at $\varepsilon = c_i \; (i = 1, \cdots, m)$.

## III. MAXIMUM MULTI-KERNEL CORRENTROPY CRITERION

The proposed MKC can be used to build new cost functions in many machine learning applications. Consider a supervised learning setting where the goal is to optimize a model $M$ that receives a random variable $X$ and outputs $Y = M(X)$ that should approximate a target variable (or teaching variable) $T$. Here $M(.)$ denotes an unknown mapping from the input to output that needs to be learned. A central problem in this learning task is the definition of a loss function (or a similarity measure) to compare $Y$ with $T$. The well-known minimum mean square error (MMSE) criterion has been the workhorse of supervised learning, which aims to minimize the MSE cost $\mathbf{E}\left[ e^2 \right]$ with $e = T - Y$ being the error variable. The combination of the linear feedforward model and MSE yields a set of equations that can be solved analytically. However, MSE is only optimal when the error variable is Gaussian distributed, which is seldom the case in real world applications. The error distributions tend to be skewed and with long tails, which create problems for MSE. Therefore, many "optimal solutions" are indeed not practical, simply because of the criterion that is used in the optimization. Many non-MSE optimization criterion were proposed in the literature to address the limitations of the MSE. The maximum correntropy criterion (MCC) is one of the hotspots of current research, which performs very well particularly when the error distribution is heavy-tailed [30]. Under the MCC, the model is optimized (or trained) to maximize the correntropy between the target $T$ and output $Y$:

$$\begin{aligned} M^* &= \arg\max_{M\in\mathbf{M}} V_\sigma(T,Y) \\ &= \arg\max_{M\in\mathbf{M}} \mathbf{E}\left[ \kappa_\sigma(e) \right] \end{aligned} \quad (12)$$

where $M^*$ denotes the optimal model and $\mathbf{M}$ stands for the hypothesis space. To improve the learning performance, the maximum mixture correntropy criterion (MMCC) was proposed in [24]. To further improve the flexibility and robustness, in the present paper, we propose the maximum multi-kernel correntropy criterion (MMKCC), where the optimal model is obtained by maximizing the MKC, that is

$$\begin{aligned} M^* &= \arg\max_{M\in\mathbf{M}} V_{\boldsymbol{\lambda},\boldsymbol{c},\boldsymbol{\sigma}}(T,Y) \\ &= \arg\max_{M\in\mathbf{M}} \mathbf{E}\left[ \sum_{i=1}^{m} \lambda_i \kappa_{\sigma_i}(e - c_i) \right] \end{aligned} \quad (13)$$

In a practical situation, given finite input-target samples $\{x_j, t_j\}_{j=1}^{N}$, the model can be trained through maximizing a sample estimator of the MKC:

$$\begin{aligned} M^* &= \arg\max_{M\in\mathbf{M}} \hat{V}_{\boldsymbol{\lambda},\boldsymbol{c},\boldsymbol{\sigma}}(T,Y) \\ &= \arg\max_{M\in\mathbf{M}} \frac{1}{N}\sum_{j=1}^{N}\sum_{i=1}^{m} \lambda_i \kappa_{\sigma_i}(e_j - c_i) \end{aligned} \quad (14)$$

where $e_j = t_j - y_j = t_j - M(x_j)$ is the $j$-th error sample.

In the following, we present a simple example to show how to solve the optimal solution under MMKCC. Consider

a linear-in-parameter (LIP) model in which the $j$-th output sample is

$$
\begin{aligned}
y_j &= \boldsymbol{h}_j \boldsymbol{\beta} \\
&= [\varphi_1(\boldsymbol{x}_j), \varphi_2(\boldsymbol{x}_j), \cdots, \varphi_L(\boldsymbol{x}_j)][\beta_1, \beta_2, \cdots, \beta_L]^T
\end{aligned} \tag{15}
$$

where $\boldsymbol{h}_j = [\varphi_1(\boldsymbol{x}_j), \varphi_2(\boldsymbol{x}_j), \cdots, \varphi_L(\boldsymbol{x}_j)] \in \mathbf{R}^L$ is the $j$-th nonlinearly mapped input vector (a row vector), with $\varphi_l(.)$ being the $l$-th nonlinear mapping function ($l = 1, 2, \cdots, L$), and $\boldsymbol{\beta} = [\beta_1, \beta_2, \cdots, \beta_L]^T \in \mathbf{R}^L$ is the output weight vector to be learned. Based on the MMKCC, the optimal weight vector $\boldsymbol{\beta}^*$ can be solved by maximizing the following objective function:

$$
\begin{aligned}
\boldsymbol{\beta}^* &= \arg\max_{\boldsymbol{\beta} \in \mathbf{R}^L} J(\boldsymbol{\beta}) \\
&= \frac{1}{N} \sum_{j=1}^{N} \sum_{i=1}^{m} \lambda_i \kappa_{\sigma_i}(e_j - c_i) - \gamma \|\boldsymbol{\beta}\|^2
\end{aligned} \tag{16}
$$

where $e_j = t_j - \boldsymbol{h}_j \boldsymbol{\beta}$, and $\gamma \geq 0$ is a regularization parameter. Setting $\partial J(\boldsymbol{\beta})/\partial \beta = \mathbf{0}$, we have

$$
\frac{1}{N} \sum_{j=1}^{N} \sum_{i=1}^{m} \frac{\lambda_i}{\sigma_i^2} \kappa_{\sigma_i}(e_j - c_i)(e_j - c_i)\boldsymbol{h}_j^T - 2\gamma\boldsymbol{\beta} = 0
$$

$$
\Rightarrow \sum_{j=1}^{N} \sum_{i=1}^{m} \frac{\lambda_i}{\sigma_i^2} \kappa_{\sigma_i}(e_j - c_i)(t_j - \boldsymbol{h}_j\boldsymbol{\beta} - c_i)\boldsymbol{h}_j^T - \gamma'\boldsymbol{\beta} = 0 \tag{17}
$$

$$
\Rightarrow \sum_{j=1}^{N} \psi(e_j)\boldsymbol{h}_j^T \boldsymbol{h}_j \beta + \gamma'\boldsymbol{\beta} = \sum_{j=1}^{N} \psi(e_j)t_j \boldsymbol{h}_j^T - \sum_{j=1}^{N} \zeta(e_j)\boldsymbol{h}_j^T
$$

where $\psi(e_j) = \sum_{i=1}^{m} \frac{\lambda_i}{\sigma_i^2} \kappa_{\sigma_i}(e_j - c_i)$, $\zeta(e_j) = \sum_{i=1}^{m} \frac{\lambda_i c_i}{\sigma_i^2} \kappa_{\sigma_i}(e_j - c_i)$, and $\gamma' = 2N\gamma$. From (17), one can easily derive

$$
\begin{aligned}
\boldsymbol{\beta} &= \left( \sum_{j=1}^{N} \psi(e_j)\boldsymbol{h}_j^T \boldsymbol{h}_j + \gamma' I \right)^{-1} \left( \sum_{j=1}^{N} \psi(e_j)t_j \boldsymbol{h}_j^T - \sum_{j=1}^{N} \zeta(e_j)\boldsymbol{h}_j^T \right) \\
&= \left( \mathbf{H}^T \boldsymbol{\Lambda} \mathbf{H} + \gamma' I \right)^{-1} \left( \mathbf{H}^T \boldsymbol{\Lambda} \mathbf{T} - \mathbf{H}^T \boldsymbol{\theta} \right)
\end{aligned} \tag{18}
$$

where $\mathbf{H} = [\boldsymbol{h}_{jl}]$ is an $N \times L$ dimensional matrix with $h_{jl} = \varphi_l(\boldsymbol{x}_j)$, $\boldsymbol{\Lambda}$ is an $N \times N$ diagonal matrix with diagonal elements $\boldsymbol{\Lambda}_{jj} = \psi(e_j)$, $\mathbf{T} = [t_1, \cdots, t_N]^T$, and $\boldsymbol{\theta} = [\zeta(e_1), \cdots, \zeta(e_N)]^T$.

The equation (18) is not a closed-form solution and it is actually a fixed-point equation because the diagonal matrix $\boldsymbol{\Lambda}$ and vector $\boldsymbol{\theta}$ on the right-hand side depend on the weight vector $\boldsymbol{\beta}$ through $e_j = t_j - \boldsymbol{h}_j \boldsymbol{\beta}$. Thus, the optimal solution of $\boldsymbol{\beta}$ can be obtained via a fixed-point iterative algorithm under MMKCC (FP-MMKCC), as described in **Algorithm 1**.

The computational complexities of some steps are given in Table I. Then, the computational complexity of the FP-MMKCC algorithm is $\left[ 2L^2N + 8LN + 21mN - 2N - L^2 + O(L^3) \right] T_{\text{FP}}$, where $T_{\text{FP}}$ is the fixed-point iteration number. Since the fixed-point iteration number $T_{\text{FP}}$ is relatively small in general, the computational complexity of the FP-MMKCC algorithm is moderate. Moreover, a sufficient condition to guarantee the convergence of the FP-MMKCC algorithm can be obtained (See APPENDIX A).

---

**Algorithm 1** FP-MMKCC algorithm

**Input:** training samples $\{\boldsymbol{x}_i, t_i\}_{i=1}^{N}$, number of nonlinear mappers $L$, mixture coefficient vector $\boldsymbol{\lambda}$, bandwidth vector $\boldsymbol{\sigma}$, center vector $\boldsymbol{c}$, regularization parameter $\gamma'$, maximum iteration number $K$, termination tolerance $\xi$ and the initial weight vector $\boldsymbol{\beta}_0=\mathbf{0}$.

**Output:** weight vector $\boldsymbol{\beta}$

1: **for all** $k = 1, 2, ..., K$ **do**
2:     Compute the errors based on $\boldsymbol{\beta}_{k-1}$: $e_i = t_i - \boldsymbol{h}_i \boldsymbol{\beta}_{k-1}$, $i = 1, 2, \cdots, N$
3:     Compute the diagonal matrix $\boldsymbol{\Lambda}$: $\boldsymbol{\Lambda}_{jj} = \sum_{i=1}^{m} \frac{\lambda_i}{\sigma_i^2} \kappa_{\sigma_i}(e_j - c_i), j = 1, 2, \cdots, N$
4:     Compute the vector $\boldsymbol{\theta}$: $\boldsymbol{\theta} = [\zeta(e_1), \cdots, \zeta(e_N)]^T$
5:     Update the weight vector $\boldsymbol{\beta}$: $\boldsymbol{\beta}_k = \left( \mathbf{H}^T \boldsymbol{\Lambda} \mathbf{H} + \gamma' \mathbf{I} \right)^{-1} \left( \mathbf{H}^T \boldsymbol{\Lambda} \mathbf{T} - \mathbf{H}^T \boldsymbol{\theta} \right)$
6:     **Until** $|J(\boldsymbol{\beta}_k) - J(\boldsymbol{\beta}_{k-1})| < \xi$
7: **end for**

---

TABLE I: Computational complexity for each iteration of the FP-MMKCC algorithm

| Step | Addition/subtraction and multiplication | Division, matrix inversion, and exponentiation |
|------|------------------------------------------|------------------------------------------------|
| 2 | $2LN$ | 0 |
| 3 | $5mN - N$ | $5mN$ |
| 4 | $6mN - N$ | $5mN$ |
| 5 | $2L^2N + 6LN - L^2$ | $O(L^3)$ |

## IV. DETERMINATION OF FREE PARAMETERS IN MMKCC

One of the most challenging problems in MMKCC is how to determine the $3m$ free parameters, namely the vectors $\boldsymbol{\lambda} = [\lambda_1, \lambda_2, \cdots, \lambda_m]^T$, $\boldsymbol{c} = [c_1, c_2, \cdots, c_m]^T$ and $\boldsymbol{\sigma} = [\sigma_1, \sigma_2, \cdots, \sigma_m]^T$. If this problem is not solved, the MMKCC will not be practical. To address this problem, we consider again the supervised learning setting in the previous section. First, we divide the MMKCC into three terms:

$$
\begin{aligned}
V_{\boldsymbol{\lambda}, \boldsymbol{c}, \boldsymbol{\sigma}}(T, Y) &= \mathbf{E} \left[ \sum_{i=1}^{m} \lambda_i \kappa_{\sigma_i}(e - c_i) \right] \\
&= \frac{1}{2} \int \left( \sum_{i=1}^{m} \lambda_i \kappa_{\sigma_i}(\varepsilon - c_i) \right)^2 d\varepsilon + \frac{1}{2} \int (p_e(\varepsilon))^2 d\varepsilon \\
&\quad - \frac{1}{2} \int \left( \sum_{i=1}^{m} \lambda_i \kappa_{\sigma_i}(\varepsilon - c_i) - p_e(\varepsilon) \right)^2 d\varepsilon
\end{aligned} \tag{19}
$$

The first term is independent of the model $M$, so we have

$$
\begin{aligned}
M^* &= \arg\max_{M \in \mathbf{M}} V_{\boldsymbol{\lambda}, \boldsymbol{c}, \boldsymbol{\sigma}}(T, Y) \\
&= \arg\max_{M \in \mathbf{M}} U_{\boldsymbol{\lambda}, \boldsymbol{c}, \boldsymbol{\sigma}}(T, Y)
\end{aligned} \tag{20}
$$

where $U_{\boldsymbol{\lambda}, \boldsymbol{c}, \boldsymbol{\sigma}}(T, Y) = \frac{1}{2} \int (p_e(\varepsilon))^2 d\varepsilon - \frac{1}{2} \int \left( \sum_{i=1}^{m} \lambda_i \kappa_{\sigma_i}(\varepsilon - c_i) - p_e(\varepsilon) \right)^2 d\varepsilon$.

To determine the free parameters, in this study we propose

$$(M^*, \boldsymbol{\lambda}^*, \boldsymbol{c}^*, \boldsymbol{\sigma}^*) = \underset{M \in \mathbf{M}, \boldsymbol{\lambda} \in \boldsymbol{\Omega_\lambda}, \boldsymbol{c} \in \boldsymbol{\Omega_c}, \boldsymbol{\sigma} \in \boldsymbol{\Omega_\sigma}}{\arg\max} U_{\boldsymbol{\lambda}, \boldsymbol{c}, \boldsymbol{\sigma}}(T, Y)$$

$$= \underset{M \in \mathbf{M}, \boldsymbol{\lambda} \in \boldsymbol{\Omega_\lambda}, \boldsymbol{c} \in \boldsymbol{\Omega_c}, \boldsymbol{\sigma} \in \boldsymbol{\Omega_\sigma}}{\arg\max} \frac{1}{2} \int (p_e(\varepsilon))^2 d\varepsilon - \frac{1}{2} \int \left( \sum_{i=1}^{m} \lambda_i \kappa_{\sigma_i}(\varepsilon - c_i) - p_e(\varepsilon) \right)^2 d\varepsilon$$

$$= \underset{M \in \mathbf{M}, \boldsymbol{\lambda} \in \boldsymbol{\Omega_\lambda}, \boldsymbol{c} \in \boldsymbol{\Omega_c}, \boldsymbol{\sigma} \in \boldsymbol{\Omega_\sigma}}{\arg\max} -\frac{1}{2} \int \left( \sum_{i=1}^{m} \lambda_i \kappa_{\sigma_i}(\varepsilon - c_i) \right)^2 d\varepsilon + \mathbf{E} \left[ \sum_{i=1}^{m} \lambda_i \kappa_{\sigma_i}(e - c_i) \right] \tag{21}$$

$$(M^*, \boldsymbol{\lambda}^*, c^*, \sigma^*) = \underset{M \in \mathbf{M}, \boldsymbol{\lambda} \in \boldsymbol{\Omega_\lambda}, \boldsymbol{c} \in \boldsymbol{\Omega_c}, \boldsymbol{\sigma} \in \boldsymbol{\Omega_\sigma}}{\arg\max} -\frac{1}{2} \boldsymbol{\lambda}^T \left( \int \tilde{\boldsymbol{g}}(\varepsilon) \tilde{\boldsymbol{g}}(\varepsilon)^T d\varepsilon \right) \boldsymbol{\lambda} + \boldsymbol{\lambda}^T \tilde{\boldsymbol{h}}$$

$$= \underset{M \in \mathbf{M}, \boldsymbol{\lambda} \in \boldsymbol{\Omega_\lambda}, \boldsymbol{c} \in \boldsymbol{\Omega_c}, \boldsymbol{\sigma} \in \boldsymbol{\Omega_\sigma}}{\arg\max} -\frac{1}{2} \boldsymbol{\lambda}^T \tilde{\mathbf{K}} \boldsymbol{\lambda} + \boldsymbol{\lambda}^T \tilde{\boldsymbol{h}} \tag{23}$$

$$\tilde{\mathbf{K}} = \begin{bmatrix} \frac{1}{\sqrt{2\pi}\sqrt{\sigma_1^2 + \sigma_1^2}} \exp(-\frac{(c_1 - c_1)^2}{2(\sigma_1^2 + \sigma_1^2)}) & \cdots & \frac{1}{\sqrt{2\pi}\sqrt{\sigma_1^2 + \sigma_m^2}} \exp(-\frac{(c_1 - c_m)^2}{2(\sigma_1^2 + \sigma_m^2)}) \\ \vdots & \ddots & \vdots \\ \frac{1}{\sqrt{2\pi}\sqrt{\sigma_m^2 + \sigma_1^2}} \exp(-\frac{(c_m - c_1)^2}{2(\sigma_m^2 + \sigma_1^2)}) & \cdots & \frac{1}{\sqrt{2\pi}\sqrt{\sigma_m^2 + \sigma_m^2}} \exp(-\frac{(c_m - c_m)^2}{2(\sigma_m^2 + \sigma_m^2)}) \end{bmatrix} \tag{24}$$

the optimization in (21), where $\boldsymbol{\Omega_\lambda}$, $\boldsymbol{\Omega_c}$ and $\boldsymbol{\Omega_\sigma}$ denote the admissible sets of the parameter vectors $\boldsymbol{\lambda}$, $\boldsymbol{c}$ and $\boldsymbol{\sigma}$.

*Remark*: It is worth noting that the objective function $U_{\boldsymbol{\lambda}, \boldsymbol{c}, \boldsymbol{\sigma}}(T, Y)$ can be expressed as

$$U_{\boldsymbol{\lambda}, \boldsymbol{c}, \boldsymbol{\sigma}}(T, Y) = \frac{1}{2} QIP(e)$$
$$- \frac{1}{2} D_{ED} \left( \sum_{i=1}^{m} \lambda_i \kappa_{\sigma_i}(\varepsilon - c_i) \| p_e(\varepsilon) \right) \tag{22}$$

where $QIP(e) = \int (p_e(\varepsilon))^2 d\varepsilon$ is the quadratic information potential (QIP) [31] of the error $e$, and $D_{ED}(. \| .)$ denotes the Euclidean distance between PDFs [32, 33], defined by $D_{ED}(p(x) \| q(x)) = \int (p(x) - q(x))^2 dx$. Therefore, maximizing the objective function $U_{\boldsymbol{\lambda}, \boldsymbol{c}, \boldsymbol{\sigma}}(T, Y)$ will try to maximize the QIP (or minimize Renyi's quadratic entropy) of the error and at the same time, minimize the Euclidean distance between the multi-Gaussian kernel function and the error's PDF.

If $N$ error samples $\{e_j\}_{j=1}^{N}$ are available, we have $\mathbf{E} \left[ \sum_{i=1}^{m} \lambda_i \kappa_{\sigma_i}(e - c_i) \right] \approx \boldsymbol{\lambda}^T \tilde{\boldsymbol{h}}$, where $\tilde{\boldsymbol{h}} = \frac{1}{N} \sum_{j=1}^{N} \tilde{\boldsymbol{g}}(e_j)$, with $\tilde{\boldsymbol{g}}(e_j) = [\kappa_{\sigma_1}(e_j - c_1), \cdots, \kappa_{\sigma_m}(e_j - c_m)]^T$. Thus by (21), we have (23), where $\tilde{\mathbf{K}}$ is expressed in (24).

According to (23), the model $M$ and $3m$ free parameters are jointly optimized via maximizing the objective function $\hat{U}_{\boldsymbol{\lambda}, \boldsymbol{c}, \boldsymbol{\sigma}}(T, Y) = -\frac{1}{2} \boldsymbol{\lambda}^T \tilde{\mathbf{K}} \boldsymbol{\lambda} + \boldsymbol{\lambda}^T \tilde{\boldsymbol{h}}$. This is a very complicated optimization problem. To simplify the optimization, one can adopt an alternative optimization method: i) given a model (hence the $N$ error samples are given), we solve the free parameters by maximizing $\hat{U}_{\boldsymbol{\lambda}, \boldsymbol{c}, \boldsymbol{\sigma}}(T, Y)$ (with error samples fixed); ii) after the free parameters are determined, we solve a new model by maximizing $\hat{U}_{\boldsymbol{\lambda}, \boldsymbol{c}, \boldsymbol{\sigma}}(T, Y)$ (with free parameters fixed).

In a practical application, there are usually two approaches to find the free parameters and the optimal model. The first
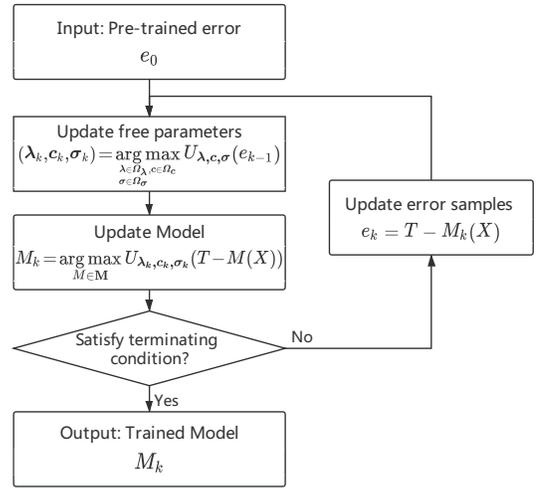


Fig. 2: Alternative optimization for model and free parameters

one is an online approach, in which the model is optimized by an iterative method and at each iteration, the $3m$ free parameters are determined based on the error samples at that iteration. The second one is a two-stage approach, which contains two stages: 1) train the model using a simple method (usually with very few free parameters) to obtain the error samples, and determine the $3m$ free parameters based on these errors; 2) train the model again under the MMKCC with the obtained free parameters, and during the training these free parameters are fixed. The above procedure can be repeated until convergence and the flow chart is shown in Fig. 2.

Next, we describe how to determine the $3m$ free parameters given a model. First, to simplify the optimization, we just apply some clustering technique such as the K-means on the error samples to obtain the center vector $\boldsymbol{c}^*$ (whose elements are the clustering centers). Then by (23), one can easily obtain
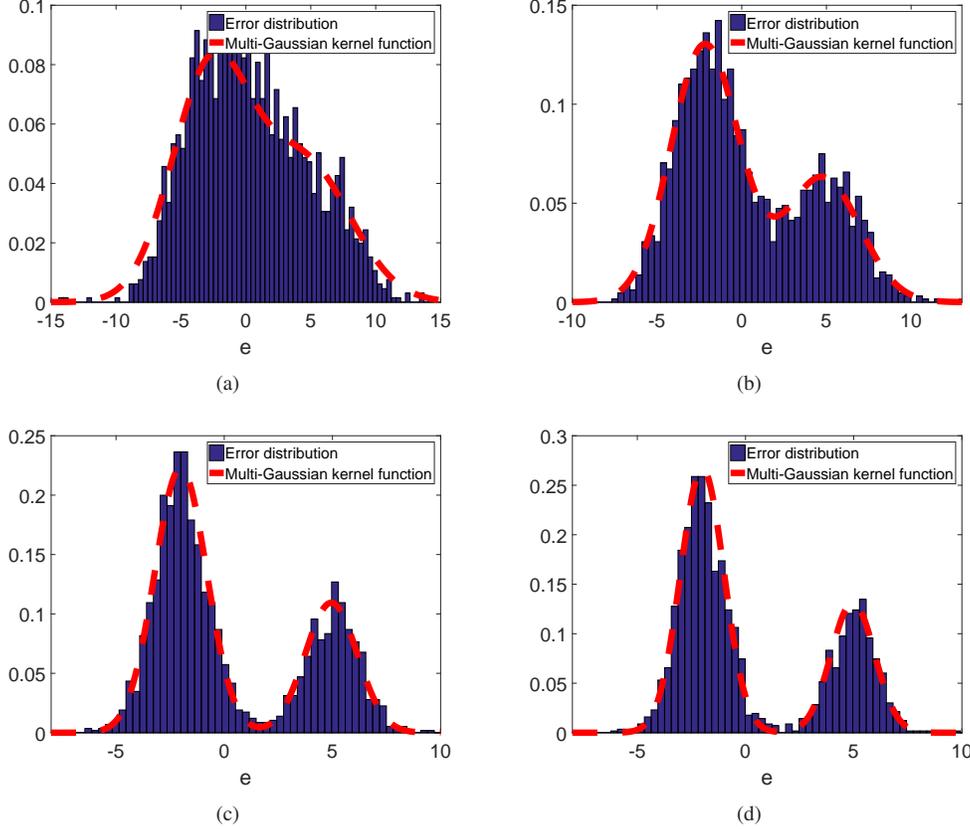
Fig. 3: Error distributions and multi-Gaussian kernel functions at different fixed-point iterations: (a) first iteration; (b)second iteration, (c)third iteration, (d)fourth iteration

the mixture coefficient vector:

$$\boldsymbol{\lambda}^* = \tilde{\mathbf{K}}^{-1}\tilde{\boldsymbol{h}} \tag{25}$$

In order to avoid numerical problem in the matrix inversion, a regularized solution can be used:

$$\boldsymbol{\lambda}^* = (\tilde{\mathbf{K}} + \eta\mathbf{I})^{-1}\tilde{\boldsymbol{h}} \tag{26}$$

where $\eta$ is a regularization parameter. Substituting (26) into (23), we solve the bandwidth vector as follows:

$$\boldsymbol{\sigma}^* = \arg\max_{\boldsymbol{\sigma}\in\boldsymbol{\Omega}_\sigma} -\frac{1}{2}\Big[\big(\tilde{\mathbf{K}} + \eta I\big)^{-1}\tilde{\boldsymbol{h}}\Big]^T \tilde{\mathbf{K}}(\tilde{\mathbf{K}} + \eta I)^{-1}\tilde{\boldsymbol{h}} \\ + \Big[\big(\tilde{\mathbf{K}} + \eta I\big)^{-1}\tilde{\boldsymbol{h}}\Big]^T \tilde{\boldsymbol{h}} \tag{27}$$

In order to reduce the computational complexity of the optimization problem in (27), one can alternately optimize every dimension of the bandwidth vector over a finite set of values. Specifically, given a finite set of bandwidths $\boldsymbol{\Omega}_\sigma$, we optimize each element of the bandwidth vector $\boldsymbol{\sigma}$ one by one and repeat this procedure until convergence. The proposed procedure for free parameters determination is summarized in **Algorithm 2**.

## V. EXPERIMENTAL RESULTS

In this section, we present experimental results to demonstrate the desirable performance of learning methods under the

---

**Algorithm 2** Determination of the free parameters

**Input:** error samples $\{e_j\}_{j=1}^N$, parameter dimension $m$, regularization parameter $\eta$, a finite set of bandwidths $\boldsymbol{\Omega}_\sigma$ and initialize $\sigma_1 = \cdots = \sigma_m = \sigma_0$.

**Output:** free parameters $\boldsymbol{\lambda}^*, \boldsymbol{c}^*, \boldsymbol{\sigma}^*$

1: Determine the center vector $\boldsymbol{c}^*$ by applying the K-means clustering on the error samples $\{e_j\}_{j=1}^N$

2: Alternately optimize every dimension of the bandwidth vector $\boldsymbol{\sigma}$ and repeat $S$ times:

3: **for all** $s = 1, 2, ..., S$ **do**

4:     **for all** $i = 1, 2, ..., m$ **do**

5:       $\sigma_i^* = \arg\max\limits_{\sigma_i\in\boldsymbol{\Omega}_\sigma} -\frac{1}{2}\Big[\big(\tilde{\mathbf{K}} + \eta I\big)^{-1}\tilde{\boldsymbol{h}}\Big]^T \tilde{\mathbf{K}}(\tilde{\mathbf{K}} + \eta I)^{-1}\tilde{\boldsymbol{h}} +$
      $\Big[\big(\tilde{\mathbf{K}} + \eta I\big)^{-1}\tilde{\boldsymbol{h}}\Big]^T \tilde{\boldsymbol{h}}$, with $\boldsymbol{c} = \boldsymbol{c}^*$ and other $m - 1$ elements of $\boldsymbol{\sigma}$ being fixed

6:     **end for**

7: **end for**

8: Compute $\boldsymbol{\lambda}^* = (\tilde{\mathbf{K}} + \eta\mathbf{I})^{-1}\tilde{\boldsymbol{h}}$ with $\boldsymbol{\sigma} = \boldsymbol{\sigma}^*$ and $\boldsymbol{c} = \boldsymbol{c}^*$
    **Return:** $\boldsymbol{\lambda}^*, \boldsymbol{\sigma}^*, \boldsymbol{c}^*$

---

proposed MMKCC criterion (i.e. the FP-MMKCC algorithm). Without explicit mention, the dimension number is $m = 2$, the regularization parameter is $\eta = 10^{-4}$ and the iteration number $S$ is $S = 3$.

TABLE II: RMSEs and computing times(sec) of different learning criteria

| | | MSE | MCC | MMCC | MMKCC |
|---|---|---|---|---|---|
| case 1) | RMSE | $0.5427 \pm 0.3175$ | $0.0881 \pm 0.0431$ | $0.0831 \pm 0.0375$ | $\mathbf{0.0342 \pm 0.0259}$ |
| | TIME(sec) | N/A | $0.0832 \pm 0.0020$ | $0.1027 \pm 0.0026$ | $0.3328 \pm 0.0070$ |
| case 2) | RMSE | $0.5031 \pm 0.2483$ | $0.0754 \pm 0.0414$ | $0.0674 \pm 0.0334$ | $\mathbf{0.0224 \pm 0.0115}$ |
| | TIME(sec) | N/A | $0.0814 \pm 0.0018$ | $0.1014 \pm 0.0024$ | $0.3415 \pm 0.0075$ |
| case 3) | RMSE | $0.5494 \pm 0.3418$ | $0.0391 \pm 0.0191$ | $0.0353 \pm 0.0176$ | $\mathbf{0.0335 \pm 0.0168}$ |
| | TIME(sec) | N/A | $0.0841 \pm 0.0022$ | $0.1021 \pm 0.0027$ | $0.3297 \pm 0.0068$ |



Fig. 4: RMSE convergence curves of different learning criteria

TABLE III: Specification of the datasets

| Datasets | Features | Observations | |
|---|---|---|---|
| | | Training | Testing |
| Servo | 5 | 83 | 83 |
| Concrete | 9 | 515 | 515 |
| Airfoil | 5 | 751 | 751 |
| Yacht | 6 | 154 | 154 |

### A. Linear Regression

First, we consider a simple linear regression example where the input-target samples are generated by a two-dimensional linear system: $t_i = \boldsymbol{\beta}^{*T}\boldsymbol{x}_i + \rho_i$, where $\boldsymbol{\beta}^* = [1, 2]^T$ is the weight vector to be estimated, and $\rho_i$ denotes an additive noise. The input samples $\{\boldsymbol{x}_i\}$ are uniformly distributed over $[-2.0, 2.0] \times [-2.0, 2.0]$. The noise $\rho_i$ comprises two mutually independent noises, namely the inner noise $B_i$ and the outlier noise $O_i$. Specifically, $\rho_i$ is given by $\rho_i = (1 - g_i)B_i + g_iO_i$, where $g_i$ is a binary variable with probability mass $\Pr\{g_i = 1\} = p$, $\Pr\{g_i = 0\} = 1 - p$, $(0 \le p \le 1)$, which is assumed to be independent of both $B_i$ and $O_i$. In this example, $p$ is set at 0.1, and the outlier $O_i$ is drawn from a zero-mean Gaussian distribution with variance 10000. As for the inner noise $B_i$, we consider three cases: 1) $B_i \sim 0.5\mathcal{N}(4.0, 1.0) + 0.5\mathcal{N}(-4.0, 1.0)$, where $\mathcal{N}(u, \sigma^2)$ denotes the Gaussian density function with mean $u$ and variance $\sigma^2$ ; 2) $B_i \sim 1/3\mathcal{N}(5.0, 1.0) + 2/3\mathcal{N}(-2.0, 1.0)$. 3) $B_i \sim 0.5\mathcal{N}(0, 1.0) + 0.5\mathcal{N}(0, 5.0)$. The root mean squared error (RMSE) is employed to measure the performance, computed by $RMSE = \sqrt{\frac{1}{2}\|\boldsymbol{\beta}_k - \boldsymbol{\beta}^*\|^2}$, where $\boldsymbol{\beta}_k$ and $\boldsymbol{\beta}^*$ denote the estimated and the target weight vectors respectively.

We compare the performance of four learning criteria, namely MSE, MCC, MMCC and MMKCC. For MSE, there is a closed-form solution, so no iteration is needed. For MCC, MMCC and MMKCC, a fixed-point iteration is used to solve the model (see the **Algorithm 1** for the fixed-point algorithm under MMKCC). The mean $\pm$ deviation results of the RMSEs and computing times over 100 Monte Carlo runs are presented in Table II. In the simulation, the sample number is $N = 400$, the fixed-point iteration number is $K = 10$, and the initial weight vector is set to $\beta_0 = [0, 0]^T$. For each learning criterion, the parameters are experimentally selected to achieve

the best results, except that the free parameters of MMKCC are determined by **Algorithm 2**. The finite set $\Omega_\sigma$ in **Algorithm 2** is equally spaced over [0.1, 2.0] with step size 0.2. The simulations are carried out with MATLAB 8.6 running in Core 4 Quad, 3.4-GHZ CPU with 20-GB RAM. From Table II, we observe: i) MCC, MMCC and MMKCC can significantly outperform MSE although they have no closed-form solution; ii) MMKCC can achieve better performance than MCC and MMCC especially for noises with multi-peak or asymmetric distributions; iii) although the MMKCC is computationally more expensive than MCC and MMCC, the computing times of three learning criteria are in the same order of magnitude. For the noise case 2), the error distributions and multi-Gaussian kernel functions (determined by **Algorithm 2** after each fixed-point iteration) at different fixed-point iterations under MMKCC are shown in Fig. 3. As expected, the multi-Gaussian kernel function matches the error distribution very well at every iteration (as discussed earlier, the free parameters in MMKCC have been optimized to minimize the Euclidean distance between the multi-Gaussian kernel function and the error PDF). The average RMSE convergence curves of three learning criteria are illustrated in Fig. 4.

### B. Non-linear regression with benchmark datasets

In the second example, we show the superior performance of the MMKCC criterion in nonlinear regression with five benchmark data sets from UCI machine learning repository [34]. The descriptions of the data sets are given in Table III. In the experiment, the training and testing samples from each data set are randomly chosen and the data values are normalized into [0, 1.0]. The robust stochastic configuration networks (RSCN) is adopted as the regression model to be trained, which is a linear-in-parameter (LIP) model with randomly generated hidden nodes [35–37]. Under the MMKCC, the

TABLE IV: RMSEs and training times(sec) of several RSCN algorithms

| Datasets | RSCN | | RSC-MCC | | RSC-MMCC | | RSC-MMKCC | |
|---|---|---|---|---|---|---|---|---|
| | RMSE | Training Time | RMSE | Training Time | RMSE | Training Time | RMSE | Training Time |
| Hardware | 0.1000± 0.0325 | 0.0129± 0.0021 | 0.0754± 0.0254 | 0.1074± 0.0256 | 0.0734± 0.0223 | 0.1024± 0.0199 | **0.0719± 0.0208** | 0.3464± 0.0389 |
| Servo | 0.1293± 0.0322 | 0.0237± 0.0020 | 0.1211± 0.0237 | 0.0462± 0.0098 | 0.1181± 0.0228 | 0.0591± 0.0139 | **0.1169± 0.0224** | 0.2896± 0.0188 |
| Yacht | 0.0484± 0.0118 | 0.1458± 0.0089 | 0.0427± 0.0149 | 0.3273± 0.0422 | 0.0400± 0.0139 | 0.3348± 0.0389 | **0.0385± 0.0166** | 0.8753± 0.0456 |
| Airfoil | 0.0923± 0.0057 | 1.7234± 0.0811 | 0.0905± 0.0065 | 3.0372± 0.1145 | 0.0900± 0.0056 | 3.0522± 0.1230 | **0.0893± 0.0055** | 3.9181± 0.1379 |

TABLE V: Testing RMSEs of TDNNs trained under different criteria

| | MSE | MCC | MMCC | MMKCC |
|---|---|---|---|---|
| RMSE | 0.0427 | 0.0309 | 0.0302 | **0.0277** |

model is trained by the fixed-point iterative algorithm in **Algorithm 1** and we call it the RSC-MMKCC algorithm. In this example, the performance of the RSC-MMKCC is compared with that of several other stochastic configuration networks (SCN) based algorithms, including RSCN[36], RSC-MCC[37] and RSC-MMCC, where the RSC-MMCC can be viewed as RSC-MMKCC with $c = 0$. The parameters of each algorithm are selected through fivefold cross-validation, except that the free parameters of MMKCC are determined by **Algorithm 2**. The finite set $\Omega_\sigma$ in **Algorithm 2** is equally spaced over [0.1, 3.0] with step size 0.1. The training and testing RMSEs over 100 runs are presented in Table IV. Clearly, the RSC-MMKCC outperforms the RSCN, RSC-MCC and RSC-MMCC for all the data sets.

*Remark*: The parameter setting method of **Algorithm 2** may have similar or worse performances compared with the cross-validation. However, the cross-validation will take a lot of time when the parameter space is very large ($3m$ parameters for MKC). Thus, the cross-validation approach for RSC-MMKCC is not practical. Actually, the proposed parameter setting method is computationally much simpler than the cross-validation but can still achieve desirable performances (RSC-MMKCC performs better than RSC-MMCC and RSC-MCC in this example).

### C. Chaotic time series prediction

In the third example, we apply different learning criteria (MSE, MCC, MMCC, MMKCC) to train a time delay neural network (TDNN)[38] to predict the Mackey-Glass chaotic time series[39]. The TDNN has a single hidden layer and six nonlinear processing elements in the hidden layer, and its inputs consist of six delayed values. A sigmoid nonlinearity was used in each of the hidden processing elements, while the output processing element was linear. The sequence for training has an additive noise, and the training samples are generated by

$$x(t) = -bx(t-1) + \frac{ax(t-\tau)}{1 + x(t-\tau)^{10}} + \rho_t \tag{28}$$

with $b = 0.1$, $a = 0.2$, $\tau = 30$ and $\rho_t \sim 0.45\mathcal{N}(-0.05, 0.05) + 0.45\mathcal{N}(0.05, 0.05) + 0.1\mathcal{N}(0, 0.2)$. The TDNN is trained to
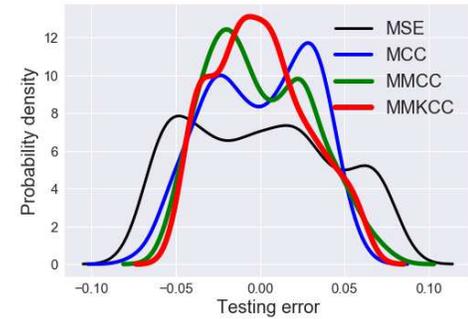


Fig. 5: Testing error PDFs of TDNNs trained under different criteria

predict the next sample of the time series by using six previous samples, with a segment of 200 samples. The trained networks are tested on clean data set (without additive noise) of length 1000. The kernel size of MCC is experimentally set at $\sigma = 2.0$, and the kernel sizes of MMCC are $\sigma_1 = 1.0$, $\sigma_2 = 2.0$, and the mixture coefficient in MMCC is $\alpha = 0.8$. For MMKCC, the finite set $\Omega_\sigma$ is equally spaced over [0.1, 3.0] with step size 0.1. The PDFs of the testing error averaged over 10 Monte Carlo runs are illustrated in Fig. 5 and the corresponding testing RMSEs are presented in Table V. Evidently, the TDNN trained under MMKCC achieves the best performance with the most concentrated error distribution and the lowest RMSE.

### D. EEG Classification

Electroencephalography (EEG) is a kind of multichannel electrophysiological signal recorded by electrodes placed on the scalp typically (also subdurally or in the cerebral cortex), which plays an important role in the brain–computer interface (BCI) systems [40, 41]. A BCI system can be defined as a system that translates the brain activity patterns into commands for an interactive application [40], and for an EEG-based BCI system, the goal is to effectively recognize the brain patterns of a user from the collected EEG signals. In view of the fact that EEG recordings are often contaminated by various artifacts, such as artifacts due to electrode displacement, motion artifacts, ocular artifacts and so on [42], the proposed RSC-

TABLE VI: Classification Accuracies of Different Algorithms on the Data Set IIa of BCI Competition IV

| Subject | KNN | SVM | RSCN | RSC-MCC | RSC-MMCC | RSC-MMKCC |
|---|---|---|---|---|---|---|
| A01 | 70.83 | 72.92 | 73.97±1.87 | 74.20±1.79 | 74.93±2.03 | **75.69±1.89** |
| A02 | 43.40 | 46.88 | 46.32±1.46 | 46.25±1.57 | 46.47±1.79 | **46.88±1.41** |
| A03 | 74.65 | 76.39 | 76.71±1.28 | 76.92±1.30 | 76.77±1.33 | **76.80±1.30** |
| A04 | 55.21 | 62.15 | 60.72±1.85 | 60.51±1.84 | 60.59±1.69 | **61.46±1.65** |
| A05 | 35.07 | 35.07 | **40.56±1.68** | 40.41±2.10 | 40.32±1.92 | 40.28±1.90 |
| A06 | 40.63 | 42.01 | 44.41±2.02 | **44.50±2.03** | 43.90±1.77 | 43.75±1.79 |
| A07 | 75.35 | 77.78 | 78.61±1.22 | 80.03±1.68 | 79.36±1.65 | **80.56±1.52** |
| A08 | 73.96 | 79.51 | 79.22±1.45 | 79.76±2.04 | **81.18±1.32** | 79.51±1.29 |
| A09 | 81.25 | 79.17 | 79.63±2.48 | 80.28±2.45 | 79.94±2.35 | **86.11±2.35** |
| Mean | 61.15 | 63.54 | 64.46±1.70 | 64.76±1.87 | 64.83±1.76 | **65.67±1.68** |

TABLE VII: The $p$-value of the Paired Sample T-Test between Classification Accuracies of MMCC and MMKCC

| Subject | A01 | A02 | A03 | A04 | A05 | A06 | A07 | A08 | A09 | Mean Accuracy |
|---|---|---|---|---|---|---|---|---|---|---|
| $p$-value | 0.0433 | 0.2206 | 0.8867 | 0.0067 | 0.8954 | 0.6432 | 0.0002 | $4.0*10^{-09}$ | $8.7*10^{-21}$ | $7.4*10^{-08}$ |

MMKCC could be a good candidate for the EEG classifier due to its excellent flexibility and robustness.

The benchmark data set adopted here is the data set IIa of BCI Competition IV [43], which consists of EEG data from 9 subjects. The BCI paradigm consisted of four different motor imagery tasks, i.e., left hand, right hand, both feet and tongue. Each subject includes two sessions, and each session is comprised of 6 runs. One run consists of 48 trials (12 for each class), yielding a total of 288 trials per session.

Considering the EEG feature extraction, the common spatial pattern (CSP) is an effective approach for multichannel EEG data concerning motor imagery tasks [44]. We adopt the CSP combined with the one-versus-one (OVO) [45] approach, which transforms the four-class classification problem into six cases of two-class classification. The first two spatial filters that correspond to the largest objective function values are used, and vice versa. Then the log variances of the spatially filtered EEG signals are used as the input features for classifiers. As a result, each trial is assigned with a 24-D feature.

Besides the aforementioned RSCN, RSC-MCC, RSC-MMCC and RSC-MMKCC, k-nearest neighbor (KNN) [46] and support vector machine (SVM) [47] are also chosen as the classifiers in EEG classification tasks for comparison. The parameters of each algorithm are selected through fivefold cross-validation, except that the free parameters of MMKCC are determined by **Algorithm 2**. Table VI shows the averaged classification accuracies of different algorithms after 30 Monte Carlo runs, and the highest accuracy for each subject is marked in bold. One can observe that the proposed RSC-MMKCC can achieve a higher classification accuracy on most subjects, and it also has the highest average classification accuracy over all the subjects. The standard deviations of KNN and SVM are zero since no random projection mechanism is adopted in them. With the Null Hypothesis $H_0 : \mu_d = 0$, where $d$ is the difference between the accuracies on pair, the $p$-values of the paired sample t-test between the classification

accuracy of MMKCC and MMCC are shown in Table VII. With the significance level $\alpha = 0.05$, we reject $H_0$ and state that we have significant evidence that accuracy difference between MMKCC and MMCC is NOT 0 for subject A01, A04, A07, A08 A09 and the mean accuracy, while the others are insignificant. Thus we can conclude that the MMKCC generally has better performance than the MMCC in statistics

## VI. CONCLUSION

A new generalized version of correntropy, called multi-kernel correntropy (MKC), was proposed in this study, where the kernel function is a mixture Gaussian kernel with different widths and centers. The original correntropy and the recently proposed mixture correntropy are both special cases of the new definition. Some important properties of the MKC were presented. In addition, a novel approach was proposed to determine the free parameters of MKC when used in supervised learning. The superior performance of the proposed learning method has been confirmed by experimental results of linear regression, nonlinear regression with benchmark datasets, chaotic time series prediction and EEG classfication.

## APPENDIX A
### CONVERGENCE ANALYSIS OF FP-MMKCC

The contraction mapping theorem (also known as the *Banach Fixed-Point Theorem*) provides an effective way to prove the convergence of a fixed-point algorithm [48, 49].

The FP-MMKCC alogrithm can be described as

$$\begin{aligned} \boldsymbol{\beta}_k &= (\mathbf{H^T \Lambda H} + \gamma' \mathbf{I})^{-1}(\mathbf{H^T \Lambda T} - \mathbf{H^T}\boldsymbol{\theta}) \\ &= (\mathbf{A}(\boldsymbol{\beta}_{k-1}) + \gamma' \mathbf{I})^{-1} \mathbf{B}(\boldsymbol{\beta}_{k-1}) \\ &= f(\boldsymbol{\beta}_{k-1}). \end{aligned} \tag{A.1}$$

According to the contraction mapping theorem, the convergence of a fixed-point algorithm is guaranteed if $\exists \delta > 0$ and

$0 < \alpha < 1$ such that if the initial weight vector $\left\|\boldsymbol{\beta}_0\right\|_p \le \delta$, and $\forall \boldsymbol{\beta} \in \{\boldsymbol{\beta} \in \mathbb{R}^L : \left\|\boldsymbol{\beta}_0\right\|_p \le \delta\}$, it holds that

$$\begin{cases} \left\|\boldsymbol{f}(\boldsymbol{\beta})\right\|_p \le \delta \\ \left\|\nabla_{\boldsymbol{\beta}}\boldsymbol{f}(\boldsymbol{\beta})\right\|_p = \left\|\dfrac{\partial \boldsymbol{f}(\boldsymbol{\beta})}{\partial \boldsymbol{\beta}^T}\right\|_p \le \alpha, \end{cases} \quad \text{(A.2)}$$

in which "$\|\cdot\|_p$" denotes an $\ell_p$-norm for a vector or an induced norm of a matrix, defined by $\left\|A\right\|_p = \max\limits_{\|X\|_p \ne 0} \left\|AX\right\|_p / \left\|X\right\|_p$, with $p \ge 1, A \in \mathbb{R}^{m \times m}, X \in \mathbb{R}^{m \times 1}$, and $\nabla_{\boldsymbol{\beta}}\boldsymbol{f}(\boldsymbol{\beta})$ denotes the $m \times m$ Jacobian matrix of $\boldsymbol{f}(\boldsymbol{\beta})$ with respect to $\boldsymbol{\beta}$, given by

$$\nabla_{\boldsymbol{\beta}}\boldsymbol{f}(\boldsymbol{\beta}) = \begin{bmatrix} \dfrac{\partial}{\partial \beta_1}\boldsymbol{f}(\boldsymbol{\beta}) & \dfrac{\partial}{\partial \beta_2}\boldsymbol{f}(\boldsymbol{\beta}) & \cdots & \dfrac{\partial}{\partial \beta_m}\boldsymbol{f}(\boldsymbol{\beta}) \end{bmatrix}, \quad \text{(A.3)}$$

in which $\beta_s$ is the $s$-th variable of $\boldsymbol{\beta}$.

To obtain a sufficient condition to guarantee the convergence of the FP-MMKCC algorithm, we put forward two theorems below. For simplicity, we denote kernel bandwidths as $\sigma_i = \mu_i \sigma$, where $\mu_i$ is a positive constant.

*Theorem 1:* If $\delta > \xi = \dfrac{\vartheta\sqrt{m}}{\lambda_{\min}(\sum\limits_{j=1}^{N}\sum\limits_{i=1}^{m}\frac{\lambda_i}{\mu_i}\boldsymbol{h}_i\boldsymbol{h}_i^T) + \lambda_r}$ and

$\sigma \ge \sigma^*$, where $\vartheta = \sum\limits_{j=1}^{N}\sum\limits_{i=1}^{m}\frac{\lambda_i}{\mu_i^3}|t_j - c_i|\left\|\boldsymbol{h}_j\right\|_1$, $\lambda_{\min}[\cdot]$ denotes the minimum eigenvalue of the matrix term and $\sigma^*$ is the solution of equation $\varphi(\sigma) = \dfrac{\vartheta\sqrt{m}}{\lambda_{\min}(\theta) + \lambda_r} = \delta, \sigma \in (0, \infty)$ with $\theta = \sum\limits_{j=1}^{N}\sum\limits_{i=1}^{m}\frac{\lambda_i}{\mu_i^3}\exp\left(-\dfrac{(\delta\|\boldsymbol{h}_j\|_1 + |t_j - c_i|)^2}{2\mu_i^2\sigma^2}\right)\boldsymbol{h}_j^T\boldsymbol{h}_j$.
Then $\left\|\boldsymbol{f}(\boldsymbol{\beta})\right\|_1 \le \delta$ for all $\boldsymbol{\beta} \in \{\boldsymbol{\beta} \in \mathbb{R}^L : \left\|\boldsymbol{\beta}\right\|_1 \le \delta\}$.

*Proof:* The induced matrix norm is compatible with the corresponding vector $\ell_p$-norm, hence

$$\begin{aligned} \left\|\boldsymbol{f}(\boldsymbol{\beta})\right\|_1 &= \left\|[\mathbf{A} + \gamma'\mathbf{I}]^{-1}\mathbf{B}\right\|_1 \\ &\le \left\|[\mathbf{A} + \gamma'\mathbf{I}]^{-1}\right\|_1 \left\|\mathbf{B}\right\|_1, \end{aligned} \quad \text{(A.4)}$$

where $\|\cdot\|_1$ is the 1-norm (also referred to as the column-sum norm), which is simply the maximum absolute column sum of the matrix. According to the matrix theory, the following inequality holds:

$$\begin{aligned} \left\|[\mathbf{A} + \gamma'\mathbf{I}]^{-1}\right\|_1 &\le \sqrt{m}\left\|[\mathbf{A} + \gamma'\mathbf{I}]^{-1}\right\|_2 \\ &= \sqrt{m}\lambda_{\max}\big[[\mathbf{A} + \gamma'\mathbf{I}]^{-1}\big], \end{aligned} \quad \text{(A.5)}$$

where $\|\cdot\|_2$ is the 2-norm (also referred to as the spectral norm), which equals the maximum eigenvalue of the matrix denoted by $\lambda_{\max}[\cdot]$. Further, we have

$$\begin{aligned} &\lambda_{\max}\big[[\mathbf{A} + \gamma'\mathbf{I}]^{-1}\big] \\ &= \frac{1}{\lambda_{\min}[\mathbf{A} + \gamma'\mathbf{I}]} \\ &= \frac{1}{\lambda_{\min}\Big[\sum\limits_{j=1}^{N}\sum\limits_{i=1}^{m}\frac{\lambda_i}{\sigma_i^2}\kappa_{\sigma_i}(e_j - c_i)\boldsymbol{h}_j^T\boldsymbol{h}_j\Big] + \gamma'} \\ &\overset{(a)}{\le} \frac{1}{\frac{1}{\sigma^2}\lambda_{\min}\Big[\frac{\lambda_i}{\mu_i^2}\sum\limits_{j=1}^{N}\sum\limits_{i=1}^{m}\kappa_{\mu_i\sigma}(\delta\|\boldsymbol{h}_j\|_1 + |t_j - c_i|)\boldsymbol{h}_j^T\boldsymbol{h}_j\Big] + \gamma'} \\ &= \frac{\sigma^3\sqrt{2\pi}}{\lambda_{\min}(\theta) + \gamma''}, \end{aligned} \quad \text{(A.6)}$$

where $\gamma'' = \sigma^3\sqrt{2\pi}\gamma'$, and (a) comes from

$$\begin{aligned} |e_j - c_i| &= |t_j - \boldsymbol{\beta}^T\boldsymbol{h}_j - c_i| \\ &\le \|\boldsymbol{\beta}\|_1\|\boldsymbol{h}_j\|_1 + |t_j - c_i| \\ &\le \delta\|\boldsymbol{h}_j\|_1 + |t_j - c_i|. \end{aligned} \quad \text{(A.7)}$$

Likewise, we have

$$\begin{aligned} \left\|\mathbf{B}\right\|_1 &= \left\|\sum\limits_{j=1}^{N}\sum\limits_{i=1}^{m}\frac{\lambda_i}{\mu_i^2\sigma^2}\kappa_{\mu_i\sigma}(e_j - c_i)[t_j - c_i]\boldsymbol{h}_j^T\right\|_1 \\ &\overset{(b)}{\le} \frac{1}{\sigma^3\sqrt{2\pi}}\sum\limits_{j=1}^{N}\sum\limits_{i=1}^{m}\frac{\lambda_i}{\mu_i^3}|t_j - c_i|\left\|\boldsymbol{h}_j^T\right\|_1, \end{aligned} \quad \text{(A.8)}$$

where (b) is because $\kappa_\sigma(x) \le \dfrac{1}{\sigma\sqrt{2\pi}}$ for any $x$.

Combining (A.4)-(A.6) and (A.8), we have

$$\left\|\boldsymbol{f}(\boldsymbol{\beta})\right\|_1 \le \frac{\vartheta\sqrt{m}}{\lambda_{\min}(\theta) + \gamma''} = \varphi(\sigma). \quad \text{(A.9)}$$

Clearly, the function $\varphi(\sigma)$ is a continuous and monotonically decreasing function of $\sigma$ over $(0, \infty)$, satisfying $\lim\limits_{\sigma \to 0}\varphi(\sigma) = \infty$, and $\lim\limits_{\sigma \to \infty}\varphi(\sigma) = \xi$. Therefore, if $\delta > \xi$, the equation $\varphi(\sigma) = \delta$ will have a unique solution $\sigma^*$ over $(0, \infty)$, and if $\sigma > \sigma^*$, we have $\varphi(\sigma) \le \delta$, which completes the proof. ∎

*Theorem 2:* If $\delta > \xi$ and $\sigma \ge max\{\sigma^*, \sigma^\dagger\}$, where $\sigma^*$ is the solution of equation $\varphi(\sigma) = \delta$, and $\sigma^\dagger$ is the solution of equation $\psi(\sigma) = \alpha(0 < \alpha < 1)$, where $\psi(\sigma) = \dfrac{\sqrt{m}(\Theta)}{(\lambda_{\min}(\theta) + \gamma'')\sigma^2}, \sigma \in (0, \infty)$ with $\Theta = \sum\limits_{j=1}^{N}\sum\limits_{i=1}^{m}\frac{\lambda_i}{\mu_i^5}\Big(\delta\|\boldsymbol{h}_j\|_1 + |t_j - c_i|\Big)\|\boldsymbol{h}_j\|_1\Big(\delta\|\boldsymbol{h}_j^T\boldsymbol{h}_j^T\|_1 + |t_j|\|\boldsymbol{h}_j\|_1\Big)$. Then it holds that $\left\|\boldsymbol{f}(\boldsymbol{\beta})\right\|_1 \le \delta$ and $\left\|\nabla_{\boldsymbol{\beta}}\boldsymbol{f}(\boldsymbol{\beta})\right\|_1 \le \alpha$ for all $\boldsymbol{\beta} \in \{\boldsymbol{\beta} \in \mathbb{R}^L : \left\|\boldsymbol{\beta}\right\|_1 \le \delta\}$.

*Proof:* By Theorem 1, we have $\left\|\boldsymbol{f}(\boldsymbol{\beta})\right\|_1 \le \delta$. To prove $\left\|\nabla_{\boldsymbol{\beta}}\boldsymbol{f}(\boldsymbol{\beta})\right\|_1 \le \alpha$, it suffices to prove

$$\forall s, \left\|\frac{\partial}{\partial \beta_s}\boldsymbol{f}(\boldsymbol{\beta})\right\|_1 \le \alpha, \quad \text{(A.10)}$$

where

$$\begin{aligned} &\left\|\frac{\partial}{\partial \beta_s}\boldsymbol{f}(\boldsymbol{\beta})\right\|_1 \\ &= \left\|\frac{\partial}{\partial \beta_s}\big([\mathbf{A} + \lambda'I]^{-1}\mathbf{B}\big)\right\|_1 \\ &= \left\|-[\mathbf{A} + \lambda'I]^{-1}\Big(\frac{\partial}{\partial \beta_s}[\mathbf{A} + \lambda'I]\Big)[\mathbf{A} + \lambda'I]^{-1}\mathbf{B}\right. \\ &\quad \left. + [\mathbf{A} + \lambda'I]^{-1}\Big(\frac{\partial}{\partial \beta_s}\mathbf{B}\Big)\right\|_1 \\ &\le \left\|[\mathbf{A} + \lambda'I]^{-1}\right\|_1\left\|\frac{\partial}{\partial \beta_s}[\mathbf{A} + \lambda'I]\right\|_1\left\|\boldsymbol{f}(\boldsymbol{\beta})\right\|_1 \\ &\quad + \left\|[\mathbf{A} + \lambda'I]^{-1}\right\|_1\left\|\frac{\partial}{\partial \beta_s}\mathbf{B}\right\|_1. \end{aligned} \quad \text{(A.11)}$$

It is easy to derive

$$\left\|\frac{\partial}{\partial\beta_s}(\mathbf{A}+\lambda'I)\right\|_1$$

$$=\left\|\sum_{j=1}^{N}\sum_{i=1}^{m}\frac{\lambda_i}{\mu_i^4\sigma^4}(e_j-c_i)h_{js}\kappa_\sigma(e_j-c_i)\boldsymbol{h}_j\boldsymbol{h}_j^T\right\|_1$$

$$\overset{(c)}{\leq}\frac{1}{\sigma^5\sqrt{2\pi}}\sum_{j=1}^{N}\sum_{i=1}^{m}\frac{\lambda_i}{\mu_i^5}\Big(\delta\|\boldsymbol{h}_j\|_1+|t_j-c_i|\Big)\|\boldsymbol{h}_j\|_1\|\boldsymbol{h}_j^T\boldsymbol{h}_j\|_1,$$
(A.12)

where (c) is due to the fact that $|(e_j-c_i)h_{js}|\leq\Big(\delta\|\boldsymbol{h}_j\|_1+|t_j-c_i|\Big)\|\boldsymbol{h}_j\|_1$ and $\kappa_\sigma(x)\leq\frac{1}{\sigma\sqrt{2\pi}}$ for any $x$, in which $h_{is}$ is the $s$-th variable of $\boldsymbol{h}_i$. Similarly, one can derive

$$\left\|\frac{\partial}{\partial\beta_s}\mathbf{B}\right\|_1\leq\frac{1}{\sqrt{2\pi}\sigma^5}\sum_{j=1}^{N}\sum_{i=1}^{m}\frac{\lambda_i}{\mu_i^5}\Big(\delta\|\boldsymbol{h}_j\|_1+|t_j-c_i|\Big)\|\boldsymbol{h}_j\|_1\|t_j\boldsymbol{h}_j^T\|_1$$
(A.13)

Then, combining (A.5), (A.6), (A.11)-(A.13) and $\|\boldsymbol{f}(\boldsymbol{\beta})\|_1\leq\delta$, we have

$$\left\|\frac{\partial}{\partial\beta_s}\boldsymbol{f}(\boldsymbol{\beta})\right\|_1\leq\psi(\sigma).$$
(A.14)

Obviously, $\psi(\sigma)$ is also a continuous and monotonically decreasing function of $\sigma$ over $(0,\infty)$, and satisfies $\lim_{\sigma\to0}\psi(\sigma)=\infty$, and $\lim_{\sigma\to\infty}\psi(\sigma)=0$. Therefore, given $0<\alpha<1$, the equation $\psi(\sigma)=\alpha$ has a unique solution $\sigma^\dagger$ over $(0,\infty)$, and if $\sigma>\sigma^\dagger$, we have $\psi(\sigma)\leq\delta$, which completes the proof. ∎

According to Theorem 2 and contraction mapping theorem, given a certain parameter $\mu_i$ and an initial weight vector satisfying $\|\boldsymbol{\beta}_0\|_1\leq\delta$, the FP-MMKCC algorithm will surely converge to a unique fixed point in the range $\boldsymbol{\beta}\in\{\boldsymbol{\beta}\in\mathbb{R}^L:\|\boldsymbol{\beta}\|_1\leq\delta\}$ provided that the kernel bandwidth $\sigma$ is larger than a certain value. Moreover, since the $\alpha(0<\alpha<1)$ is the Lipschitz constant in the contraction mapping theorem, its value guarantees the convergence speed.

## REFERENCES

[1] J.-H. Lin, T. M. Sellke, and E. J. Coyle, "Adaptive stack filtering under the mean absolute error criterion," *IEEE transactions on acoustics, speech, and signal processing*, vol. 38, no. 6, pp. 938–954, 1990.

[2] E. J. Coyle and J.-H. Lin, "Stack filters and the mean absolute error criterion," *IEEE Transactions on Acoustics, Speech, and Signal Processing*, vol. 36, no. 8, pp. 1244–1254, 1988.

[3] S.-C. Pei and C.-C. Tseng, "Least mean p-power error criterion for adaptive fir filter," *IEEE Journal on Selected Areas in Communications*, vol. 12, no. 9, pp. 1540–1547, 1994.

[4] W. Liu, P. Pokharel, and J. Principe, "Error entropy, correntropy and m-estimation," in *2006 16th IEEE Signal Processing Society Workshop on Machine Learning for Signal Processing*. IEEE, 2006, pp. 179–184.

[5] M. O. Sayin, N. D. Vanli, and S. S. Kozat, "A novel family of adaptive filtering algorithms based on the log-arithmic cost," *IEEE Transactions on signal processing*, vol. 62, no. 17, pp. 4411–4424, 2014.

[6] X. Chen, J. Yang, J. Liang, and Q. Ye, "Recursive robust least squares support vector regression based on maximum correntropy criterion," *Neurocomputing*, vol. 97, pp. 63–73, 2012.

[7] Y. Feng, X. Huang, L. Shi, Y. Yang, and J. A. Suykens, "Learning with the maximum correntropy criterion induced losses for regression." *Journal of Machine Learning Research*, vol. 16, pp. 993–1034, 2015.

[8] G. Xu, B.-G. Hu, and J. C. Principe, "Robust c-loss kernel classifiers," *IEEE transactions on neural networks and learning systems*, vol. 29, no. 3, pp. 510–522, 2018.

[9] M. N. Syed, J. C. Principe, and P. M. Pardalos, "Correntropy in data classification," in *Dynamics of Information Systems: Mathematical Foundations*. Springer, 2012, pp. 81–117.

[10] A. Singh, R. Pokharel, and J. Principe, "The c-loss function for pattern classification," *Pattern Recognition*, vol. 47, no. 1, pp. 441–453, 2014.

[11] W. Shi, Y. Gong, X. Tao, and N. Zheng, "Training dcnn by combining max-margin, max-correlation objectives, and correntropy loss for multilabel image classification," *IEEE transactions on neural networks and learning systems*, vol. 29, no. 7, pp. 2896–2908, 2018.

[12] L.-R. Ren, Y.-L. Gao, J.-X. Liu, J. Shang, and C.-H. Zheng, "Correntropy induced loss based sparse robust graph regularized extreme learning machine for cancer classification," *BMC bioinformatics*, vol. 21, no. 1, pp. 1–22, 2020.

[13] R. He, B.-G. Hu, W.-S. Zheng, and X.-W. Kong, "Robust principal component analysis based on maximum correntropy criterion," *IEEE Transactions on Image Processing*, vol. 20, no. 6, pp. 1485–1494, 2011.

[14] N. Zhou, Y. Xu, H. Cheng, Z. Yuan, and B. Chen, "Maximum correntropy criterion based sparse subspace learning for unsupervised feature selection," *IEEE Transactions on Circuits and Systems for Video Technology*, 2017.

[15] N. Yu, M.-J. Wu, J.-X. Liu, C.-H. Zheng, and Y. Xu, "Correntropy-based hypergraph regularized nmf for clustering and feature selection on multi-cancer integrated data," *IEEE Transactions on Cybernetics*, 2020.

[16] S. Zhao, B. Chen, and J. C. Principe, "Kernel adaptive filtering with maximum correntropy criterion," in *The 2011 International Joint Conference on Neural Networks*. IEEE, 2011, pp. 2012–2017.

[17] B. Chen, L. Xing, H. Zhao, N. Zheng, and J. C. Principe, "Generalized correntropy adaptive filtering," *IEEE Transactions on Signal Processing*, vol. 64, no. 13, pp. 3376–3387, 2016.

[18] W. Ma, H. Qu, G. Gui, L. Xu, J. Zhao, and B. Chen, "Maximum correntropy criterion based sparse adaptive filtering algorithms for robust channel estimation under non-gaussian environments," *Journal of the Franklin Institute*, vol. 352, no. 7, pp. 2708–2727, 2015.

[19] Z. Wu, S. Peng, B. Chen, and H. Zhao, "Robust hammerstein adaptive filtering under maximum correntropy criterion," *Entropy*, vol. 17, no. 10, pp. 7149–7166, 2015.

[20] Z. Wu, J. Shi, X. Zhang, W. Ma, and B. Chen, "Kernel recursive maximum correntropy," *Signal Processing*, vol. 117, pp. 11–16, 2015.

[21] G. T. Cinar and J. C. Príncipe, "Hidden state estimation using the correntropy filter with fixed point update and adaptive kernel size," in *Neural Networks (IJCNN), The 2012 International Joint Conference on.* IEEE, 2012, pp. 1–6.

[22] X. Liu, B. Chen, B. Xu, Z. Wu, and P. Honeine, "Maximum correntropy unscented filter," *International Journal of Systems Science*, vol. 48, no. 8, pp. 1607–1615, 2017.

[23] W. Liu, P. P. Pokharel, and J. C. Príncipe, "Correntropy: Properties and applications in non-gaussian signal processing," *IEEE Transactions on Signal Processing*, vol. 55, no. 11, pp. 5286–5298, 2007.

[24] B. Chen, X. Wang, N. Lu, S. Wang, J. Cao, and J. Qin, "Mixture correntropy for robust learning," *Pattern Recognition*, vol. 79, pp. 318–327, 2018.

[25] M. Gönen and E. Alpaydın, "Multiple kernel learning algorithms," *Journal of machine learning research*, vol. 12, no. Jul, pp. 2211–2268, 2011.

[26] G. R. Lanckriet, N. Cristianini, P. Bartlett, L. E. Ghaoui, and M. I. Jordan, "Learning the kernel matrix with semidefinite programming," *Journal of Machine learning research*, vol. 5, no. Jan, pp. 27–72, 2004.

[27] Z. Wang, S. Chen, and T. Sun, "Multik-mhks: a novel multiple kernel learning algorithm," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 30, no. 2, pp. 348–353, 2007.

[28] M. Yukawa, "Multikernel adaptive filtering," *IEEE Transactions on Signal Processing*, vol. 60, no. 9, pp. 4672–4682, 2012.

[29] B. Chen, X. Liu, H. Zhao, and J. C. Principe, "Maximum correntropy kalman filter," *Automatica*, vol. 76, pp. 70–77, 2017.

[30] F. Wang, Y. He, S. Wang, and B. Chen, "Maximum total correntropy adaptive filtering against heavy-tailed noises," *Signal Processing*, vol. 141, pp. 84–95, 2017.

[31] J. C. Principe, *Information theoretic learning: Renyi's entropy and kernel perspectives.* Springer Science & Business Media, 2010.

[32] I. Santamaría, C. Pantaleón, L. Vielva, and J. C. Principe, "Adaptive blind equalization through quadratic pdf matching," in *2002 11th European Signal Processing Conference.* IEEE, 2002, pp. 1–4.

[33] A. R. Heravi and G. A. Hodtani, "A new information theoretic relation between minimum error entropy and maximum correntropy," *IEEE Signal Processing Letters*, vol. 25, no. 7, pp. 921–925, 2018.

[34] A. Frank and A. Asuncion, "Uci machine learning repository [http://archive. ics. uci. edu/ml]. irvine, ca: University of california," *School of information and computer science*.

[35] D. Wang and M. Li, "Stochastic configuration networks: Fundamentals and algorithms," *IEEE transactions on cybernetics*, vol. 47, no. 10, pp. 3466–3479, 2017.

[36] ——, "Robust stochastic configuration networks with kernel density estimation for uncertain data regression," *Information Sciences*, vol. 412, pp. 210–222, 2017.

[37] M. Li, C. Huang, and D. Wang, "Robust stochastic configuration networks with maximum correntropy criterion for uncertain data regression," *Information Sciences*, vol. 473, pp. 73–86, 2019.

[38] A. Waibel, T. Hanazawa, G. Hinton, K. Shikano, and K. J. Lang, "Phoneme recognition using time-delay neural networks," *Backpropagation: Theory, Architectures and Applications*, pp. 35–61, 1995.

[39] J.-M. Kuo, "Nonlinear dynamic modeling with artificial neural networks," Ph.D. dissertation, Citeseer, 1993.

[40] F. Lotte, L. Bougrain, A. Cichocki, M. Clerc, M. Congedo, A. Rakotomamonjy, and F. Yger, "A review of classification algorithms for eeg-based brain–computer interfaces: a 10 year update," *Journal of neural engineering*, vol. 15, no. 3, p. 031005, 2018.

[41] B. He, H. Yuan, J. Meng, and S. Gao, "Brain–computer interfaces," in *Neural engineering.* Springer, 2020, pp. 131–183.

[42] M. K. Islam, A. Rastegarnia, and Z. Yang, "Methods for artifact detection and removal from scalp eeg: A review," *Neurophysiologie Clinique/Clinical Neurophysiology*, vol. 46, no. 4-5, pp. 287–305, 2016.

[43] M. Tangermann, K.-R. Müller, A. Aertsen, N. Birbaumer, C. Braun, C. Brunner, R. Leeb, C. Mehring, K. J. Miller, G. Mueller-Putz *et al.*, "Review of the bci competition iv," *Frontiers in neuroscience*, vol. 6, p. 55, 2012.

[44] F. Lotte and C. Guan, "Regularizing common spatial patterns to improve bci designs: unified theory and new algorithms," *IEEE Transactions on biomedical Engineering*, vol. 58, no. 2, pp. 355–362, 2010.

[45] A. Rocha and S. K. Goldenstein, "Multiclass from binary: Expanding one-versus-all, one-versus-one and ecoc-based approaches," *IEEE Transactions on Neural Networks and Learning Systems*, vol. 25, no. 2, pp. 289–302, 2013.

[46] S. Sun and R. Huang, "An adaptive k-nearest neighbor algorithm," in *2010 seventh international conference on fuzzy systems and knowledge discovery*, vol. 1. IEEE, 2010, pp. 91–94.

[47] J. A. Suykens and J. Vandewalle, "Least squares support vector machine classifiers," *Neural processing letters*, vol. 9, no. 3, pp. 293–300, 1999.

[48] Y. Zhang, B. Chen, X. Liu, Z. Yuan, and J. Principe, "Convergence of a fixed-point minimum error entropy algorithm," *Entropy*, vol. 17, no. 8, pp. 5549–5560, 2015.

[49] B. Chen, J. Wang, H. Zhao, N. Zheng, and J. C. Principe, "Convergence of a fixed-point algorithm under maximum correntropy criterion," *IEEE Signal Processing Letters*, vol. 22, no. 10, pp. 1723–1727, 2015.