

Constrained Maximum Correntropy Adaptive Filtering

Siyuan Peng, *Student Member, IEEE*, Badong Chen, *Senior Member, IEEE*, Lei Sun, *Senior Member, IEEE*,
Wee Ser, *Senior Member, IEEE*, and Zhiping Lin, *Senior Member, IEEE*

Abstract—Constrained adaptive filtering algorithms including constrained least mean square (CLMS), constrained affine projection (CAP) and constrained recursive least squares (CRLS) have been extensively studied in many applications. Most existing constrained adaptive filtering algorithms are developed under the mean square error (MSE) criterion, which is an ideal optimality criterion under Gaussian noises. This assumption however fails to model the behavior of non-Gaussian noises found in practice. Motivated by the robustness and simplicity of maximum correntropy criterion (MCC) for non-Gaussian impulsive noises, this paper proposes a new adaptive filtering algorithm called constrained maximum correntropy criterion (CMCC). Specifically, CMCC incorporates a linear constraint into a MCC filter to solve a constrained optimization problem explicitly. The proposed adaptive filtering algorithm is easy to implement and has low computational complexity, and in terms of convergence accuracy (say lower mean square deviation) and stability, it can significantly outperform those MSE based constrained adaptive algorithms in presence of heavy-tailed impulsive noises. Additionally, the mean square convergence behaviors are studied under energy conservation relation, and a sufficient condition to ensure the mean square convergence and the steady-state mean square deviation (MSD) of the proposed algorithm are obtained. Simulation results confirm the theoretical predictions under both Gaussian and non-Gaussian noises, and demonstrate the excellent performance of the novel algorithm by comparing it with other conventional methods.

Index Terms—adaptive filtering, constrained maximum correntropy criterion, non-Gaussian signal processing, convergence analysis.

I. INTRODUCTION

CONSTRAINED adaptive filtering algorithms have been successfully applied in domains of signal processing and communications, such as system identification, blind interference suppression, array signal processing, and spectral analysis [1]–[4]. The main advantage of constrained adaptive filters is that they have an error-correcting feature that can prevent the accumulation of errors (e.g. the quantization errors in a digital implementation). As a well-known linearly-constrained adaptive filtering algorithm, the *constrained least mean square* (CLMS) [5] is a simple stochastic-gradient based adaptive algorithm, originally conceived as an adaptive solution to

a linearly-constrained minimum-variance (LCMV) filtering problem in antenna array processing [6]. Other stochastic-gradient based linearly-constrained adaptive algorithms were also developed [7]–[11]. Although the LMS-type algorithms are simple and computationally efficient, they may suffer from low convergence speed especially when the input signal is correlated. In order to improve the convergence rate, the *constrained recursive least squares* (CRLS) algorithm was derived in [12], at the cost of higher computational complexity. Some improvements of the CRLS can be found in [13], [14]. Several *constrained affine projection* (CAP) algorithms were also developed [15], [16].

Most of the existing constrained adaptive filtering algorithms have been developed based on the common mean square error (MSE) criterion due to its attractive features, such as mathematical tractability, computational simplicity and optimality under Gaussian assumption [17], [18]. However, Gaussian assumption does not always hold in real-world environments, even though it is justified for many natural noises. When the signals are disturbed by non-Gaussian noises, the MSE based algorithms may perform poorly or encounter the instability problem [19], [23]. From a statistical viewpoint, the MSE is insufficient to capture all possible information in non-Gaussian signals. In practical situations, non-Gaussian noises are frequently encountered. For example, some sources of non-Gaussian impulsive noises are ill synchronization in digital recording, motor ignition noise in internal combustion engines, scratches on recording disks, and lighting spikes in natural phenomena [20]–[22].

To deal with the non-Gaussian noise problem (which usually causes large outliers), various alternative optimization criteria have been proposed to replace the MSE criterion for developing robust adaptive filtering algorithms in the literature. In recent years, maximum correntropy criterion (MCC) has been successfully applied in diverse domains due to its simplicity and robustness [19], [23]–[29]. As a nonlinear and local similarity measure directly related to the probability of how similar two random variables are in the bisector neighborhood of the joint space controlled by the kernel bandwidth, correntropy is insensitive to large outliers, and is frequently used as a powerful method to handle non-Gaussian impulsive noises in various applications of engineering. For instance, Singh et al. [30] and Zhao et al. [25] utilized the correntropy as a cost function to develop robust adaptive filtering algorithm for signal processing, and Chen et al. extended the original correntropy by using the generalized Gaussian density (GGD) function as the kernel, and proposed

S. Peng, Z. Lin and W. Ser are with the School of Electrical and Electronic Engineering, Nanyang Technological University, 639798 Singapore e-mail: (PENG0074@e.ntu.edu.sg, EZPLin@ntu.edu.sg, ewser@ntu.edu.sg).

B. Chen is with the Institute of Artificial Intelligence and Robotics, Xi'an Jiaotong University, Xi'an 710049 China e-mail: (Corresponding author, chenbd@mail.xjtu.edu.cn).

L. Sun is with the School of Information and Electronics, Beijing Institute of Technology, Beijing 100081, China e-mail: (bitsunlei@126.com).

Manuscript received Dec. 6, 2016.

a generalized correntropy for robust adaptive filtering [31]. He et al. presented a MCC-based rotationally invariant principal component analysis (PCA) algorithm for image processing [32], and also incorporated the correntropy induced metric (CIM) into MCC to develop an effective sparse representation algorithm for robust face recognition [33]. Bessa et al. adopted MCC to train neural networks for wind prediction in power system [34]. Hasanbelliu et al. utilized information theoretic measures (entropy and correntropy) to develop two algorithms that can deal with both rigid and non-rigid point set registration with different computational complexities and accuracies [35]. However, constrained adaptive filtering based on MCC has not been studied yet in the literature. In this work, a constrained maximum correntropy criterion (CMCC) adaptive filtering algorithm is proposed for signal processing especially in presence of heavy-tailed impulsive noises.

Our main contributions in this paper are summarized as follows:

- First, we develop the CMCC adaptive filtering algorithm by incorporating a linear constraint into the MCC, instead of the traditional MSE criterion, to solve a constrained optimization problem explicitly. The computational complexity analysis is also presented.
- Second, based on the energy conservation relation [36]–[39], we analyze the mean square convergence behaviors of the proposed algorithm, and present particularly a sufficient condition to guarantee the mean square convergence and the steady-state mean square deviation (MSD) in the cases of Gaussian and non-Gaussian noises.
- Finally, we confirm the validity of theoretical expectations experimentally, and illustrate the desirable performance (e.g. lower MSD) of CMCC by comparing it with other methods in linear-phase system identification and beamforming application.

The rest of the paper is organized as follows. In Section II, after briefly reviewing the MCC, we develop the CMCC algorithm and analyze the computational complexity. In Section III, we study the mean square convergence of the proposed algorithm. Simulation results are then presented in Section IV. Finally, Section V gives the conclusion and discusses some work in the future. Some derivations are relegated to the Appendix.

II. CMCC ALGORITHM

A. Maximum Correntropy Criterion

As a similarity measure between two random variables X and Y , correntropy is defined by [23], [27], [29]–[31]

$$V(X, Y) = E[\kappa(X, Y)] = \int \kappa(x, y) dF_{XY}(x, y) \quad (1)$$

where $E[\cdot]$ denotes the *expectation operator*, $\kappa(\cdot, \cdot)$ is a *shift-invariant Mercer kernel*, and $F_{XY}(x, y)$ stands for the joint distribution function of (X, Y) . It takes the advantage of a kernel trick that nonlinearly maps the input space to a higher dimensional feature space. In the present work, without mentioning otherwise, the kernel function of correntropy $\kappa(\cdot, \cdot)$

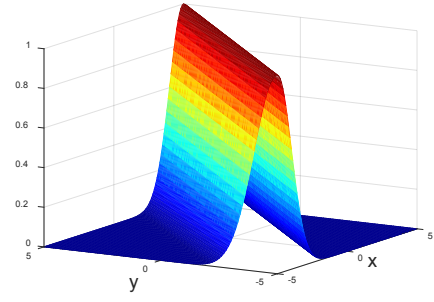


Fig. 1: MCC cost function in the joint space ($\sigma = 1.0$).

is the Gaussian kernel, given by

$$\kappa_{\sigma}(x - y) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(x - y)^2}{2\sigma^2}\right) \quad (2)$$

where $\sigma > 0$ is the kernel bandwidth parameter. In most practical situations, the joint distribution $F_{XY}(x, y)$ is usually unknown, and only a finite number of data samples $\{(x(n), y(n))\}_{n=1}^N$ are available. In these cases, the correntropy can be estimated by

$$\hat{V}_{N,\sigma} = \frac{1}{N} \sum_{n=1}^N \kappa_{\sigma}(x(n) - y(n)) \quad (3)$$

Under the *maximum correntropy criterion* (MCC), an adaptive filter will be trained by maximizing the correntropy between the desired response and filter output, formulated by

$$\max_W J_{MCC} = \frac{1}{N} \sum_{n=1}^N \kappa_{\sigma}(e(n)) \quad (4)$$

where $e(n)$ denotes the error between the desired response and filter output, and W stands for the filter weight vector. Fig. 1 shows the MCC cost function $\kappa_{\sigma}(x - y)$ in the joint space of x and y . As one can see clearly, the MCC is a local similarity measure, whose value is mainly decided by the kernel function along the line $x = y$. Furthermore, from a view of geometric meaning, we can divide the space in three regions, namely Euclidean region, transition region and rectification region. The MCC behaves like 2-norm distance in the Euclidean region, similarly like a 1-norm distance in the transition region and eventually approaches a zero-norm in the rectification region, which also interprets the robustness of correntropy for outliers [23], [29].

B. CMCC Algorithm

Consider a linear unknown system, with an M -dimensional weight vector $W^* = [w_1^*(n), w_2^*(n), \dots, w_M^*(n)]^T$ that needs to be estimated. The measured output $d(n)$ of the unknown system at instant n is assumed to be

$$d(n) = y^*(n) + v(n) = W^{*T}X(n) + v(n) \quad (5)$$

where $y^*(n) = W^{*T}X(n)$ denotes the actual output of the unknown system, $X(n) = [x_1(n), x_2(n), \dots, x_M(n)]^T$ is the input vector, with $[\cdot]^T$ being the transpose operator, and $v(n)$ stands for an interference or measurement noise.

Suppose the estimator is another M -dimensional linear filter, with an adaptive weight vector $W(n)$. Then the instantaneous prediction error at instant n is

$$e(n) = d(n) - y(n) = d(n) - W^T(n-1)X(n) \quad (6)$$

where $y(n) = W^T(n-1)X(n)$ denotes the output of the adaptive filter. For a constrained adaptive filter, a linear constraint will be imposed upon the filter weight vector as

$$C^T W = f \quad (7)$$

where C is an $M \times K$ constraint matrix, and f is a vector containing K constraint values. The CLMS algorithm is derived by solving the following optimization problem [5], [37]:

$$\begin{aligned} \min_W & E \left[(d(n) - W^T(n-1)X(n))^2 \right] \\ \text{subject to} & C^T W = f \end{aligned} \quad (8)$$

leading to the following weight update equation:

$$\begin{aligned} W(n) = & P[W(n-1) + \eta(d(n) - \\ & W^T(n-1)X(n))X(n)] + Q \end{aligned} \quad (9)$$

where $\eta > 0$ is the step-size parameter, $P = \mathbf{I}_M - C(C^T C)^{-1}C^T$ with \mathbf{I}_M being an $M \times M$ identity matrix, and $Q = C(C^T C)^{-1}f$.

In this work, we use the MCC instead of MSE to develop a constrained adaptive filtering algorithm. Similar to (8), we propose the following CMCC optimization problem

$$\begin{aligned} \max_W & E[\kappa_\sigma(d(n) - W^T(n-1)X(n))] \\ \text{subject to} & C^T W = f \end{aligned} \quad (10)$$

and accordingly the CMCC cost J_{CMCC} is

$$\begin{aligned} J_{CMCC} = & E[\kappa_\sigma(d(n) - W^T(n-1)X(n))] + \\ & \xi^T(n) (C^T W(n-1) - f) \end{aligned} \quad (11)$$

where $\xi(n)$ is a $K \times 1$ Lagrange multiplier vector. A stochastic-gradient based algorithm can thus be derived as (see Appendix A for a detailed derivation)

$$\begin{aligned} W(n) = & P[W(n-1) + \eta g(e(n))(d(n) - \\ & W^T(n-1)X(n))X(n)] + Q \end{aligned} \quad (12)$$

where $g(e(n))$ is a nonlinear function of $e(n)$, given by

$$g(e(n)) = \exp\left(-\frac{e^2(n)}{2\sigma^2}\right) \quad (13)$$

The above algorithm is referred to as the CMCC algorithm, whose pseudocodes are presented in Table I.

C. Computational Complexity

The computational complexity of the proposed CMCC algorithm and other constrained adaptive algorithms-CLMS, CAP and CRLS, in terms of the total number of required additions and multiplications at each iteration, are shown in Table II, where Γ_g is a constant associated with the complexity of the nonlinear function $g(e(n))$. Obviously, the computational complexity of these algorithms are $O(M^2)$. When Γ_g is small,

TABLE I: CMCC Algorithm

Parameters:	η, σ, C and f
Initialization :	$P = \mathbf{I}_M - C(C^T C)^{-1}C^T$ $Q = C(C^T C)^{-1}f$ $W(0) = Q$
Update:	$y(n) = W^T(n-1)X(n)$ $e(n) = d(n) - y(n)$ $W(n) = P[W(n-1) + \eta g(e(n))e(n)X(n)] + Q$

TABLE II: Computational Complexity of CMCC, CLMS, CAP and CRLS

Algorithm	Computational Complexity
CMCC	$2M^2 + 5M + \Gamma_g$
CLMS	$2M^2 + 5M + 1$
CAP	$2M^2 + (2L + 3)M + 1$
CRLS	$7M^2 + (6K^2 + 9K + 5)M + 3K$

it can be seen that the proposed algorithm has lower computational cost than CRLS due to calculating the covariance matrix \mathbf{R} per iteration for CRLS, also has lower computational cost than CAP (especially when the sliding window length L is large). Generally speaking, the computational complexity of CMCC is almost the same as that of the CLMS.

III. CONVERGENCE ANALYSIS

In this section, we analyze the mean square convergence behaviors of the proposed CMCC algorithm. First, we give the following assumptions:

- 1) The input sequence $\{X(n)\}$ is independent multivariate Gaussian, with zero-mean and the positive-definite covariance matrix of the input sequence $\mathbf{R} = E[X(n)X^T(n)]$.
- 2) The noise $\{v(n)\}$ is zero-mean, independent, identically distributed, and independent of any other signals in the system.
- 3) The error nonlinearity $g(e(n))$ is asymptotically uncorrelated with $\{X(n)X^T(n)\}$ at steady-state.
- 4) The filter is long enough such that the a priori error $e_a(n) = (W^* - W(n))^T X(n)$ is zero-mean Gaussian.

The independence assumptions 1) and 2) are common in the literature of adaptive filtering [36]–[39]. When the filter is long enough, assumption 3) will become realistic and valid. Assumption 4) is reasonable by the central limit theorem, and also remains valid in the whole stage of adaptation (see [18], [38] for more detailed explanation about assumptions 3) and 4)).

A. Mean Square Stability

Let us define the weight error vector:

$$\tilde{W}(n) = W(n) - W_{opt} \quad (14)$$

where W_{opt} stands for the optimal solution of the CMCC optimization problem under the above assumptions, given by (see Appendix B for the detailed derivation)

$$W_{opt} = W^* + \mathbf{R}_g^{-1} C (C^T \mathbf{R}_g^{-1} C)^{-1} (f - C^T W^*) \quad (15)$$

where $\mathbf{R}_g = E[g(e(n))X(n)X^T(n)]$ denotes a weighted autocorrelation matrix of the input vector. We also define

$$\varepsilon_w = W^* - W_{opt} \quad (16)$$

Substituting (5), (14) and (16) into (12) yields

$$\begin{aligned} \tilde{W}(n) &= P[W(n-1) + \eta g(e(n))(d(n) - \\ &\quad W^T(n-1)X(n))X(n)] + Q - W_{opt} \\ &= P[\mathbf{I}_M - \eta g(e(n))X(n)X^T(n)] \tilde{W}(n-1) + \\ &\quad \eta g(e(n))v(n)PX(n) + \eta g(e(n))PX(n) \times \\ &\quad X^T(n)\varepsilon_w + PW_{opt} - W_{opt} + Q \end{aligned} \quad (17)$$

Due to $PW_{opt} - W_{opt} + Q = \mathbf{0}_{M \times 1}$ (Here $\mathbf{0}_{M \times 1}$ denotes the $M \times 1$ zero vector), we can rewrite (17) as

$$\begin{aligned} \tilde{W}(n) &= P[\mathbf{I}_M - \eta g(e(n))X(n)X^T(n)] \times \\ &\quad \tilde{W}(n-1) + \eta g(e(n))v(n)PX(n) + \\ &\quad \eta g(e(n))PX(n)X^T(n)\varepsilon_w \end{aligned} \quad (18)$$

Note that matrix P is idempotent, namely $P = P^2$ and $P = P^T$. Multiplying both sides of (18) by P and after some straightforward matrix manipulations, we can obtain

$$P\tilde{W}(n) = \tilde{W}(n) \quad (19)$$

Combining (18) and (19), we have

$$\begin{aligned} \tilde{W}(n) &= P\tilde{W}(n-1) - \eta g(e(n))PX(n)X^T(n)\tilde{W}(n-1) + \\ &\quad \eta g(e(n))v(n)PX(n) + \eta g(e(n))PX(n)X^T(n)\varepsilon_w \\ &= (\mathbf{I}_M - \eta g(e(n))PX(n)X^T(n)P) \tilde{W}(n-1) + \\ &\quad \eta g(e(n))v(n)PX(n) + \eta g(e(n))PX(n)X^T(n)\varepsilon_w \end{aligned} \quad (20)$$

Under assumptions 1), 2) and 3), taking the expectations of the squared Euclidean norms of both sides of (20) leads to the following energy conservation relation:

$$\begin{aligned} E[\|\tilde{W}(n)\|^2] &= E[\|\tilde{W}(n-1)\|_{\mathbf{H}}^2] + \eta^2 E[g^2(e(n))] \times \\ &\quad E[v^2(n)] E[X^T(n)PX(n)] + \eta^2 E[g^2(e(n))] \\ &\quad \times \varepsilon_w^T E[X(n)X^T(n)PX(n)X^T(n)] \varepsilon_w \end{aligned} \quad (21)$$

where $E[\|\tilde{W}(n)\|^2]$ is called the *weight error power (WEP)* at iteration n , $\|\tilde{W}(n-1)\|_{\mathbf{H}}^2 = \tilde{W}^T(n-1)\mathbf{H}\tilde{W}(n-1)$, and

$$\begin{aligned} \mathbf{H} &= \mathbf{I}_M - 2\eta E[g(e(n))] \mathbf{P} \mathbf{R} \mathbf{P} + \eta^2 E[g^2(e(n))] \times \\ &\quad PE[X(n)X^T(n)PX(n)X^T(n)]P \end{aligned}$$

Since $P = P^2$, we derive

$$\begin{aligned} E[X^T(n)PX(n)] &= E[X^T(n)PPX(n)] \\ &= \text{tr}\{\mathbf{P} \mathbf{R} \mathbf{P}\} \\ &= \text{tr}\{\Upsilon\} \end{aligned} \quad (22)$$

where $\text{tr}\{\cdot\}$ stands for the *trace operator*, and $\Upsilon = \mathbf{P} \mathbf{R} \mathbf{P}$. According to the Isserlis' theorem [40] for Gaussian vectors $\tilde{h}_1, \tilde{h}_2, \tilde{h}_3$ and \tilde{h}_4 , we have

$$\begin{aligned} E[\tilde{h}_1 \tilde{h}_2^T \tilde{h}_3 \tilde{h}_4^T] &= E[\tilde{h}_1 \tilde{h}_2^T] E[\tilde{h}_3 \tilde{h}_4^T] + E[\tilde{h}_1 \tilde{h}_3^T] E[\tilde{h}_2 \tilde{h}_4^T] \\ &\quad + E[\tilde{h}_1 \tilde{h}_4^T] E[\tilde{h}_2 \tilde{h}_3^T] \end{aligned} \quad (23)$$

With $\tilde{h}_1 = X(n)$, $\tilde{h}_2 = X(n)$, $\tilde{h}_3 = PX(n)$ and $\tilde{h}_4 = X(n)$, we obtain

$$\begin{aligned} E[X(n)X^T(n)PX(n)X^T(n)] &= \mathbf{P} \mathbf{R} \mathbf{P} + \mathbf{P} \mathbf{R} \mathbf{P} + E[X^T(n)PX(n)] \mathbf{R} \\ &= \text{tr}\{\Upsilon\} \mathbf{R} + 2\mathbf{P} \mathbf{R} \mathbf{P} \end{aligned} \quad (24)$$

Since $\mathbf{P} \mathbf{R} \varepsilon_w = \mathbf{0}_{M \times 1}$, Substituting (22) and (24) into (21), we get

$$\begin{aligned} E[\|\tilde{W}(n)\|^2] &= E[\|\tilde{W}(n-1)\|_{\mathbf{H}}^2] + \eta^2 E[g^2(e(n))] \times \\ &\quad \text{tr}\{\Upsilon\} (\varepsilon_w^T \mathbf{R} \varepsilon_w + E[v^2(n)]) \end{aligned} \quad (25)$$

and

$$\begin{aligned} \mathbf{H} &= \mathbf{I}_M - 2\eta E[g(e(n))] \mathbf{P} \mathbf{R} \mathbf{P} + \eta^2 E[g^2(e(n))] \times \\ &\quad (\text{tr}\{\Upsilon\} \mathbf{P} \mathbf{R} \mathbf{P} + 2\mathbf{P} \mathbf{R} \mathbf{P} \mathbf{R} \mathbf{P}) \end{aligned}$$

Let λ_i ($i = 1, \dots, M - K$) be the eigenvalues of the matrix Υ . A sufficient condition for the mean square stability can be obtained as [5], [24], [37]

$$\begin{aligned} |1 - 2\eta E[g(e(n))] \lambda_i + \eta^2 E[g^2(e(n))] \text{tr}\{\Upsilon\} \lambda_i + \\ 2\eta^2 E[g^2(e(n))] \lambda_i^2| < 1 \\ i = 1, \dots, M - K \end{aligned} \quad (26)$$

After some simple manipulations, we have

$$0 < \eta < \frac{2E[g(e(n))]}{[2\lambda_{max} + \text{tr}\{\Upsilon\}] E[g^2(e(n))]} \quad (27)$$

where λ_{max} denotes the largest eigenvalue of the matrix Υ . Due to $E[g(e(n))] \geq E[g^2(e(n))] > 0$, one can obtain a stronger condition to guarantee the mean square stability:

$$0 < \eta \leq \frac{2}{2\lambda_{max} + \text{tr}\{\Upsilon\}} \quad (28)$$

Remark: Since we only derive (27) and (28) under the steady-state assumption, we cannot solve the problem of how to select the best step-size for a specific application. However, the condition provides a possible range for choosing a step-size for CMCC algorithm.

B. Steady-state mean square deviation (MSD)

Assume that \mathbf{T} is an arbitrary symmetric nonnegative definite matrix. Under assumptions 1), 2) and 3), one can derive the following relation by taking the expectations of the squared-weighted Euclidean norms of both sides of (20):

$$\begin{aligned} E[\|\tilde{W}(n)\|_{\mathbf{T}}^2] &= E[\|\tilde{W}(n-1)\|_{\mathbf{T}}^2] + \eta^2 E[g^2(e(n))] \times \\ &\quad E[v^2(n)] E[X^T(n) \mathbf{T} P X(n)] + \eta^2 \times \\ &\quad E[g^2(e(n))] \varepsilon_w^T E[X(n)X^T(n) \times \\ &\quad \mathbf{T} P X(n)X^T(n)] \varepsilon_w \end{aligned} \quad (29)$$

in which

$$\begin{aligned} \mathbf{U} &= E \left[(\mathbf{I}_M - \eta g(e(n))X(n)X^T(n)) \mathbf{P} \mathbf{T} \mathbf{P} \times \right. \\ &\quad \left. (\mathbf{I}_M - \eta g(e(n))X(n)X^T(n)) \right] \\ &= \mathbf{P} \mathbf{T} \mathbf{P} - \eta E[g(e(n))] \mathbf{R} \mathbf{P} \mathbf{T} \mathbf{P} - \eta \times \\ &\quad E[g(e(n))] \mathbf{P} \mathbf{T} \mathbf{P} \mathbf{R} + \eta^2 E[g^2(e(n))] \times \\ &\quad E[X(n)X^T(n) \mathbf{P} \mathbf{T} \mathbf{P} X(n)X^T(n)] \end{aligned} \quad (30)$$

In the same way as for (22) and (24), we derive

$$E[X^T(n) \mathbf{P} \mathbf{T} \mathbf{P} X(n)] = \text{tr}\{\mathbf{T} \Upsilon\} \quad (31)$$

$$\begin{aligned} E[X(n)X^T(n) \mathbf{P} \mathbf{T} \mathbf{P} X(n)X^T(n)] \\ = \text{tr}\{\mathbf{T} \Upsilon\} \mathbf{R} + 2\mathbf{R} \mathbf{P} \mathbf{T} \mathbf{P} \mathbf{R} \end{aligned} \quad (32)$$

Thus we can rewrite (30) as

$$\begin{aligned} \mathbf{U} &= (\mathbf{I}_M - \eta E[g(e(n))] \mathbf{R}) \mathbf{P} \mathbf{T} \mathbf{P} (\mathbf{I}_M - \eta E[g(e(n))] \mathbf{R}) + \\ &\quad \eta^2 E[g^2(e(n))] \text{tr}\{\mathbf{T} \Upsilon\} \mathbf{R} + 2\eta^2 E[g^2(e(n))] \times \\ &\quad \mathbf{R} \mathbf{P} \mathbf{T} \mathbf{P} \mathbf{R} - \eta^2 E^2[g(e(n))] \mathbf{R} \mathbf{P} \mathbf{T} \mathbf{P} \mathbf{R} \end{aligned} \quad (33)$$

From [41], some useful properties can be obtained, that is,

$$\text{vec}\{\mathbf{BCD}\} = (\mathbf{D}^T \otimes \mathbf{B}) \text{vec}\{\mathbf{C}\}$$

and

$$\text{tr}\{\mathbf{B}^T \mathbf{C}\} = \text{vec}^T\{\mathbf{C}\} \text{vec}\{\mathbf{B}\}$$

where $\text{vec}\{\cdot\}$ denotes the vectorization operator, \otimes stands for the Kronecker product. With the vectorization and the above properties, we have

$$\text{vec}\{\mathbf{U}\} = \mathbf{F} \mathbf{t} \quad (34)$$

where

$$\begin{aligned} \mathbf{F} &= (\mathbf{I}_M - \eta E[g(e(n))] \mathbf{R}) \mathbf{P} \otimes (\mathbf{I}_M - \eta E[g(e(n))] \mathbf{R}) \mathbf{P} \\ &\quad + 2\eta^2 E[g^2(e(n))] (\mathbf{R} \mathbf{P} \otimes \mathbf{R} \mathbf{P}) + \eta^2 E[g^2(e(n))] \times \\ &\quad \text{vec}\{\mathbf{R}\} \text{vec}\{\Upsilon\} - \eta^2 E^2[g(e(n))] (\mathbf{R} \mathbf{P} \otimes \mathbf{R} \mathbf{P}) \end{aligned}$$

and $\mathbf{t} = \text{vec}\{\mathbf{T}\}$. Combining (31), (32) and (34), we can rewrite (29) as

$$\begin{aligned} E[\|\tilde{W}(n)\|_{\mathbf{t}}^2] &= E[\|\tilde{W}(n-1)\|_{\mathbf{t}}^2] + \eta^2 E[g^2(e(n))] \times \\ &\quad (\varepsilon_w^T \mathbf{R} \varepsilon_w + E[v^2(n)]) \text{vec}^T\{\Upsilon\} \mathbf{t} \end{aligned} \quad (35)$$

We define the steady-state MSD as follows:

$$S = \lim_{n \rightarrow \infty} E[\|\tilde{W}(n)\|^2] \quad (36)$$

Assume that the filter is stable and achieves the steady-state, i.e. $\lim_{n \rightarrow \infty} E[\|\tilde{W}(n)\|^2] = \lim_{n \rightarrow \infty} E[\|\tilde{W}(n-1)\|^2]$. By (35), we have

$$\begin{aligned} \lim_{n \rightarrow \infty} E[\|\tilde{W}(n)\|_{(\mathbf{I}_{M^2} - \mathbf{F})}^2] &= \lim_{n \rightarrow \infty} \eta^2 E[g^2(e(n))] \times \\ &\quad (\varepsilon_w^T \mathbf{R} \varepsilon_w + E[v^2(n)]) \times \\ &\quad \text{vec}^T\{\Upsilon\} \mathbf{t} \end{aligned} \quad (37)$$

Therefore, by selecting an appropriate $\mathbf{t} = (\mathbf{I}_{M^2} - \mathbf{F})^{-1} \text{vec}\{\mathbf{I}_M\}$, we can obtain

$$\begin{aligned} S &= \eta^2 (\varepsilon_w^T \mathbf{R} \varepsilon_w + E[v^2(n)]) \text{vec}^T\{\Upsilon\} \times \\ &\quad \lim_{n \rightarrow \infty} (\mathbf{I}_{M^2} - \mathbf{F})^{-1} \text{vec}\{\mathbf{I}_M\} E[g^2(e(n))] \end{aligned} \quad (38)$$

Based on assumption 3), we can rewrite (15) as following:

$$W_{opt} = W^* + \mathbf{R}^{-1} C (C^T \mathbf{R}^{-1} C)^{-1} (f - C^T W^*) \quad (39)$$

and accordingly

$$\varepsilon_w = \mathbf{R}^{-1} C (C^T \mathbf{R}^{-1} C)^{-1} (C^T W^* - f) \quad (40)$$

In order to obtain the theoretical value of the steady-state MSD, we also need to evaluate the values of $\lim_{n \rightarrow \infty} E[g(e(n))]$ and $\lim_{n \rightarrow \infty} E[g^2(e(n))]$. We consider two cases below:

- 1) If $v(n)$ is zero-mean Gaussian distributed with variance σ_v^2 , then

$$\lim_{n \rightarrow \infty} E[g(e(n))] \approx \frac{\sigma}{\sqrt{\sigma^2 + \varepsilon_w^T \mathbf{R} \varepsilon_w + \sigma_v^2}} \quad (41)$$

$$\lim_{n \rightarrow \infty} E[g^2(e(n))] \approx \frac{\sigma}{\sqrt{\sigma^2 + 2\varepsilon_w^T \mathbf{R} \varepsilon_w + 2\sigma_v^2}} \quad (42)$$

Thus

$$\begin{aligned} S &\approx \eta^2 (\varepsilon_w^T \mathbf{R} \varepsilon_w + \sigma_v^2) \text{vec}^T\{\Upsilon\} (\mathbf{I}_{M^2} - \mathbf{F})^{-1} \times \\ &\quad \text{vec}\{\mathbf{I}_M\} \frac{\sigma}{\sqrt{\sigma^2 + 2\varepsilon_w^T \mathbf{R} \varepsilon_w + 2\sigma_v^2}} \end{aligned} \quad (43)$$

- 2) If $v(n)$ is non-Gaussian, then by Taylor expansion we have

$$\begin{aligned} \lim_{n \rightarrow \infty} E[g(e(n))] &\approx E\left[\exp\left(-\frac{v^2(n)}{2\sigma^2}\right)\right] + \frac{1}{2} \varepsilon_w^T \mathbf{R} \varepsilon_w \times \\ &\quad E\left[\left(\frac{v^2(n)}{\sigma^4} - \frac{1}{\sigma^2}\right) \exp\left(-\frac{v^2(n)}{2\sigma^2}\right)\right] \end{aligned} \quad (44)$$

$$\begin{aligned} \lim_{n \rightarrow \infty} E[g^2(e(n))] &\approx E\left[\exp\left(-\frac{v^2(n)}{\sigma^2}\right)\right] + \varepsilon_w^T \mathbf{R} \varepsilon_w \times \\ &\quad E\left[\left(\frac{2v^2(n)}{\sigma^4} - \frac{1}{\sigma^2}\right) \exp\left(-\frac{v^2(n)}{\sigma^2}\right)\right] \end{aligned} \quad (45)$$

It follows that

$$\begin{aligned} S &\approx \eta^2 (\varepsilon_w^T \mathbf{R} \varepsilon_w + E[v^2(n)]) \text{vec}^T\{\Upsilon\} (\mathbf{I}_{M^2} - \mathbf{F})^{-1} \times \\ &\quad \text{vec}\{\mathbf{I}_M\} \left(E\left[\exp\left(-\frac{v^2(n)}{\sigma^2}\right)\right] + \varepsilon_w^T \mathbf{R} \varepsilon_w \times \right. \\ &\quad \left. E\left[\left(\frac{2v^2(n)}{\sigma^4} - \frac{1}{\sigma^2}\right) \exp\left(-\frac{v^2(n)}{\sigma^2}\right)\right] \right) \end{aligned} \quad (46)$$

Remark: It is worth noting that (43) and (46) have been derived by using the approximation $W(n) \approx W_{opt}$ at the steady state. In addition, the theoretical value for non-Gaussian noise case has been derived by taking the Taylor expansion of $g(e(n))$ around $v(n)$ and omitting the higher-order terms. If the noise power is very large, the approximation is not accurate and hence, the derived values at steady state may deviate seriously from the actual results. The detailed derivations for (41) to (46) can be found in Appendix C.

IV. SIMULATION RESULTS

In this section, we present simulation results to confirm the theoretical conclusions drawn in the previous section, and illustrate the superior performance of the proposed CMCC algorithm compared with the traditional CLMS algorithm [5], CAP algorithm [11] and CRLS algorithm [12]. The selection of kernel bandwidth is also discussed in the end.

A. Non-Gaussian Noise Models

Generally speaking, the non-Gaussian noise distributions can be divided into two categories: light-tailed (e.g. binary, uniform, etc.) and heavy-tailed (e.g. Laplace, Cauchy, mixed Gaussian, alpha-stable, etc.) distributions [31], [38], [39], [42], [43]. In the following experiments, five common non-Gaussian noise models including binary noise, Laplace noise, Cauchy noise, Mixed Gaussian noise, and alpha-stable noise, are selected for performance evaluation. Descriptions of these non-Gaussian noises are as following:

- 1) Binary noise model: Standard binary noise takes the values of either $v = 1$ or $v = -1$, with probability mass function $Pr\{v = 1\} = Pr\{v = -1\} = 0.5$.
- 2) Laplace noise model: The Laplace noise is distributed with probability density function (PDF):

$$p(v) = \frac{1}{2} \exp^{-|v|} \quad (47)$$

- 3) Cauchy noise model: The PDF of the Cauchy noise is

$$p(v) = \frac{1}{\pi(1+v^2)} \quad (48)$$

- 4) Mixed Gaussian noise model: The mixed Gaussian noise model is given by:

$$(1 - \theta)\mathcal{N}(\lambda_1, v_1^2) + \theta\mathcal{N}(\lambda_2, v_2^2) \quad (49)$$

where $\mathcal{N}(\lambda_i, v_i^2)$ ($i = 1, 2$) denote the Gaussian distributions with mean values λ_i and variances v_i^2 , and θ is the mixture coefficient. Usually one can set θ to a small value and $v_2^2 \gg v_1^2$ to represent the impulsive noises (or large outliers). Therefore, we define the mixed Gaussian noise parameter vector as $V_{mix} = (\lambda_1, \lambda_2, v_1^2, v_2^2, \theta)$.

- 5) Alpha-stable noise model: The characteristic function of the alpha-stable noise is defined as:

$$\psi(t) = \exp\{j\delta t - \gamma|t|^\alpha[1 + j\beta \operatorname{sgn}(t)S(t, \alpha)]\} \quad (50)$$

in which

$$S(t, \alpha) = \begin{cases} \tan(\frac{\alpha\pi}{2}) & \text{if } \alpha \neq 1 \\ \frac{\pi}{2} \log|t| & \text{if } \alpha = 1 \end{cases} \quad (51)$$

From (50), one can observe that a stable distribution is completely determined by four parameters: 1) the characteristic factor α ; 2) the symmetry parameter β ; 3) the dispersion parameter γ ; 4) the location parameter δ . So we define the alpha-stable noise parameter vector as $V_{alpha} = (\alpha, \beta, \gamma, \delta)$.

It is worth mentioning that, in the case of $\alpha = 2$, the alpha-stable distribution coincides with the Gaussian distribution, while $\alpha = 1, \delta = 0$ is the same as the Cauchy distribution.

B. Validation of Steady-state MSD

In this experiment, we show the values of the theoretical and simulated steady-state MSDs of the CMCC in a linear channel with weight vector ($M = 7$)

$$W^* = [0.332, -0.040, -0.094, 0.717, -0.652, -0.072, 0.580]^T \quad (52)$$

Assume that $K = 3$, C is full-rank, and the input covariance matrix \mathbf{R} is positive-definite with $\operatorname{tr}\{\mathbf{R}\} = M$ [13]. The input vectors are zero-mean multivariate Gaussian, and the disturbance noises considered include Gaussian noise, binary noise (light-tailed disturbance) and Laplace noise (heavy-tailed disturbance). Fig. 2 shows the theoretical and simulated steady-state MSDs with different step-sizes, and Fig. 3 presents the theoretical and simulated steady-state MSDs with different noise variances. If not mentioned otherwise, simulation results are averaged over 500 independent Monte Carlo runs, and in each simulation, 5000 iterations are run to ensure the algorithms to reach the steady state, and the steady-state MSDs are obtained as averages over the last 200 iterations. Evidently, the steady-state MSDs are increasing with the step-size and noise variances increasing. In addition, the steady-state MSDs obtained from simulations match well with those theoretical results (computed by (41) for Gaussian noise and (46) for Non-Gaussian noise).

C. Linear System Identification

We consider a linear system identification problem where the length of the adaptive filter is equal to that of the unknown system impulse response. Assume that the weight vector W^* of the unknown system, the constraint parameters C and f , the input vectors, and the input covariance matrix \mathbf{R} are the same as the previous experiment. In the simulations below, without mentioning otherwise, 500 independent Monte Carlo simulations are performed and in each simulation, 3000 iterations are run to ensure the algorithms to reach the steady state. The sliding data length for CAP is set to 4, and the forgetting factor for CRLS is set to 0.998. The kernel bandwidth for CMCC is $\sigma = 2.0$.

First, we illustrate the performance of the proposed CMCC compared with CLMS, CAP and CRLS in four noise distributions. Simulation results are shown in Fig. 4. In the simulation, the mixed Gaussian noise parameters are set as $V_{mix} = (0, 0, 0.01, 100, 0.05)$, the alpha-stable noise parameters are set as $V_{alpha} = (1.5, 0, 0.4, 0)$, the laplace noise is zero-mean with standard deviation 5, and the cauchy noise is reduced to $\frac{1}{10}$. The step-sizes are chosen such that all the algorithms have almost the same initial convergence speed. As one can see clearly, the CMCC algorithm significantly outperforms other algorithms in terms of stability, and achieves much lower steady-state MSD.

Second, we demonstrate how the kernel bandwidth σ will influence the convergence performance of CMCC. Fig. 5 shows the convergence curves of CMCC with different σ , where the mixed gaussian noise is chosen for measurement noise and the noise parameters are the same as the previous simulation. The step-sizes are set at $\eta =$

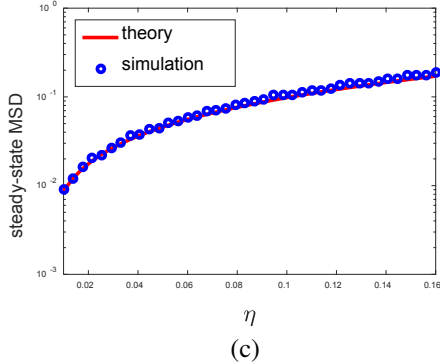
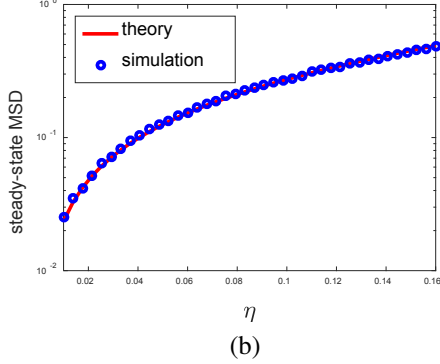
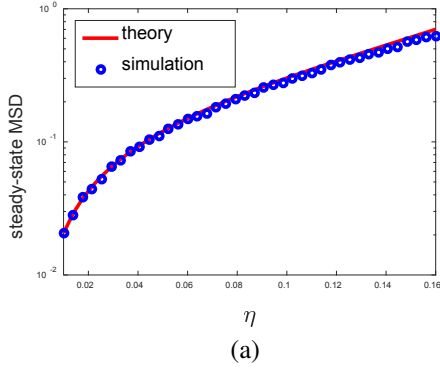


Fig. 2: Theoretical and simulated steady-state MSDs with different step-sizes η : (a) Gaussian noise ($\sigma = 8.0$, $\sigma_v^2 = 0.81$); (b) Binary noise ($\sigma = 2.0$, $\sigma_v^2 = 1.0$); (c) Laplace noise ($\sigma = 1.0$, $\sigma_v^2 = 1.0$).

0.06, 0.012, 0.01, 0.01, 0.01 for $\sigma = 0.5, 2.0, 8.0, 16.0, 32.0$ respectively. Obviously, the kernel bandwidth has significant influence on the convergence behavior. In this example, the proposed algorithm achieves the lowest steady-state MSD when $\sigma = 2.0$. If the kernel bandwidth is too larger (e.g. $\sigma = 32.0$) or too small (e.g. $\sigma = 0.5$), the convergence performance of CMCC will become poor. We provide some useful properties later for kernel bandwidth selection in practical applications.

Third, we investigate the stability problem of the CMCC in different step-sizes η . Fig. 6 illustrates the convergence performance with different step-sizes, and accordingly Fig. 7 shows the performance evolution curve. The noise is still the mixed Gaussian noise with same parameters. From simulation

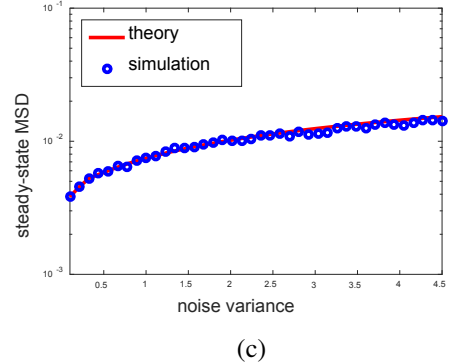
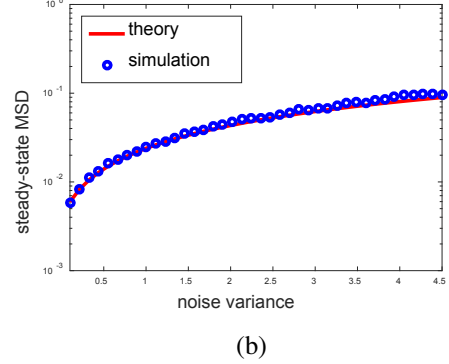
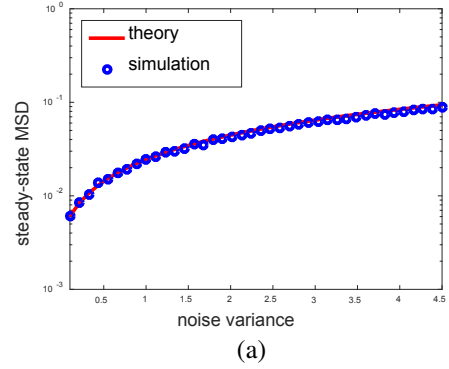


Fig. 3: Theoretical and simulated steady-state MSDs with different noise variance σ_v^2 : (a) Gaussian noise ($\eta = 0.01$, $\sigma = 8.0$); (b) Binary noise ($\eta = 0.01$, $\sigma = 6.0$); (c) Laplace noise ($\eta = 0.01$, $\sigma = 0.8$).

results, one can observe clearly that: 1) when the step-size is very large (such as $\eta \geq 0.5$), the CMCC will be divergent, which confirms the validity of the theoretical analysis of mean square stability in section III; 2) As the step-size increases, the mean and variance of MSD of the proposed algorithm become larger. Simulation results show that a larger step-size leads to a more unstable algorithm, and even make the new algorithm to become diverge. Additionally, in this simulation, we calculate the value of $\frac{2}{2\lambda_{max} + tr\{\Upsilon\}}$ (by (28)) to 0.278, not larger than 0.4, which also illustrates the effectiveness of (28).

D. Beamforming Application

In this scenario, we consider a uniform linear array consisting of $M = 7$ omnidirectional sensors with an element spacing

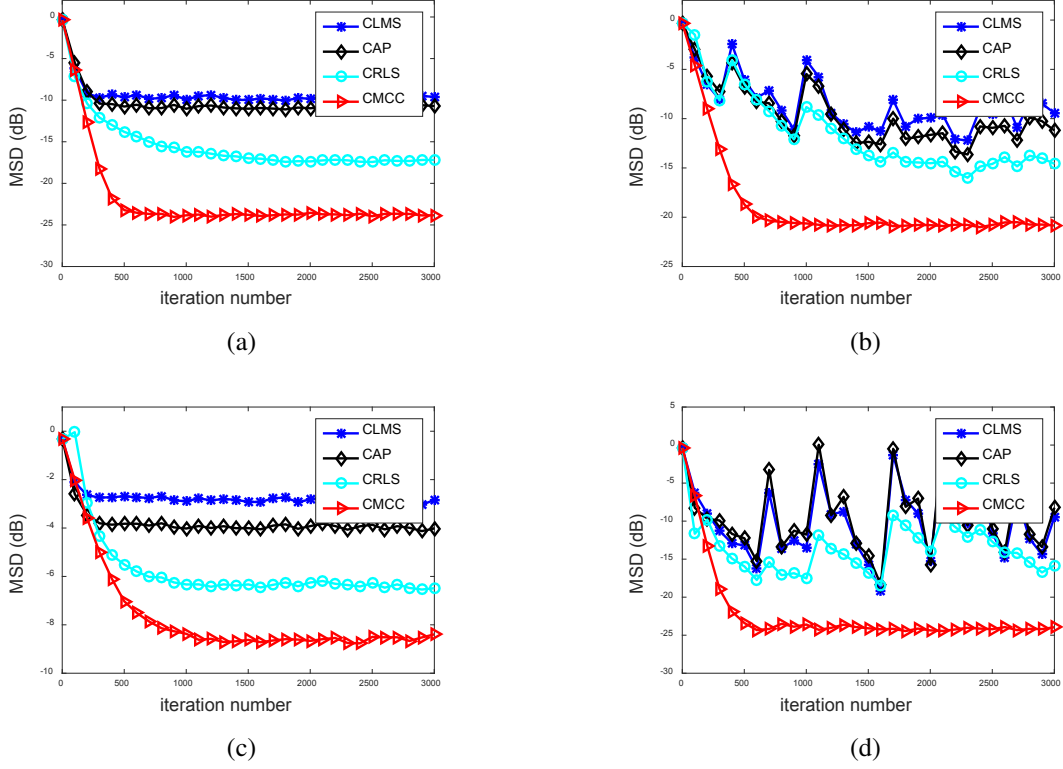


Fig. 4: Convergence curves of CLMS, CAP, CRLS and CMCC in different noises: (a) Mixed Gaussian noise; (b) Alpha-stable noise; (c) Laplace noise; (d) Cauchy noise

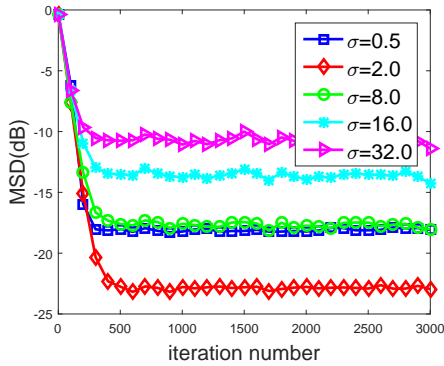


Fig. 5: Convergence curves of CMCC with different σ

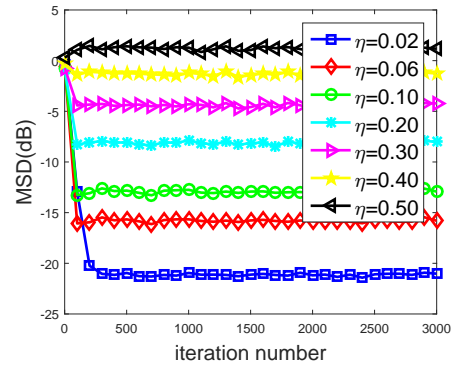


Fig. 6: Convergence curves of CMCC with different η

of half wavelength. We also assume that there are four users. Among them, the signal of one user is of interest, and is presumed to arrive at the direction-of-arrival (DOA) of $\varphi_d = 0^\circ$, while the other three signals are considered as interferers with DOAs of $\varphi_1 = -25^\circ$, $\varphi_2 = 30^\circ$, $\varphi_3 = 60^\circ$, respectively. We choose the constraint matrix $C = [\mathbf{I}_{M-1}, \mathbf{0}, -\mathbf{J}_{M-1}]$ with \mathbf{J} being a reversal matrix of size (an identity matrix with all rows in reversed order), and the response vector $\mathbf{f} = \mathbf{0}$ [13]. The measurement noise $v(n)$ is the additive non-Gaussian noise, and the measured output of the unknown system is set to $d(n) = v(n)$ [11]. In the following simulations, simulation results are averaged over 1000 independent Monte Carlo runs,

and in each simulation, 3000 iterations are run to ensure the algorithms to reach the steady state, and the steady-state MSDs are obtained as averages over the last 200 iterations. The *signal-to-noise ratio* (SNR) is set to 0 dB, and the *interference-to-noise ratio* (INR) is set to 10 dB. The sliding data length for CAP is set to 4, and the forgetting factor for CRLS is set to 0.999. The kernel bandwidth σ is set at 20.

The convergence curves of CLMS, CAP, CRLS and CMCC in alpha-stable noise are illustrated in Fig. 8, and accordingly, the beampatterns of different methods are given in Fig. 9. The noise parameters are set at $V_{alpha} = (1.2, 0, 1.6, 0)$, and other parameters are chosen such that all algorithms have almost the same initial convergence rate. As one can see that, compared

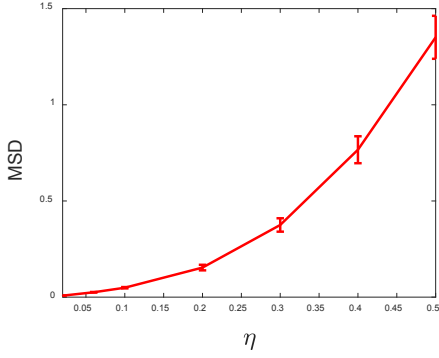


Fig. 7: Performance evolution curve of CMCC with different η .

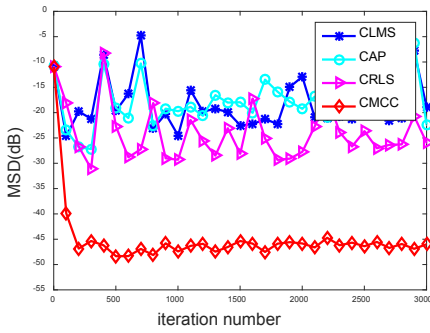


Fig. 8: Convergence curves of CLMS, CAP, CRLS and CMCC.

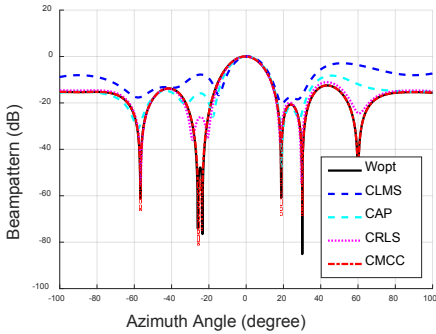


Fig. 9: Beampatterns of CLMS, CAP, CRLS and CMCC.

with these traditional constrained adaptive filtering algorithms, the proposed algorithm performs best in all scenarios in term of MSD and beampattern shape. Furthermore, it has similar performance to the optimal beamformer after convergence.

Fig. 10 shows the steady-state MSDs of CLMS, CAP, CRLS and CMCC with different $\alpha = (0.6, 0.8, 1.0, 1.2, 1.4, 1.6)$ and different $\gamma = (1.2, 1.4, 1.6, 1.7, 1.8, 1.9)$ in 3-D space. Other parameters are the same as in the previous simulation for all algorithms. As expected, the proposed algorithm can achieve much better steady-state performance than CLMS, CAP and CRLS in all cases.

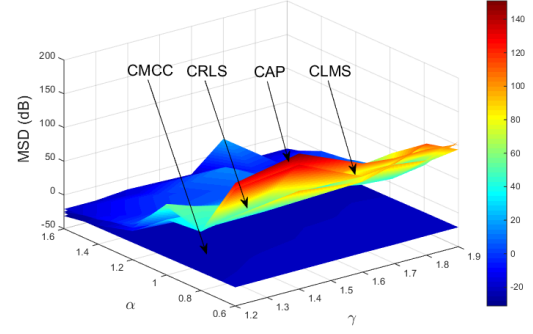


Fig. 10: Steady-state MSDs of CLMS, CAP, CRLS and CMCC in 3-D space.

E. Parameter Selection

The kernel bandwidth σ is an important free parameter in CMCC since it controls all robust properties of correntropy. An appropriate kernel bandwidth can provide an effective mechanism to eliminate the effect of outliers and noise.

According to the previous studies, some useful tricks for kernel bandwidth selection are as follows [23]–[30]:

- 1) If the data are plentiful, a small σ should be used so that high precision can be achieved; however, the kernel bandwidth must be selected to make a compromise between estimation efficiency and outlier rejection if the data are small.
- 2) As σ increases, the contribution of the higher-order moments decays faster, and the second-order moment plays a key role. Therefore, a large σ is frequently appropriate for Gaussian noises, while a small σ is usually adapt to non-Gaussian impulsive noises.
- 3) For a given noise environment, there is a relatively large range of σ that provides nearly optimal performance.

Currently, Silverman's rule, one of the most widely used methods in kernel density estimation, is often used to estimate σ . However, the limitation is that this method cannot obtain the best possible value. Therefore, in a practical application, σ is manually selected or optimized by trials and errors.

V. CONCLUSION

In this paper, we have developed the constrained maximum correntropy criterion (CMCC) adaptive filtering algorithm by incorporating a linear constraint into the maximum correntropy criterion. We also studied the mean square convergence performance including the mean square stability and the steady-state mean square deviation (MSD) of the proposed algorithm. Simulation results have confirmed the theoretical conclusions and shown that the new algorithm can significantly outperform the traditional CLMS, CAP and CRLS algorithms when the noise is of heavy-tailed non-Gaussian distribution.

A main benefit of correntropy is that the kernel bandwidth controls all its properties. However, in practical applications, the kernel bandwidth is manually selected by scanning the performance. Therefore, how to select an optimal kernel bandwidth is a big challenge for future study. On the other

hand, CMCC belongs to the family of stochastic gradient based algorithms, which usually suffers from slow convergence. It is expected to solve this problem by investigating the MCC-based constrained affine projection algorithm and recursive maximum correntropy algorithm.

APPENDIX A DERIVATION OF (12)

Based on (11), we can easily derive the following instantaneous weight update equation [24], [29]

$$\begin{aligned} W(n) &= W(n-1) + \eta \frac{\partial J_{CMCC}}{\partial W} \Big|_{W=W(n-1)} \\ &= W(n-1) + \eta g(e(n)) (d(n) - W^T(n-1)X(n)) \\ &\quad \times X(n) + \eta C\xi(n) \end{aligned} \quad (53)$$

where the weight vector W is initialized at a vector satisfying $W(0) = C(C^T C)^{-1}f$. Due to

$$\begin{aligned} f &= C^T W(n) \\ &= C^T [W(n-1) + \eta C\xi(n) + \eta g(e(n)) \times \\ &\quad (d(n) - W^T(n-1)X(n)) X(n)] \end{aligned} \quad (54)$$

we have

$$\begin{aligned} \xi(n) &= \frac{1}{\eta} (C^T C)^{-1} [f - C^T W(n-1) - \eta g(e(n)) \times \\ &\quad (d(n) - W^T(n-1)X(n)) C^T X(n)] \end{aligned} \quad (55)$$

Substituting (55) into (53), and after some simple vector manipulations, we derive

$$\begin{aligned} W(n) &= P [W(n-1) + \eta g(e(n))(d(n) - \\ &\quad W^T(n-1)X(n))X(n)] + Q \end{aligned} \quad (56)$$

which is the CMCC algorithm.

APPENDIX B DERIVATION OF (15)

Setting $\frac{\partial J_{CMCC}}{\partial W} \Big|_{W=W(n-1)} = \mathbf{0}_{M \times 1}$, one can derive the optimal weight vector W_{opt} under CMCC as follows:

$$\begin{aligned} E [g(e(n))(d(n) - W_{opt}^T X(n))X(n)] + C\xi(n) &= \mathbf{0}_{M \times 1} \\ \Rightarrow E [g(e(n))X(n)X^T(n)] W_{opt} &= E [g(e(n))d(n)X(n)] + \\ &\quad C\xi(n) \\ \Rightarrow \mathbf{R}_g W_{opt} &= \mathbf{P}_g + C\xi(n) \\ \Rightarrow W_{opt} &= \mathbf{R}_g^{-1} \mathbf{P}_g + \mathbf{R}_g^{-1} C\xi(n) \end{aligned} \quad (57)$$

where $\mathbf{P}_g = E [g(e(n))d(n)X(n)]$ is a weighted cross-correlation vector between the measured output and the input vector. Since

$$\begin{aligned} C^T W_{opt} &= f \\ \Rightarrow C^T [\mathbf{R}_g^{-1} \mathbf{P}_g + \mathbf{R}_g^{-1} C\xi(n)] &= f \\ \Rightarrow \xi(n) &= [C^T \mathbf{R}_g^{-1} C]^{-1} (f - C^T \mathbf{R}_g^{-1} \mathbf{P}_g) \end{aligned} \quad (58)$$

one can rewrite (57) as

$$\begin{aligned} W_{opt} &= \mathbf{R}_g^{-1} \mathbf{P}_g + \mathbf{R}_g^{-1} C [C^T \mathbf{R}_g^{-1} C]^{-1} \times \\ &\quad (f - C^T \mathbf{R}_g^{-1} \mathbf{P}_g) \end{aligned} \quad (59)$$

Under the assumptions 1) and 2), we derive by using (5)

$$\begin{aligned} d(n) &= W^{*T} X(n) + v(n) \\ \Rightarrow d(n)X^T(n) &= W^{*T} X(n)X^T(n) + v(n)X^T(n) \\ \Rightarrow g(e(n))d(n)X^T(n) &= g(e(n))W^{*T} X(n)X^T(n) + \\ &\quad v(n)g(e(n))X^T(n) \\ \Rightarrow \mathbf{P}_g &= \mathbf{R}_g W^* \\ \Rightarrow W^* &= \mathbf{R}_g^{-1} \mathbf{P}_g \end{aligned} \quad (60)$$

Therefore, combining (59) and (60), we obtain

$$W_{opt} = W^* + \mathbf{R}_g^{-1} C (C^T \mathbf{R}_g^{-1} C)^{-1} (f - C^T W^*) \quad (61)$$

APPENDIX C DERIVATION OF (41)~(46)

Here we consider two cases below:

1) Gaussian noise case

Since $e(n) = e_a(n) + v(n)$, in this case $e(n)$ is also zero-mean Gaussian. Let σ_e^2 be the variance of the error $e(n)$. Then we have

$$\sigma_e^2 = E [e_a^2(n)] + \sigma_v^2 \quad (62)$$

Using (16) and the approximation $W(n) \approx W_{opt}$ at the steady-state, we obtain

$$e_a(n) \approx (W^* - W_{opt})^T X(n) = \varepsilon_w X(n) \quad (63)$$

Therefore

$$\sigma_e^2 \approx \varepsilon_w^T \mathbf{R} \varepsilon_w + \sigma_v^2 \quad (64)$$

It follows that

$$\begin{aligned} \lim_{n \rightarrow \infty} E [g(e(n))] &= \lim_{n \rightarrow \infty} \frac{1}{\sqrt{2\pi}\sigma_e} \int_{-\infty}^{\infty} \exp\left(-\frac{e^2(n)}{2\sigma_e^2}\right) \times \\ &\quad \exp\left(-\frac{e^2(n)}{2\sigma_e^2}\right) de(n) \\ &= \frac{\sigma}{\sqrt{\sigma^2 + \sigma_e^2}} \\ &\approx \frac{\sigma}{\sqrt{\sigma^2 + \varepsilon_w^T \mathbf{R} \varepsilon_w + \sigma_v^2}} \end{aligned} \quad (65)$$

$$\begin{aligned} \lim_{n \rightarrow \infty} E [g^2(e(n))] &= \lim_{n \rightarrow \infty} \frac{1}{\sqrt{2\pi}\sigma_e} \int_{-\infty}^{\infty} \exp\left(-\frac{e^2(n)}{\sigma^2}\right) \times \\ &\quad \exp\left(-\frac{e^2(n)}{2\sigma_e^2}\right) de(n) \\ &= \frac{\sigma}{\sqrt{\sigma^2 + 2\sigma_e^2}} \\ &\approx \frac{\sigma}{\sqrt{\sigma^2 + 2\varepsilon_w^T \mathbf{R} \varepsilon_w + 2\sigma_v^2}} \end{aligned} \quad (66)$$

Substituting (65) and (66) into (38), we obtain

$$\begin{aligned} S &\approx \eta^2 (\varepsilon_w^T \mathbf{R} \varepsilon_w + \sigma_v^2) \frac{\text{vec}^T \{\Upsilon\} (\mathbf{I}_{M^2} - \mathbf{F})^{-1} \times \\ &\quad \text{vec}\{\mathbf{I}_M\}}{\sqrt{\sigma^2 + 2\varepsilon_w^T \mathbf{R} \varepsilon_w + 2\sigma_v^2}} \end{aligned} \quad (67)$$

2) Non-Gaussian noise case

Taking the Taylor expansion of $g(e(n))$ with respect to $e_a(n)$ around $v(n)$, we have

$$\begin{aligned} g(e(n)) &= g(e_a(n) + v(n)) \\ &= g(v(n)) + g'(v(n))e_a(n) + \\ &\quad \frac{1}{2}g''(v(n))e_a^2(n) + o(e_a^2(n)) \end{aligned} \quad (68)$$

where

$$g(v(n)) = \exp\left(-\frac{v^2(n)}{2\sigma^2}\right) \quad (69)$$

$$g'(v(n)) = -\frac{v(n)}{\sigma^2} \exp\left(-\frac{v^2(n)}{2\sigma^2}\right) \quad (70)$$

$$g''(v(n)) = \left(\frac{v^2(n)}{\sigma^4} - \frac{1}{\sigma^2}\right) \exp\left(-\frac{v^2(n)}{2\sigma^2}\right) \quad (71)$$

Thus

$$\begin{aligned} E[g(e(n))] &\approx E[g(v(n))] + \frac{1}{2}E[g''(v(n))]E[e_a^2(n)] \\ &= E\left[\exp\left(-\frac{v^2(n)}{2\sigma^2}\right)\right] + \frac{1}{2}\varepsilon_w^T \mathbf{R} \varepsilon_w \times \\ &\quad E\left[\left(\frac{v^2(n)}{\sigma^4} - \frac{1}{\sigma^2}\right) \exp\left(-\frac{v^2(n)}{2\sigma^2}\right)\right] \end{aligned} \quad (72)$$

$$\begin{aligned} E[g^2(e(n))] &\approx E[g(v(n))] + E[e_a^2(n)] \times \\ &\quad E[g(v(n))g''(v(n)) + g'^2(v(n))] \\ &= E\left[\exp\left(-\frac{v^2(n)}{\sigma^2}\right)\right] + \varepsilon_w^T \mathbf{R} \varepsilon_w \times \\ &\quad E\left[\left(\frac{2v^2(n)}{\sigma^4} - \frac{1}{\sigma^2}\right) \exp\left(-\frac{v^2(n)}{\sigma^2}\right)\right] \end{aligned} \quad (73)$$

Substituting (72) and (73) into (38) yields

$$\begin{aligned} S &\approx \eta^2 (\varepsilon_w^T \mathbf{R} \varepsilon_w + E[v^2(n)]) \mathbf{vec}^T\{\Upsilon\} (\mathbf{I}_{M^2} - \mathbf{F})^{-1} \times \\ &\quad \mathbf{vec}\{\mathbf{I}_M\} \left(E\left[\exp\left(-\frac{v^2(n)}{\sigma^2}\right)\right] + \varepsilon_w^T \mathbf{R} \varepsilon_w \times \right. \\ &\quad \left. E\left[\left(\frac{2v^2(n)}{\sigma^4} - \frac{1}{\sigma^2}\right) \exp\left(-\frac{v^2(n)}{\sigma^2}\right)\right] \right) \end{aligned} \quad (74)$$

ACKNOWLEDGMENTS

This work was supported by 973 Program (No. 2015CB351703) and National Natural Science Foundation of China (No. 61372152).

REFERENCES

- [1] M. L. R. de Campos, S. Werner, and J. A. Apolinrio, Jr., "Constrained Adaptive Filters," in *Adaptive Antenna Arrays: Trends and Applications*, S. Chandra, Ed. New York: Springer-Verlag, 2004.
- [2] M. Moinuddin, A. Zerguine, and A. U. H. Sheikh, "Multiple-Access Interference Plus Noise-Constrained Least Mean Square (MNCLMS) Algorithm for CDMA Systems," *IEEE Trans. Circuits Sys. I, Regular papers*, vol. 55, no. 9, pp. 2870-2883, Oct. 2008.
- [3] S. D. Lin, and M. Barkat, "The performance of the SMI method in the constrained LMS array and the Griffiths array," *IEEE Trans. Antennas. Propagation*, vol. 38, no. 11, pp. 1878-1882, Nov. 2002.
- [4] L. G. Wang, D. F. Liu, and Q. M. Wang, "Geometric Method of Fully Constrained Least Squares Linear Spectral Mixture Analysis," *IEEE Trans. Geosci. Remote sens.*, vol. 51, no. 6, pp. 3558-3566, June. 2013.
- [5] O. L. Frost, III, "An algorithm for linearly constrained adaptive array processing," *Proc. IEEE*, vol. 60, no. 8, pp. 926-935, Aug. 1972.
- [6] H. L. Van Trees, *Detection, Estimation, and Modulation Theory Part IV: Optimum Array Processing*, New York: Wiley, 2002.
- [7] L. C. Godara, and C. Antonio, "Analysis of constrained LMS algorithm with application to adaptive beamforming using perturbation sequences," *IEEE Trans. Antennas. Propagation*, vol. 34, no. 3, pp. 368-379, Mar. 1986.
- [8] R. B. Staszewski, K. Muhammad, and P. T. Balsara, "A constrained asymmetry LMS algorithm for PRML disk drive read channels," *IEEE Trans. Circuits Sys. II, Analog and Digital Signal Processing*, vol. 48, no. 8, pp. 793-798, Aug. 2001.
- [9] W. Hu, and W. P. Tay, "Multi-Hop Diffusion LMS for Energy-Constrained Distributed Estimation," *IEEE Trans. Signal Process.*, vol. 63, no.15, pp. 4022-4036, Aug. 2015.
- [10] X. Li, and W. K. Jenkins, "The comparison of the constrained and unconstrained frequency-domain block-LMS adaptive algorithms," *IEEE Trans. Signal Process.*, vol. 44, no. 1, Jul. 1996.
- [11] S. Werner, J. A. Apolinrio, Jr., M. L. R. de Campos, and P. S. R. Diniz, "Low-complexity constrained affine-projection algorithms," *IEEE Trans. Signal Process.*, vol. 53, no. 12, pp. 4545-4555, Dec.2005.
- [12] L. S. Resende, J. M. T. Romano, and M. G. Bellanger, "A fast least squares algorithm for linearly constrained adaptive filtering," *IEEE Trans. Signal Process.*, vol. 44, no. 5, pp. 1168-1174, May 1996.
- [13] R. Arablouei and K. Dogancay, "Reduced-complexity constrained recursive least-squares adaptive filtering algorithm," *IEEE Trans. Signal Process.*, vol. 60, no. 12, pp. 6687-6692, Dec. 2012.
- [14] R. Arablouei and K. Dogancay, "Linearly-constrained recursive total least-squares algorithm," *IEEE Signal Process. Letters*, vol. 19, no. 12, pp. 821-824, 2012.
- [15] K. Lee, Y. Baek, and Y. Park, "Nonlinear acoustic echo cancellation using a nonlinear postprocessor with a linearly constrained affine projection algorithm," *IEEE Trans. Circuits Syst. II Exp. Briefs*, vol. 62, no. 9, pp. 881-885, Sep. 2015.
- [16] S. Werner, J. A. Jr., and M. de Campos, "The data-selective constrained affine projection algorithm," in *Proc. IEEE Int. Conf. on Acous. Speech and Signal Process. (ICASSP)*, vol. 6, 2001, pp. 3745-3748.
- [17] S. Haykin, *Adaptive Filtering Theory*, 3rd ed., NY: Prentice Hall, 1996.
- [18] A. H. Sayed, *Fundamentals of Adaptive Filtering*, John Wiley & Sons, 2003.
- [19] B. Chen, Y. Zhu, J. Hu and J. C. Principe, *System Parameter Identification: Information Criteria and Algorithms*, Newnes, 2013.
- [20] K. N. Plataniotis, D. Androutsos and A. N. Venetsanopoulos, "Nonlinear filtering of non-Gaussian noise," *J. Intelligent and Robotic System*, vol. 19, no. 2, pp. 207-231, June. 1997.
- [21] B. Weng and K. E. Barner, "Nonlinear system identification in impulsive environments," *IEEE Trans. Signal Process*, vol. 53, no. 7, pp. 2588-2594, July. 2005.
- [22] T. W. Hilands and S. Thomopoulos, "Nonlinear filtering methods for Harmonic retrieval and model order selection in Gaussian and non-Gaussian noise," *IEEE Trans. Signal Processing*, vol. 45, no. 4, pp. 982995, Apr. 1997.
- [23] J. C. Principe, *Information Theoretic Learning: Renyi's Entropy and Kernel Perspectives*, Springer: New York, NY, USA, 2010.
- [24] W. Ma, H. Qua, G. Gui, L. Xu, J. H. 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, Apr. 2015.
- [25] S. Zhao, B. Chen, and J. C. Principe, "Kernel adaptive filtering with maximum correntropy criterion," in *Proc. Int. Joint Conf. Neural Networks (IJCNN)*, 2011, pp. 2012-2017.
- [26] L. Shi and Y. Lin, "Convex combination of adaptive filters under the maximum correntropy criterion in impulsive interference," *IEEE Signal Process. Lett.*, vol. 21, no. 11, pp. 1385-1388, July 2014.
- [27] W. Liu, P. P. Pokharel and J. C. Principe, "Correntropy: properties and applications in non-Gaussian signal processing," *IEEE Trans. Signal Processing*, vol. 55, no. 11, pp. 5286-5298, Nov. 2007.
- [28] B. Chen and J. C. Principe, "Maximum correntropy estimation is a smoothed MAP estimation," *IEEE Signal Processing Letters*, vol. 19, no. 8, pp. 491-494, Aug. 2012.
- [29] 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, Oct. 2015.

- [30] A. Singh and J. C. Principe, "Using correntropy as a cost function in linear adaptive filters," in *Proc. Int. Joint Conf. Neural Networks (IJCNN)*, 2009, pp. 2950-2955.
- [31] B. Chen, L. Xing, H. Zhao, N. Zheng, J. C. Principe, "Generalized correntropy for robust adaptive filtering," *IEEE Trans. Signal Processing*, vol. 64, no. 13, pp. 3376-3387, July. 2016.
- [32] 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, June. 2011.
- [33] R. He, W. Zheng, and B. Hu, "Maximum correntropy criterion for robust face recognition," *IEEE Transactions Pattern Analysis and Machine Intelligence*, vol. 33, no. 8, pp. 1561-1576, Feb. 2011.
- [34] R. J. Bessa, V. Miranda and J. Gama, "Entropy and correntropy against minimum square error in offline and online three-day ahead wind power forecasting," *IEEE Transactions on Power Systems*, vol. 24, no. 4, pp. 1657-1666, Sep. 2009.
- [35] E. Hasanbelliu, L. S. Giraldo, and J. C. Principe, "Information theoretic shape matching," *IEEE Transactions Pattern Analysis and Machine Intelligence*, vol. 36, no. 12, pp. 2436 - 2451, Dec. 2014.
- [36] B. Chen, Y. Zhu, and J. Hu, "Mean-square convergence analysis of ADALINE training with minimum error entropy criterion," *IEEE Trans. Neural Networks*, vol. 21, no. 7, pp. 1168-1179, July 2010.
- [37] R. Arablouei and K. Dogancay, "On the mean-square performance of the constrained LMS algorithm," *Signal Processing*, vol. 117, pp. 192-197, Dec. 2015.
- [38] T. Y. Al-Naffouri and A. H. Sayed, "Adaptive filters with error nonlinearities: Mean-square analysis and optimum design," *EURASIP J. Appl. Signal Process.*, vol. 4, pp. 192-205, 2001.
- [39] B. Chen, L. Xing, J. Liang, N. Zheng, and J. C. Principe, "Steady-state Mean-square error analysis for adaptive filtering under the maximum correntropy criterion," *IEEE Signal Process. Letters*, vol. 21, no. 7, pp. 880-884, 2014.
- [40] L. Isserlis, "On a formula for the product-moment coefficient of any order of a normal frequency distribution in any number of variables," *Biometrika*, vol. 12, no. 1/2, pp. 134-139, Nov. 1918.
- [41] K. M. Abadir and J. R. Magnus, *Matrix Algebra*, NY: Cambridge Univ. Press, 2005.
- [42] P. G. Georgiou, P. Tsakalides and C. Kyriakakis. "Alpha-stable modeling of noise and robust time-delay estimation in presence of impulsive noise," *IEEE Trans. on Multimedia*, vol. 1, no. 3, pp. 291-301, Sept. 1999.
- [43] C. L. Nikias and M. Shao, *Signal Processing with Alpha-Stable Distributions and Applications*, John Wiley & Sons Canada, 1995.