

Independent Vector Extraction for Fast Joint Blind Source Separation and Dereverberation

Rintaro Ikeshita, *Member, IEEE*, and Tomohiro Nakatani, *Fellow, IEEE*

Abstract—We address a blind source separation (BSS) problem in a noisy reverberant environment in which the number of microphones M is greater than the number of sources of interest, and the other noise components can be approximated as stationary and Gaussian distributed. Conventional BSS algorithms for the optimization of a multi-input multi-output convolutional beamformer have suffered from a huge computational cost when M is large. We here propose a computationally efficient method that integrates a weighted prediction error (WPE) dereverberation method and a fast BSS method called independent vector extraction (IVE), which has been developed for less reverberant environments. We show that, given the power spectrum for each source, the optimization problem of the new method can be reduced to that of IVE by exploiting the stationary condition, which makes the optimization easy to handle and computationally efficient. An experiment of speech signal separation shows that, compared to a conventional method that integrates WPE and independent vector analysis, our proposed method achieves much faster convergence while maintaining its separation performance.

Index Terms—Blind source separation, dereverberation, independent vector analysis, block coordinate descent method

I. INTRODUCTION

When multiple speech signals are observed by distant microphones (e.g., in a conference room), they are contaminated with reverberation and background noise. The problem of extracting each speech signal and removing the reverberation and background noise from only the observed signal is called (convolutive) blind source separation or extraction (BSE) [1]–[3]. Here, we consider BSE in the short-term Fourier transform (STFT) domain under the following two conditions:

- The reverberation time (RT_{60}) is larger than the frame length of the STFT, and the mixture should be treated as a convolutive mixture in the STFT domain as well.
- The number of microphones M is greater than that of speech signals K and there can be background noise.

To cope with reverberation, one can apply a dereverberation method [4] such as weighted prediction error (WPE) [5]–[7] as preprocessing of BSE for instantaneous mixtures in the STFT domain (called BSE-inst in this paper). We then apply some BSE-inst method such as independent vector analysis (IVA) [8]–[10] and independent vector extraction (IVE) [11]–[17] developed for less reverberant environments, to extract K speech signals. Such a cascade configuration of WPE and IVA/IVE has a low computational cost, but the WPE dereverberation filter is estimated without considering the separation attained by IVA/IVE following WPE.

To jointly optimize the WPE dereverberation and separation filters through a unified optimization, methods that integrate

WPE and several BSE-inst methods have been proposed [6], [7], [18]–[20], and it has been reported that these methods can give higher separation performance than the cascade configuration of WPE and BSE-inst (see, e.g., [18]). However, the computational cost of optimizing both WPE and BSE-inst models becomes huge when M is large.

To reduce the computational cost of the conventional joint optimization methods while maintaining their separation performance, we propose a new BSE method called *IVE for convolutive mixtures (IVE-conv)*, which integrates WPE and IVE (Section III). We show that, given source power spectra, the IVE-conv optimization problem can be reduced to the IVE optimization problem by exploiting the stationary condition, and this reduction is not computationally intensive (Section IV-A). The IVE optimization problem can be solved fast [13]–[17], and so can the IVE-conv optimization problem (Section IV-B). We also propose another new algorithm for IVE-conv that alternately optimizes WPE and IVE (Section IV-C). Similar algorithms have already been developed in [6], [7], but our proposed one significantly reduces the computational time complexity of the conventional ones. In a numerical experiment in which two speech signals are extracted from mixtures, we show the effectiveness of our new approach.

II. BLIND SOURCE EXTRACTION PROBLEM

Let M be the number of microphones. Suppose that an observed mixture $\mathbf{x} := \{\mathbf{x}(f, t)\}_{f, t} \subset \mathbb{C}^M$ in the STFT domain is a convolutive mixture of K nonstationary source signals and $N_z := M - K$ background noise signals:¹

$$\mathbf{x}(f, t) = \sum_{\tau=0}^{N_\tau} \left[\sum_{i=1}^K \mathbf{a}_i(f, \tau) s_i(f, t - \tau) + A_z(f, \tau) \mathbf{z}(f, t - \tau) \right],$$

$$\mathbf{a}_i(f, \tau) \in \mathbb{C}^M, \quad s_i(f, t) \in \mathbb{C}, \quad i \in \{1, \dots, K\}, \quad (1)$$

$$A_z(f, \tau) \in \mathbb{C}^{M \times N_z}, \quad \mathbf{z}(f, t) \in \mathbb{C}^{N_z}. \quad (2)$$

Here, $f = 1, \dots, F$ and $t = 1, \dots, T$ denote the frequency bin and time frame indexes, respectively. Also, $s_i(f, t) \in \mathbb{C}$ and $\mathbf{z}(f, t) \in \mathbb{C}^{N_z}$ are the signals of the target source $i = 1, \dots, K$ and the background noises, respectively. $\{\mathbf{a}_i(f, \tau)\}_{\tau=0}^{N_\tau}$ and $\{A_z(f, \tau)\}_{\tau=0}^{N_\tau}$ are the acoustic transfer functions (ATFs) for the corresponding sources, where $N_\tau + 1$ is the length of the ATFs. The BSE problem addressed in this paper is defined as the problem of estimating the sources of interest, i.e., $\{s_i(f, t)\}_{i, f, t}$. We assume that K is given and the background noises are more stationary than the sources of interest.

¹The assumption that the dimension of the noise signal is $M - K$ concerns the rigorous development of efficient algorithms and can be violated to some extent when applied in practice (see numerical experiments in Section V).

III. PROBABILISTIC MODEL

We present the proposed IVE-conv model that integrates WPE [5]–[7] and IVE [11]–[17]. Let $\hat{\mathbf{x}}(f, t) \in \mathbb{C}^{M+L}$ with $L = M(D_2 - D_1 + 1)$ and $0 < D_1 \leq D_2$ be given by

$$\hat{\mathbf{x}}(f, t) = [\mathbf{x}(f, t)^\top, \mathbf{x}(f, t - D_1)^\top, \dots, \mathbf{x}(f, t - D_2)^\top]^\top,$$

where \top is the transpose of a vector. Suppose that there exists a convolutional filter $\hat{W}(f) \in \mathbb{C}^{(M+L) \times M}$ satisfying

$$s_i(f, t) = \hat{\mathbf{w}}_i(f)^h \hat{\mathbf{x}}(f, t) \in \mathbb{C}, \quad i \in \{1, \dots, K\}, \quad (3)$$

$$\mathbf{z}(f, t) = \hat{W}_z(f)^h \hat{\mathbf{x}}(f, t) \in \mathbb{C}^{N_z}, \quad (4)$$

$$\hat{W}(f) = [\hat{\mathbf{w}}_1(f), \dots, \hat{\mathbf{w}}_K(f), \hat{W}_z(f)] \in \mathbb{C}^{(M+L) \times M}, \quad (5)$$

where h denotes the conjugate transpose. As pointed out in [21], [22], convolutional filter $\hat{W}(f)$ can be decomposed into the WPE prediction matrix $G(f) \in \mathbb{C}^{L \times M}$ and the ICA separation matrix $W(f) \in \mathbb{C}^{M \times M}$:

$$\hat{W}(f) = \begin{bmatrix} W(f) \\ -G(f)W(f) \end{bmatrix} = \begin{bmatrix} I_M \\ -G(f) \end{bmatrix} W(f). \quad (6)$$

Here, $I_d \in \mathbb{C}^{d \times d}$ is the identity matrix.

We also assume that the original source signals are mutually independent and that the target source (resp. noise) signals obey time-dependent (resp. time-independent) complex Gaussian distributions in the same way as in IVE [11]–[17]:

$$\mathbf{s}_i(t) := [s_i(1, t), \dots, s_i(F, t)]^\top \in \mathbb{C}^F, \quad (7)$$

$$\mathbf{s}_i(t) \sim \mathcal{CN}(\mathbf{0}_F, v_i(t)I_F), \quad v_i(t) \in \mathbb{R}_{>0}, \quad (8)$$

$$\mathbf{z}(f, t) \sim \mathcal{CN}(\mathbf{0}_{N_z}, \Omega(f)), \quad \Omega(f) \in \mathcal{S}_{++}^{N_z}, \quad (9)$$

$$\{\mathbf{s}_i(t), \mathbf{z}(f, t)\}_{i, f, t} \text{ are mutually independent.} \quad (10)$$

Here, $\mathbf{0}_d \in \mathbb{C}^d$ is the zero vector, \mathcal{S}_{++}^d denotes the set of all Hermitian positive definite matrices of size $d \times d$, and $\mathbb{R}_{>0} = \mathcal{S}_{++}^1$. Assumption (9) that the background noise signal is stationary and Gaussian distributed is essential for developing computationally efficient algorithms. In Section V, we will experimentally show that this assumption can be violated to some extent when applied in practice.

The IVE-conv model is defined by (3)–(10). The parameters $\hat{W} := \{\hat{W}(f)\}_f$, $v := \{v_i(t)\}_{i, t}$, and $\Omega := \{\Omega(f)\}_f$ can be estimated based on maximum likelihood, which is equivalent to minimizing $\hat{g}(\hat{W}, \Omega, v) := -\frac{1}{T} \log p(\mathbf{x})$:

$$\begin{aligned} \hat{g}(\hat{W}, \Omega, v) &= \sum_{f=1}^F \sum_{i=1}^K \left[\hat{\mathbf{w}}_i(f)^h \hat{R}_i(f) \hat{\mathbf{w}}_i(f) + \frac{1}{T} \sum_{t=1}^T \log v_i(t) \right] \\ &+ \sum_{f=1}^F \text{tr}(\hat{W}_z(f)^h \hat{R}_z(f) \hat{W}_z(f) \Omega(f)^{-1}) \\ &- \sum_{f=1}^F \log \det(W(f)^h W(f) \Omega(f)^{-1}), \end{aligned} \quad (11)$$

$$\hat{R}_i(f) = \frac{1}{T} \sum_{t=1}^T \frac{\hat{\mathbf{x}}(f, t) \hat{\mathbf{x}}(f, t)^h}{v_i(t)}, \quad i \in \{1, \dots, K, z\},$$

where we define $v_z(t) = 1$ for all $t = 1, \dots, T$ (see, e.g., [19] for the derivation of \hat{g}). If $L = 0$ and $\hat{W}(f) = W(f)$, then objective function \hat{g} has the same form as the counterparts

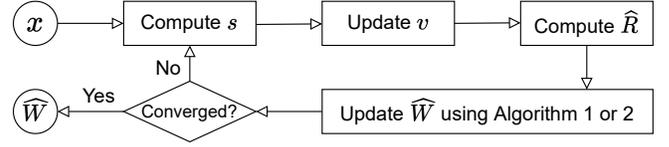


Fig. 1: Flowchart of IVE-conv

of ICA, IVA, and IVE, which has been discussed extensively in the literature [13]–[17], [23]–[27]. For $L \geq 1$, $K = M$, and $N_z = 0$, the optimization problem has been discussed explicitly in [18], [19] and implicitly in [6], [7], [20].

Remark 1. The proposed IVE-conv is an integration of WPE and IVE. If we replace IVE with ICA, IVA, or independent low-rank matrix analysis (ILRMA) [28], then the IVE-conv turns out to be the method that integrates WPE with ICA [6], WPE with IVA (IVA-conv) [18], or WPE with ILRMA [19], [20], respectively. In this sense, the novelty of the IVE-conv model might seem limited. However, if M gets large, computationally efficient algorithms can be developed only for IVE-conv, which is our main contribution.²

IV. OPTIMIZATION ALGORITHM

To obtain a local optimal solution for the minimization problem of (11), two block coordinate descent (BCD [34]) algorithms summarized in Table I will be developed. All the algorithms shown in Table I update v and (\hat{W}, Ω) alternately. The flowchart of IVE-conv is shown in Figure 1.

When (\hat{W}, Ω) are kept fixed, v can be optimized as

$$v_i(t) = \frac{1}{F} \|\mathbf{s}_i(t)\|_2^2 = \frac{1}{F} \mathbf{s}_i(t)^h \mathbf{s}_i(t). \quad (12)$$

In what follows, we will develop two BCDs to optimize (\hat{W}, Ω) while keeping v fixed. Because this subproblem can be addressed independently for each frequency bin, we focus only on optimizing $\hat{W}(f)$ and $\Omega(f)$, and the frequency bin index f is dropped off to ease the notation. Also, we will denote the submatrices of \hat{W} and \hat{R}_i , $i \in \{1, \dots, K, z\}$ as

$$\hat{W} = \begin{bmatrix} W \\ -GW \end{bmatrix} = \begin{bmatrix} W \\ \bar{W} \end{bmatrix} = \begin{bmatrix} \mathbf{w}_1 & \cdots & \mathbf{w}_K & W_z \\ \bar{\mathbf{w}}_1 & \cdots & \bar{\mathbf{w}}_K & \bar{W}_z \end{bmatrix}, \quad (13)$$

$$\mathbf{w}_i \in \mathbb{C}^M, \quad \bar{\mathbf{w}}_i \in \mathbb{C}^L, \quad W_z \in \mathbb{C}^{M \times N_z}, \quad \bar{W}_z \in \mathbb{C}^{L \times N_z},$$

$$\hat{R}_i = \begin{bmatrix} R_i & \bar{P}_i^h \\ \bar{P}_i & \bar{R}_i \end{bmatrix} \in \mathcal{S}_{++}^{M+L}, \quad \bar{P}_i \in \mathbb{C}^{L \times M}, \quad \bar{R}_i \in \mathcal{S}_{++}^L.$$

A. Reduction from IVE-conv to IVE when v is kept fixed

Before developing the algorithms, we show that the problem of minimizing \hat{g} with respect to \hat{W} and Ω (when source power spectra $v = \{v_i(t)\}_{i, t}$ are kept fixed), i.e.,

$$(\hat{W}, \Omega) \in \underset{(\hat{W}, \Omega)}{\text{argmin}} \hat{g}(\hat{W}, \Omega, v), \quad (14)$$

²This letter is based on our work [29] reported in a domestic workshop in which an algorithm similar to but less efficient than Algorithm 1 (proposed in Section IV-B) was first presented. Recently, as follow-up research of our previous work [29], a method has been developed [30] that replaces the IVE-conv spectrum model (7)–(8) with a model using nonnegative matrix factorization (NMF) [31]–[33]. In contrast, here, we develop a more efficient Algorithm 1 in a rigorous way by providing new insight into the IVE-conv optimization problem in Section IV-A. In addition, Algorithm 2 proposed in Section IV-C is completely new.

can be reduced to problem (17) below that has been addressed in the study of IVE [13]–[17].

Every optimal \bar{W} (the lower part of \hat{W}) in problem (14) satisfies the stationary condition [35], which is computed as

$$\begin{aligned} \frac{\partial \hat{g}}{\partial \bar{\mathbf{w}}_i^*} = \mathbf{0}_L &\iff \bar{P}_i \mathbf{w}_i + \bar{R}_i \bar{\mathbf{w}}_i = \mathbf{0}_L \in \mathbb{C}^L, \\ &\iff \bar{\mathbf{w}}_i = -\bar{R}_i^{-1} \bar{P}_i \mathbf{w}_i \in \mathbb{C}^L, \end{aligned} \quad (15)$$

$$\begin{aligned} \frac{\partial \hat{g}}{\partial \bar{W}_z^*} = O &\iff \bar{P}_z W_z + \bar{R}_z \bar{W}_z = O \in \mathbb{C}^{L \times N_z}, \\ &\iff \bar{W}_z = -\bar{R}_z^{-1} \bar{P}_z W_z \in \mathbb{C}^{L \times N_z}, \end{aligned} \quad (16)$$

where $*$ denotes the element-wise conjugate. Eqs. (15) and (16) imply that the optimal \bar{W} is a function of W and that the variable \bar{W} can be removed from \hat{g} by substituting (15) and (16). In other words, problem (14) is equivalent to the following problem through (15) and (16):

$$(W, \Omega) \in \underset{(W, \Omega)}{\operatorname{argmin}} g(W, \Omega, v), \quad (17)$$

$$g = \sum_{i=1}^K \mathbf{w}_i^h V_i \mathbf{w}_i + \operatorname{tr} (W_z^h V_z W_z \Omega^{-1}) - \log \det (W^h W \Omega^{-1}),$$

$$V_i := R_i - \bar{P}_i^h \bar{R}_i^{-1} \bar{P}_i \in \mathcal{S}_{++}^M, \quad i \in \{1, \dots, K, z\}. \quad (18)$$

Since problem (17) is nothing but the problem addressed in the study of IVE, we can directly apply efficient algorithms that have been developed for IVE [13]–[17]. Our new algorithm developed in Section IV-B is based on this observation.

B. Algorithm 1: Update each convolutional filter one by one

To solve problem (14), we propose a cyclic BCD algorithm that updates $\hat{\mathbf{w}}_1 \rightarrow (\hat{W}_z, \Omega) \rightarrow \dots \rightarrow \hat{\mathbf{w}}_K \rightarrow (\hat{W}_z, \Omega)$ one by one by solving the following subproblems:

$$\hat{\mathbf{w}}_i \in \underset{\hat{\mathbf{w}}_i}{\operatorname{argmin}} \hat{g}(\hat{\mathbf{w}}_1, \dots, \hat{\mathbf{w}}_K, \hat{W}_z, \Omega, v), \quad (19)$$

$$(\hat{W}_z, \Omega) \in \underset{(\hat{W}_z, \Omega)}{\operatorname{argmin}} \hat{g}(\hat{\mathbf{w}}_1, \dots, \hat{\mathbf{w}}_K, \hat{W}_z, \Omega, v). \quad (20)$$

From the observation given in Section IV-A, these subproblems can be equivalently transformed to

$$\mathbf{w}_i \in \underset{\mathbf{w}_i}{\operatorname{argmin}} \mathbf{w}_i^h V_i \mathbf{w}_i - \log \det (W^h W), \quad (21)$$

$$(W_z, \Omega) \in \underset{(W_z, \Omega)}{\operatorname{argmin}} g_z(W_z, \Omega), \quad (22)$$

$$g_z(W_z, \Omega) = \operatorname{tr} (W_z^h V_z W_z \Omega^{-1}) - \log \det (W^h W \Omega^{-1})$$

through (15) and (16), respectively. Here, V_i and V_z are defined by (18). As shown in [27], problem (21) can be solved as

$$\mathbf{u}_i \leftarrow (W^h V_i)^{-1} \mathbf{e}_i \in \mathbb{C}^M, \quad (23)$$

$$\mathbf{w}_i \leftarrow \mathbf{u}_i (\mathbf{u}_i^h V_i \mathbf{u}_i)^{-\frac{1}{2}} \in \mathbb{C}^M, \quad (24)$$

where \mathbf{e}_i is the i -th column of I_M . On the other hand, as shown in [16, Proposition 4], problem (22) can be solved as

$$W_z \leftarrow \begin{bmatrix} (W_s^h V_z E_s)^{-1} (W_s^h V_z E_z) \\ -I_{N_z} \end{bmatrix} \in \mathbb{C}^{M \times N_z}, \quad (25)$$

$$\Omega \leftarrow W_z^h V_z W_z \in \mathcal{S}_{++}^{N_z}, \quad (26)$$

where $W_s := [\mathbf{w}_1, \dots, \mathbf{w}_K] \in \mathbb{C}^{M \times K}$, $E_s \in \mathbb{C}^{M \times K}$ is the first K columns of I_M , and $E_z \in \mathbb{C}^{M \times N_z}$ is the last N_z columns of I_M , i.e., $[E_s, E_z] = I_M$.

Remark 2. The update formula for \mathbf{w}_i , i.e., (15), (18), (23), and (24), has already been developed in our previous paper [18], [29] in a different manner. In this subsection, we reveal that it can also be developed by exploiting the stationary condition. The efficient update formula for \hat{W}_z , i.e., (16), (18), and (25), is newly developed based on the stationary Gaussian assumption of the background noises. There is no need to update Ω as it does not affect the behavior of the algorithm.

C. Algorithm 2: Alternate update of WPE and ICA

In Section III, we recalled by (6) that convolutional filter \hat{W} can be decomposed into WPE prediction matrix G and ICA separation matrix W . Here, we develop a new cyclic BCD that updates $G \rightarrow \mathbf{w}_1 \rightarrow W_z \rightarrow \dots \rightarrow \mathbf{w}_K \rightarrow W_z$ one by one by solving the following subproblems:

$$G \in \underset{G}{\operatorname{argmin}} \hat{g}(G, \mathbf{w}_1, \dots, \mathbf{w}_K, W_z, \Omega, v), \quad (27)$$

$$\mathbf{w}_i \in \underset{\mathbf{w}_i}{\operatorname{argmin}} \hat{g}(G, \mathbf{w}_1, \dots, \mathbf{w}_K, W_z, \Omega, v), \quad (28)$$

$$(W_z, \Omega) \in \underset{(W_z, \Omega)}{\operatorname{argmin}} \hat{g}(G, \mathbf{w}_1, \dots, \mathbf{w}_K, W_z, \Omega, v). \quad (29)$$

When $K = M$ and there are no noise components, problems (27) and (28) have already been discussed in [6], [7], [18], [20]. However, the conventional algorithms to solve (27) suffer from a huge computational cost as shown in Table I. We thus propose a more computationally efficient algorithm.

1) *Algorithm to solve problems (28) and (29):* We first explain how to solve problems (28) and (29). By substituting Eq. (6) into objective function \hat{g} , these problems can be simply expressed as problems (21) and (22), respectively, except that V_i is replaced by the following V'_i for each $i \in \{1, \dots, K, z\}$:

$$V'_i = \begin{bmatrix} I_M \\ -G \end{bmatrix}^h \hat{R}_i \begin{bmatrix} I_M \\ -G \end{bmatrix} \in \mathcal{S}_{++}^M. \quad (30)$$

Thus, in the same way as in the previous subsection, problem (28) can be solved as (23)–(24), where V_i is replaced by V'_i . Also, problem (29) can be solved as (25)–(26), where V_z is replaced by V'_z .

2) *Algorithm to solve problem (27):* We next propose an algorithm to solve (27) with less computational time complexity than conventional ones. Every optimal $G \in \mathbb{C}^{L \times M}$ of problem (27) (when W , Ω , and v are kept fixed) satisfies the stationary condition, which can be computed as

$$O_{L,M} = \frac{\partial \hat{g}}{\partial G^*} = - \frac{\partial \hat{g}}{\partial \bar{W}^*} \Big|_{\bar{W} = -GW} W^h, \quad (31)$$

$$\iff \begin{cases} G \mathbf{w}_i = \bar{R}_i^{-1} \bar{P}_i \mathbf{w}_i, & i = 1, \dots, K, \\ G W_z = \bar{R}_z^{-1} \bar{P}_z W_z, \end{cases} \quad (32)$$

$$\iff G = [\bar{R}_1^{-1} \bar{P}_1 \mathbf{w}_1 \mid \dots \mid \bar{R}_K^{-1} \bar{P}_K \mathbf{w}_K \mid \bar{R}_z^{-1} \bar{P}_z W_z] W^{-1}. \quad (33)$$

Here, we used (15) and (16) to derive (32). Because problem (27) is (strictly) convex, the update formula (33) gives the (unique) global optimal solution. The computational time

TABLE I: Optimization process of BCD

Method	Reference	Optimization process ¹⁾	Computational time complexity
Conventional	[18], [19] [6], [7], [18], [20]	$v \rightarrow \hat{w}_1 \rightarrow \dots \rightarrow \hat{w}_K \rightarrow \hat{w}_{K+1} \rightarrow \dots \rightarrow \hat{w}_M$ $v \rightarrow G \rightarrow w_1 \rightarrow \dots \rightarrow w_K \rightarrow w_{K+1} \rightarrow \dots \rightarrow w_M$	$O(ML^2FT + ML^3F)$ $O(ML^2FT + M^3L^3F)$ in [18]
Proposed	§IV-B (Algorithm 1) §IV-C (Algorithm 2)	$v \rightarrow \hat{w}_1 \rightarrow (\hat{W}_z, \Omega) \rightarrow \dots \rightarrow \hat{w}_K \rightarrow (\hat{W}_z, \Omega)$ $v \rightarrow G \rightarrow w_1 \rightarrow (W_z, \Omega) \rightarrow \dots \rightarrow w_K \rightarrow (W_z, \Omega)$	$O((K+1)L^2FT + (K+1)L^3F)$

¹⁾ We use the notations $\hat{W}_z = [\hat{w}_{K+1}, \dots, \hat{w}_M] \in \mathbb{C}^{(M+L) \times (M-K)}$ and $W_z = [w_{K+1}, \dots, w_M] \in \mathbb{C}^{M \times (M-K)}$.

²⁾ The IVA and IVE source models can be freely changed to the ICA and ILRMA source models, and so we discuss only the IVA or IVE models.

TABLE II: Methods tested in experiment

Method	Description
IVE [13], [16]	Identical to IVE-conv-(Alg1) with $L = 0$
IVA-conv [18]	An integration of WPE and IVA, which is identical to IVE-conv-(Alg1) with $K = M$, $D_1 = 2$, and $D_2 = 5$.
IVE-conv-(Alg1)	IVE-conv with $D_1 = 2$ and $D_2 = 5$ using Algorithm 1.
IVE-conv-(Alg2)	IVE-conv with $D_1 = 2$ and $D_2 = 5$ using Algorithm 2. For every five updates to v and W , we updated G once.

complexity to calculate (33) is shown in Table I, which is much smaller than that of the conventional methods.

V. EXPERIMENT

In this experiment, we evaluated the signal extraction and runtime performance of the four methods described in Table II.

Dataset: We generated synthesized convolutive noisy mixtures of two speech signals. We obtained speech signals from the test set of the TIMIT corpus [36] and concatenated them so that the length of each signal exceeded 10 seconds. We obtained point-source noise signals recorded in a cafe (CAF) and a pedestrian area (PED) from the third ‘CHiME’ Speech Separation and Recognition Challenge (CHiME-3) [37]. Note that the noise signals are nonstationary, but are considered to be more stationary than speech signals. We obtained RIR data recorded in room OFC from the RWCP Sound Scene Database in Real Acoustical Environments [38]. The reverberation time (RT_{60}) of room OFC is 780 ms.

The generated mixtures consisted of $K = 2$ speech signals and six noise signals randomly chosen from the above dataset. The SNR of each mixture was adjusted to $\text{SNR} = 10 \log_{10} \frac{(\lambda_1^{(s)} + \lambda_2^{(s)})/2}{\lambda_1^{(n)} + \dots + \lambda_6^{(n)}} = 5$ or 10 [dB], where $\lambda_i^{(s)}$ and $\lambda_j^{(n)}$ denote the sample variances of the i -th speech signal ($i = 1, 2$) and the j -th noise signal ($j = 1, \dots, 6$).

Criteria: Using *museval* [39], we measured the signal-to-distortion ratio (SDR) [40] between the separated and oracle spatial images of the speech signals at the first microphone. The oracle spatial images were obtained by truncating the RIRs at 32 ms (i.e., the points after 32 ms were replaced by 0) and convolving them with the speech signals.

Conditions: The sampling rate was 16 kHz, the frame length was 2048 (128 ms), and the frame shift was 512 (32 ms).

Initialization: For all methods, we initialized the convolutional filter as $W(f) = -I_M$ and $\hat{W}(f) = G(f) = O$, and then updated $W_z(f)$ once using (25) before the optimization.

A. Experimental results

Figure 2 shows the convergence of the SDR when each method was applied. Compared to IVE, which does not

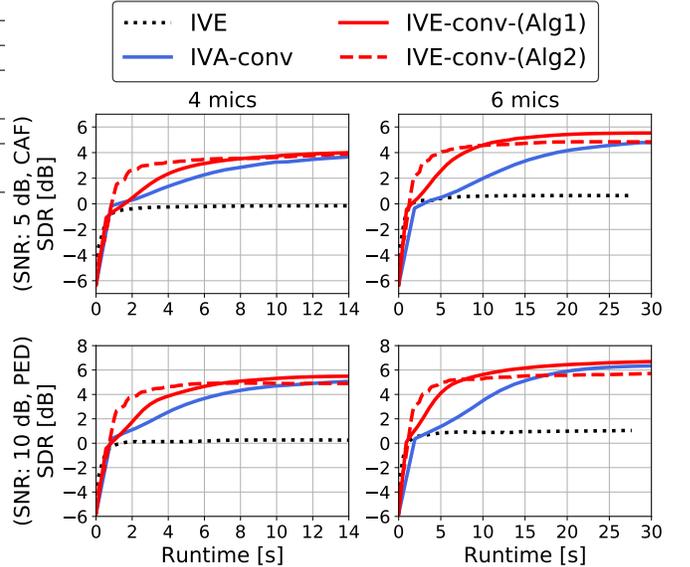


Fig. 2: SDR [dB] performance as a function of runtime. The noise condition was CAF with $\text{SNR} = 5$ [dB] (top) or PED with $\text{SNR} = 10$ [dB] (bottom), and the number of microphones was $M = 4$ (left) or 6 (right). Results shown were averaged over 50 mixtures and obtained by running the algorithms on a PC with ‘‘Intel(R) Core(TM) i7-7820 CPU @ 3.60 GHz’’ using a single thread. The average length of the mixture signals is 12.51 sec. The separated spatial image was obtained by $(W(f)^{-h} e_i)(\hat{w}_i(f)^h \hat{x}(f, t)) \in \mathbb{C}^M$ for each source $i = 1, 2$.

handle reverberation, both IVA-conv and IVE-conv showed the higher SDRs. Although the SDR performance at the convergence points is comparable, the convergence of the proposed IVE-conv is much faster than that of IVA-conv since the computational cost to update \hat{W}_z is much lower. This fast convergence behavior is important in practice, since using more microphones can improve the SDR at the expense of increased runtime as observed in Fig. 2. IVE-conv-(Alg2) converged faster than IVE-conv-(Alg1), but gave a slightly lower SDR.

VI. CONCLUSION

To achieve joint source separation and dereverberation with a small computational cost, we proposed IVE-conv, which is an integration of IVE and WPE. We also developed two efficient BCD algorithms for optimizing IVE-conv. The experimental results showed that IVE-conv yields significantly faster convergence than the integration of IVA and WPE while maintaining its separation performance.

REFERENCES

- [1] M. S. Pedersen, J. Larsen, U. Kjems, and L. C. Parra, "Convolutional blind source separation methods," *Springer handbook of speech processing*, pp. 1065–1094, 2008.
- [2] P. Comon and C. Jutten, *Handbook of Blind Source Separation: Independent component analysis and applications*. Academic press, 2010.
- [3] A. Cichocki and S. Amari, *Adaptive blind signal and image processing: learning algorithms and applications*. John Wiley & Sons, 2002.
- [4] P. A. Naylor and N. D. Gaubitch, *Speech dereverberation*. Springer Science & Business Media, 2010.
- [5] T. Nakatani, T. Yoshioka, K. Kinoshita, M. Miyoshi, and B.-H. Juang, "Speech dereverberation based on variance-normalized delayed linear prediction," *IEEE Trans. Audio, Speech, Language Process.*, vol. 18, no. 7, pp. 1717–1731, 2010.
- [6] T. Yoshioka, T. Nakatani, M. Miyoshi, and H. G. Okuno, "Blind separation and dereverberation of speech mixtures by joint optimization," *IEEE Trans. Audio, Speech, Language Process.*, vol. 19, no. 1, pp. 69–84, 2011.
- [7] T. Yoshioka and T. Nakatani, "Generalization of multi-channel linear prediction methods for blind MIMO impulse response shortening," *IEEE Trans. Audio, Speech, Language Process.*, vol. 20, no. 10, pp. 2707–2720, 2012.
- [8] T. Kim, H. T. Attias, S.-Y. Lee, and T.-W. Lee, "Blind source separation exploiting higher-order frequency dependencies," *IEEE Trans. Audio, Speech, Language Process.*, vol. 15, no. 1, pp. 70–79, 2007.
- [9] A. Hiroe, "Solution of permutation problem in frequency domain ICA, using multivariate probability density functions," in *Proc. ICA*, 2006, pp. 601–608.
- [10] Y.-O. Li, T. Adali, W. Wang, and V. D. Calhoun, "Joint blind source separation by multiset canonical correlation analysis," *IEEE Trans. Signal Process.*, vol. 57, no. 10, pp. 3918–3929, 2009.
- [11] Z. Koldovský and P. Tichavský, "Gradient algorithms for complex non-Gaussian independent component/vector extraction, question of convergence," *IEEE Trans. Signal Process.*, vol. 67, no. 4, pp. 1050–1064, 2018.
- [12] J. Janský, J. Málek, J. Čmejla, T. Kounovský, Z. Koldovský, and J. Žďánský, "Adaptive blind audio source extraction supervised by dominant speaker identification using x-vectors," in *Proc. ICASSP*, 2020, pp. 676–680.
- [13] R. Scheibler and N. Ono, "Independent vector analysis with more microphones than sources," in *Proc. WASPAA*, 2019, pp. 185–189.
- [14] —, "Fast independent vector extraction by iterative SINR maximization," in *Proc. ICASSP*, 2020, pp. 601–605.
- [15] —, "MM algorithms for joint independent subspace analysis with application to blind single and multi-source extraction," *arXiv:2004.03926v1*, 2020.
- [16] R. Ikeshita, T. Nakatani, and S. Araki, "Block coordinate descent algorithms for auxiliary-function-based independent vector extraction," *arXiv:2010.08959v1*, 2020.
- [17] —, "Overdetermined independent vector analysis," in *Proc. ICASSP*, 2020, pp. 591–595.
- [18] T. Nakatani, R. Ikeshita, K. Kinoshita, H. Sawada, and S. Araki, "Computationally efficient and versatile framework for joint optimization of blind speech separation and dereverberation," in *Proc. Interspeech*, 2020.
- [19] R. Ikeshita, N. Ito, T. Nakatani, and H. Sawada, "Independent low-rank matrix analysis with decorrelation learning," in *Proc. WASPAA*, 2019, pp. 288–292.
- [20] H. Kagami, H. Kameoka, and M. Yukawa, "Joint separation and dereverberation of reverberant mixtures with determined multichannel non-negative matrix factorization," in *Proc. ICASSP*, 2018, pp. 31–35.
- [21] C. Boeddeker, T. Nakatani, K. Kinoshita, and R. Haeb-Umbach, "Jointly optimal dereverberation and beamforming," in *Proc. ICASSP*, 2020, pp. 216–220.
- [22] T. Nakatani, C. Boeddeker, K. Kinoshita, R. Ikeshita, and R. Haeb-Umbach, "Jointly optimal denoising, dereverberation and source separation," *IEEE/ACM Trans. Audio, Speech, Language Process.*, vol. 28, pp. 2267–2282, 2020.
- [23] D.-T. Pham and J.-F. Cardoso, "Blind separation of instantaneous mixtures of nonstationary sources," *IEEE Trans. Signal Process.*, vol. 49, no. 9, pp. 1837–1848, 2001.
- [24] S. Dégerine and A. Zaïdi, "Determinant maximization of a nonsymmetric matrix with quadratic constraints," *SIAM J. Optim.*, vol. 17, no. 4, pp. 997–1014, 2006.
- [25] A. Yeredor, B. Song, F. Roemer, and M. Haardt, "A "sequentially drilled" joint congruence (SeDJoCo) transformation with applications in blind source separation and multiuser MIMO systems," *IEEE Trans. Signal Process.*, vol. 60, no. 6, pp. 2744–2757, 2012.
- [26] N. Ono and S. Miyabe, "Auxiliary-function-based independent component analysis for super-Gaussian sources," in *Proc. LVA/ICA*, 2010, pp. 165–172.
- [27] N. Ono, "Stable and fast update rules for independent vector analysis based on auxiliary function technique," in *Proc. WASPAA*, 2011, pp. 189–192.
- [28] D. Kitamura, N. Ono, H. Sawada, H. Kameoka, and H. Saruwatari, "Determined blind source separation unifying independent vector analysis and nonnegative matrix factorization," *IEEE/ACM Trans. Audio, Speech, Language Process.*, vol. 24, no. 9, pp. 1622–1637, 2016.
- [29] R. Ikeshita and T. Nakatani, "Independent vector extraction," in *Proc. ASJ Spring Meeting*, 2020, (in Japanese).
- [30] M. Togami and R. Scheibler, "Over-determined speech source separation and dereverberation," in *Proc. APSIPA*, 2020, pp. 705–710.
- [31] D. D. Lee and H. S. Seung, "Learning the parts of objects by non-negative matrix factorization," *Nature*, vol. 401, no. 6755, p. 788, 1999.
- [32] C. Févotte, N. Bertin, and J.-L. Durrieu, "Nonnegative matrix factorization with the Itakura-Saito divergence: With application to music analysis," *Neural Comput.*, vol. 21, no. 3, pp. 793–830, 2009.
- [33] P. Smaragdis and J. C. Brown, "Non-negative matrix factorization for polyphonic music transcription," in *Proc. WASPAA*, 2003, pp. 177–180.
- [34] P. Tseng, "Convergence of a block coordinate descent method for nondifferentiable minimization," *Journal of optimization theory and applications*, vol. 109, no. 3, pp. 475–494, 2001.
- [35] J. Nocedal and S. Wright, *Numerical optimization*. Springer Science & Business Media, 2006.
- [36] J. Garofolo, L. Lamel, W. Fisher, J. Fiscus, D. Pallett, N. Dahlgren, and V. Zue, "TIMIT Acoustic-Phonetic Continuous Speech Corpus LDC93S1. Web Download. Philadelphia: Linguistic Data Consortium," 1993.
- [37] J. Barker, R. Marxer, E. Vincent, and S. Watanabe, "The third 'CHiME' speech separation and recognition challenge: Dataset, task and baselines," in *Proc. ASRU*, 2015, pp. 504–511.
- [38] S. Nakamura, K. Hiyane, F. Asano, T. Nishiura, and T. Yamada, "Acoustical sound database in real environments for sound scene understanding and hands-free speech recognition," in *LREC*, 2000.
- [39] F.-R. Stöter, A. Liutkus, and N. Ito, "The 2018 signal separation evaluation campaign," in *Proc. LVA/ICA*, 2018, pp. 293–305.
- [40] E. Vincent, R. Gribonval, and C. Févotte, "Performance measurement in blind audio source separation," *IEEE Trans. Audio, Speech, Language Process.*, vol. 14, no. 4, pp. 1462–1469, 2006.