

Design of Spectrogram-Consistency Regularization Term Dependent on Observation in Independent Low-Rank Matrix Analysis for Blind Source Separation

Takaaki Kojima*, Norihiro Takamune*, Daichi Kitamura†, and Hiroshi Saruwatari*

* The University of Tokyo, Japan

† National Institute of Technology, Kagawa College, Japan

Abstract—Independent low-rank matrix analysis (ILRMA) is the state-of-the-art technique for blind source separation under the overdetermined condition. Recently, some methods that focus on the spectrogram consistency of the separated signals to improve the separation performance have been proposed. One of such methods introduces into ILRMA the regularization term that directly induces the separated signals to be consistent. Although in the conventional method, the regularization term is designed to be independent of the observed signals, we design the regularization term to be dependent on the observed signals. For the obtained cost function, we derive a new update rule on the basis of the majorization-minimization algorithm with the auxiliary function that is an extension of that derived in the conventional method. Finally, a numerical experiment is conducted to verify the separation performance.

I. INTRODUCTION

Blind source separation (BSS) is a technique for estimating individual sources from an observed mixture without knowledge of the characteristics of sources or mixing conditions. Independent component analysis (ICA) [1] is a representative approach to solving BSS under the overdetermined condition (the number of microphones M is greater than or equal to that of sources N). In ICA, the demixing matrix is estimated under the assumption that individual sources are statistically independent and the mixing system is instantaneous. Since the sources are mixed by time-domain convolution for acoustic signals, ICA is often applied in the time-frequency domain via short-time Fourier transform (STFT) under the assumption that the mixing system can be approximated by instantaneous mixing in the time-frequency domain. This method is called frequency-domain ICA (FDICA) [2] and estimates frequency-wise demixing matrices by independently applying ICA to complex-valued signals in each frequency bin. ICA has arbitrariness to the permutation of the estimated signals; therefore, the permutations should be aligned for frequency bins in FDICA (permutation problem). Several permutation solvers have been studied to solve this problem [3]–[5]. In contrast,

independent vector analysis (IVA) [6] and independent low-rank matrix analysis (ILRMA) [7] can estimate the frequency-wise permutations simultaneously with the demixing matrix. These techniques assume a more complex source model than FDICA and can therefore capture dependences between frequency bins and align the frequency-wise permutations. Among these methods, it has been reported that ILRMA tends to show a higher separation performance than FDICA and IVA [7]. However, it has been reported that there is room for performance improvement even in ILRMA [8].

The spectrogram consistency (or simply *consistency*) [9]–[11] has been studied as a new promising clue to align the frequency-wise permutations [12]–[14]. In some BSS methods in the time-frequency domain, such as FDICA, IVA, and ILRMA, assumptions on the spectrograms of the separated signals are made. However, these assumptions are not guaranteed to be satisfied for the final output in the time domain obtained as the inverse STFT (ISTFT) of the spectrograms of the separated signals, i.e., the spectrograms obtained as the STFT of the final output are often different from the spectrograms of the separated signals before ISTFT. This is because the operation of STFT after ISTFT is a projection operator on the image space of the STFT as a linear operator. The image space of the STFT is called the consistent subspace in this paper. For example, if the frequency-wise permutations are not aligned, the spectrogram is changed greatly by the projection and becomes blurred in both time and frequency directions as shown in Fig. 1. Then, the performance of BSS is expected to be improved by making assumptions on the projected spectrograms or by reducing the difference between the spectrograms before and after the projection. It has been reported that the performance of FDICA without a permutation solver can be improved by incorporating the projection operator into the cost function of FDICA [12].

As for ILRMA with the spectrogram consistency considered, consistent ILRMA (Consist. ILRMA) [13] and ILRMA with spectrogram-consistency regularization [14] have been proposed. In Consist. ILRMA [13], the update rules of ILRMA are modified to approximate the projected spectrograms instead of the spectrograms of the separated signals (before the

This work was supported by JST Moonshot R&D Grant Number JP-MJMS2011 (for algorithm development) and Tateisi Science and Technology Foundation (for practical experiment).

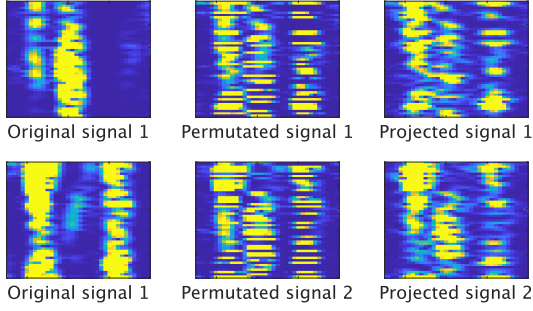


Fig. 1: Examples of the spread by projecting inconsistent spectrograms onto the consistent subspace. Original speech signals (left) were permuted randomly (middle) with frequency bins and projected onto the consistent subspace (right).

projection) by nonnegative matrix factorization (NMF) [15]. In [14], a regularization term that induces the spectrograms of the separated signals (before the projection) themselves to be close to the consistent subspace is designed for ILRMA. ILRMA with the spectrogram-consistency regularization has the advantage of being formulated as an optimization problem, whereas Consist. ILRMA is heuristic and unfortunately not formulated as an optimization problem. The regularization term proposed in [14] is designed to be independent of the observed signals; in other words, it is determined solely by the STFT conditions (e.g., the window length, window function, and hopsize). Therefore, neither the specific structure nor correlation of the observed signals can be fully utilized in [14]. In this paper, we design a new regularization term for ILRMA on the basis of the distances between the spectrograms of the separated signals and the projected spectrograms of the separated signals onto the consistent subspace. The proposed regularization term depends on the observed signals; thus, the separation performance is expected to be improved more than in [14]. For the new cost function, we unfortunately cannot use the same technique of [14] to design the auxiliary function for the majorization-minimization (MM) algorithm [16]. To address this problem, we extend the auxiliary-function-design methodology from [14] by utilizing the mathematical finding in this paper and design an auxiliary function on the basis of this finding. We derive a new update rule by applying vector-wise coordinate descent (VCD) [17] to the designed auxiliary function. Finally, we verify the separation performance of the proposed method by a simulation experiment.

II. PRELIMINARIES

In this paper, we consider a time-domain signal to be a real-valued signal of length L . The complex-valued spectrogram of a time-domain signal $\tilde{\psi} \in \mathbb{R}^L$ is represented by the supervector $\vec{\Psi} \in \mathbb{C}^{IJ}$, where I and J are the window length and the number of time frames in STFT, respectively. Here, $(I(j-1)+i)$ th element of $\vec{\Psi}$ is the i th frequency bin of the spectrum in the j th frame, where $i = 1, \dots, I$ and $j = 1, \dots, J$ are the indices of the frequency bins and time frames, respectively.

Hereafter, we assume the practical case where $IJ > L$ holds. Let τ be the hopsize of STFT, $\mathbf{h} \in \mathbb{R}^I$ be the window function. The length of the signal after zero-padding is denoted by $L' = L^{(\text{pre})} + L + L^{(\text{post})}$, where $L^{(\text{pre})} (< I)$ and $L^{(\text{post})} (< I)$ are the length of zero-padding at the beginning and end of the time signal, respectively. Note that $L' = \tau(J-1) + I$ holds. The $L' \times L$ matrix representing zero-padding is given by

$$\mathbf{Z} = [\mathbf{O}_{L \times L^{(\text{pre})}} \quad \mathbf{E}_L \quad \mathbf{O}_{L \times L^{(\text{post})}}]^T \in \mathbb{R}^{L' \times L}, \quad (1)$$

where $\mathbf{O}_{q \times q'}$ is the $q \times q'$ zero matrix and \mathbf{E}_q is the $q \times q$ identity matrix for any nonnegative integer q, q' . Here, \cdot^T represents the transpose. That is, the signal after zero padding can be represented by $\mathbf{Z}\tilde{\psi}$. The $(I(j-1)+1)$ th to (Ij) th elements of $\tilde{\psi}$, which is the spectrum in the j th frame, are obtained by multiplying the $(\tau(j-1)+1)$ th to $(\tau(j-1)+I)$ th samples of the zero-padded signal $\mathbf{Z}\tilde{\psi}$ by the window function \mathbf{h} and applying discrete Fourier transform (DFT). Thus, the vector created by arranging the $(I(j-1)+1)$ th to (Ij) th element of $\vec{\Psi}$ is represented by $\mathbf{F}\Delta_j\mathbf{Z}\tilde{\psi}$, where $\mathbf{F} \in \mathbb{C}^{I \times I}$ denotes the representation matrix of DFT and $\Delta_j \in \mathbb{R}^{I \times L'}$ is defined as

$$\Delta_j = [\mathbf{O}_{I \times \tau(j-1)} \quad \text{diag}\{\mathbf{h}\} \quad \mathbf{O}_{I \times (L' - \tau(j-1) - I)}]. \quad (2)$$

Here, $\text{diag}\{\mathbf{h}\}$ represents the diagonal matrix whose diagonal elements are \mathbf{h} . Therefore, $\vec{\Psi}$ is represented by $\vec{\Psi} = \mathcal{F}\tilde{\psi}$, where $\mathcal{F} \in \mathbb{C}^{IJ \times L}$ is defined as

$$\mathcal{F} = \begin{bmatrix} \mathbf{F}\Delta_1\mathbf{Z} \\ \vdots \\ \mathbf{F}\Delta_J\mathbf{Z} \end{bmatrix}. \quad (3)$$

In this paper, we consider that the STFT conditions satisfy the condition that \mathcal{F} has full column rank so that ISTFT exists. We use the Moore–Penrose inverse \cdot^\dagger of \mathcal{F} as ISTFT, i.e., \mathcal{F}^\dagger is a representation matrix of ISTFT. Note that $\mathcal{F}^\dagger\mathcal{F}$ is equal to \mathbf{E}_L and $\mathcal{F}\mathcal{F}^\dagger$ is an orthogonal projection matrix onto the image space of \mathcal{F} . A spectrogram $\vec{\Psi}' \in \mathbb{C}^{IJ}$ is called *consistent* when $\vec{\Psi}'$ is in the image space of \mathcal{F} , i.e., $\vec{\Psi}' = \mathcal{F}\mathcal{F}^\dagger\vec{\Psi}'$ holds.

III. CONVENTIONAL METHODS

A. ILRMA[7]

Let $s_{ij,n}$, $x_{ij,m}$, and $y_{ij,n}$ be the (i, j) th bin in the spectrogram of the n th source, the m th observed, and n th separated signals, respectively, where $n = 1, \dots, N$ and $m = 1, \dots, M$ denote the indices of the source and observed signals, respectively. We define $\mathbf{s}_{ij} = (s_{ij,1} \cdots s_{ij,N})^T \in \mathbb{C}^N$, $\mathbf{x}_{ij} = (x_{ij,1} \cdots x_{ij,M})^T \in \mathbb{C}^M$, and $\mathbf{y}_{ij} = (y_{ij,1} \cdots y_{ij,N})^T \in \mathbb{C}^N$. When the mixing system is time-invariant and the reverberation time is sufficiently shorter than the length of the STFT window, we can assume instantaneous mixing in the time-frequency domain as $\mathbf{x}_{ij} = \mathbf{A}_i\mathbf{s}_{ij}$, where $\mathbf{A}_i \in \mathbb{C}^{M \times N}$ is the mixing matrix. In this paper, we assume that $M = N$ and \mathbf{A}_i is regular. Under these assumptions, there exists the demixing matrix $\mathbf{W}_i = (\mathbf{w}_{i,1} \cdots \mathbf{w}_{i,N})^H \simeq \mathbf{A}_i^{-1}$ that satisfies $\mathbf{y}_{ij} = \mathbf{W}_i\mathbf{x}_{ij}$. Let $w_{i,nm}$ be the m th elements of $\mathbf{w}_{i,n}$.

In ILRMA, $y_{ij,n}$ is assumed to be generated from a complex Gaussian distribution with a mean of zero and the variance $\rho_{ij,n}$ as follows:

$$p(y_{ij,n}) = \frac{1}{\pi \rho_{ij,n}} \exp\left(-\frac{|y_{ij,n}|^2}{\rho_{ij,n}}\right). \quad (4)$$

The variance $\rho_{ij,n}$ is assumed to have a low-rank structure and modeled by NMF [15] as

$$\rho_{ij,n} = \sum_{k=1}^K t_{ik,n} v_{kj,n}, \quad (5)$$

where $t_{ik,n} \geq 0$ and $v_{kj,n} \geq 0$ represent the basis and the activation, respectively, and $k = 1, \dots, K$ denotes the index of the basis. The demixing matrix $\{\mathbf{W}_i\}$ and the NMF variables $\{t_{ik,n}\}$, $\{v_{kj,n}\}$ are obtained from maximum likelihood estimation by minimizing the following negative log-likelihood as a cost function:

$$\begin{aligned} \mathcal{L} = & \frac{1}{2} \sum_{i,j,n} \left(\frac{|\mathbf{w}_{i,n}^H \mathbf{x}_{ij}|^2}{\sum_k t_{ik,n} v_{kj,n}} + \log \sum_k t_{ik,n} v_{kj,n} \right) \\ & - \frac{J}{2} \sum_i \log |\det \mathbf{W}_i|^2 + \text{const.}, \end{aligned} \quad (6)$$

where const. denotes the terms independent of $\{\mathbf{W}_i\}$, $\{t_{ik,n}\}$ and $\{v_{kj,n}\}$. To optimize (6), the demixing matrix $\{\mathbf{W}_i\}$ and the NMF variables $\{t_{ik,n}\}$, $\{v_{kj,n}\}$ are updated alternately. The demixing matrix $\{\mathbf{W}_i\}$ is updated with iterative projection (IP) [18]. The NMF variables $\{t_{ik,n}\}$, $\{v_{kj,n}\}$ are updated on the basis of the MM algorithm [16]. All of these update rules ensure the monotonic nonincrease in the cost function (6).

B. ILRMA with spectrogram-consistency regularization independent of observation [14]

In [14], the regularization term based on spectrogram consistency is introduced to ILRMA. Let $\vec{\mathbf{Y}}_n$ be the supervectorized spectrogram of the n th separated signal whose $(I(j-1)+i)$ th element is $y_{ij,n}$, and the degree to which $\vec{\mathbf{Y}}_n$ deviates from the consistent subspace is formulated as [14]

$$\begin{aligned} \mathcal{E}(\vec{\mathbf{Y}}_n) = & \|\vec{\mathbf{Y}}_n - \mathcal{F} \mathcal{F}^\dagger \vec{\mathbf{Y}}_n\|_2^2 \\ = & \left\| (\mathbf{E}_{IJ} - \mathcal{F} \mathcal{F}^\dagger) \sum_m \mathcal{W}_{nm} \mathcal{F} \tilde{\mathbf{x}}_m \right\|_2^2, \end{aligned} \quad (7)$$

where $\mathcal{W}_{nm} \in \mathbb{C}^{IJ \times IJ}$ is a diagonal matrix whose $(I(j-1)+i)$ th diagonal element is $w_{i,nm}^*$ and $\tilde{\mathbf{x}}_m \in \mathbb{R}^L$ is the m th observed signal in the time domain. Here, \cdot^* denotes the complex conjugate. To induce consistency for arbitrary observed signals, the case where $\tilde{\mathbf{x}}_m$ is a white Gaussian signal without correlation between channels is considered, and the expected value $\sum_n \mathbb{E}[\mathcal{E}(\vec{\mathbf{Y}}_n)]$ is used as a regularization term in [14]. In this paper, the regularization used in [14] is called observation-independent consistency regularization (OICR).

The regularization term can be represented as the quadratic form of $w_{i,n}$ as follows:

$$\sum_n \mathbb{E}[\mathcal{E}(\vec{\mathbf{Y}}_n)] = \sum_{n,i,i'} \mathbf{w}_{i,n}^H \mathbf{G}_{ii'}^{(\text{ind})} \mathbf{w}_{i',n}, \quad (8)$$

where $\mathbf{G}_{ii'}^{(\text{ind})} \in \mathbb{C}^{M \times M}$ is only determined by the STFT conditions. Since $\tilde{\mathbf{x}}_m$ is assumed to be independent and identically distributed with respect to m , $\mathbf{G}_{ii'}^{(\text{ind})} = g_{ii'}^{(\text{ind})} \mathbf{E}_M$ ($g_{ii'}^{(\text{ind})} \in \mathbb{C}$) holds. The cost function of ILRMA with OICR is formulated as follows:

$$\mathcal{L}_{\text{OICR}} = \mathcal{L} + \beta \sum_{n,i,i'} g_{ii'}^{(\text{ind})} \mathbf{w}_{i,n}^H \mathbf{w}_{i',n}, \quad (9)$$

where $\beta > 0$ is the weight of the regularization term in ILRMA with OICR. Since the added regularization term (8) is independent of the NMF variables, the update rules for the NMF variables $\{t_{ik,n}\}$ and $\{v_{kj,n}\}$ are the same as in ILRMA. For $\{\mathbf{W}_i\}$, an auxiliary function is designed for the MM algorithm [16] and updated using VCD [17] instead of IP (see details in [14]). The update rule of $\{\mathbf{W}_i\}$ also guarantees the monotonic nonincrease in the cost function (9).

IV. PROPOSED METHOD

A. Motivation and approach

Since the regularization term in [14] is formulated to be independent of the to-be-separated observed signals, the regularization cannot be appropriate for specific observed signals. Therefore, if we can introduce a regularization term that depends on the specific structure of the to-be-separated observed signals, further improvement in separation performance is expected. In this paper, we propose the use of $\sum_n \mathcal{E}(\vec{\mathbf{Y}}_n)$ itself calculated with the to-be-separated observed signals as a regularization term for ILRMA. The regularization term in the conventional method has a simple structure because there is no correlation between the observed channels. Unfortunately, the proposed regularization term does not always have a simple structure because there exists a correlation between the observed channels. Thus, we cannot exactly follow the technique for designing the auxiliary function as in [14] for the newly obtained optimization problem. We find out the relationship between two matrices based on the positive semidefiniteness that allows us to extend the technique in [14], and we design an auxiliary function for the proposed cost function. On the basis of the MM algorithm, we derive an update rule by applying VCD [17] to the proposed auxiliary function.

B. Design of regularization term and update rules

By expanding (7), the proposed regularization term is also deformed into the following quadratic form of $w_{i,n}$:

$$\sum_n \mathcal{E}(\vec{\mathbf{Y}}_n) = \sum_{n,i,i'} \mathbf{w}_{i,n}^H \mathbf{G}_{ii'} \mathbf{w}_{i',n}, \quad (10)$$

where $\mathbf{G}_{ii'} \in \mathbb{C}^{M \times M}$ is dependent on not only the STFT conditions but also the to-be-separated observed signals. The specific representation of $\mathbf{G}_{ii'}$ will be omitted for space limitation. The cost function of the proposed method is the weighted sum of the cost function of ILRMA (6) and the regularization term (10) as follows:

$$\mathcal{L}_{\text{Prop}} = \mathcal{L} + \tilde{\gamma} \sum_{n,i,i'} \mathbf{w}_{i,n}^H \mathbf{G}_{ii'} \mathbf{w}_{i',n}, \quad (11)$$

where $\tilde{\gamma} \geq 0$ is the weight of the regularization term in the proposed method. Then, we derive update rules for the cost function (11). Since the regularization term (10) is independent of the NMF variables $\{t_{ik,n}\}$ and $\{v_{kj,n}\}$ as in [14], the update rules for the NMF variables are the same as in ILRMA. Therefore, it is sufficient to derive an update rule for $\{w_{i,n}\}$.

The cost function (11) has a correlation term between frequencies and this makes it difficult to optimize $\{w_{i,n}\}$ simultaneously. Therefore, we focus on one frequency bin i and treat the variables related to the other frequency bins as constants. That is, $\{w_{i,n}\}$ is updated one by one for each frequency bin in a block coordinate descent manner. The cost function is rewritten with regard to the frequency i as follows:

$$\begin{aligned} \mathcal{L}_{\text{Prop}} = & J \left[\sum_n w_{i,n}^H \frac{1}{J} \sum_j \frac{x_{ij} x_{ij}^H}{\sum_k t_{ik,n} v_{kj,n}} w_{i,n} - \log |\det \mathbf{W}_i|^2 \right] \\ & + \tilde{\gamma} \sum_n \begin{bmatrix} w_{i,n} \\ w_{i,n}^* \end{bmatrix}^H \begin{bmatrix} G_{ii} & G_{i\hat{i}} \\ G_{\hat{i}i} & G_{\hat{i}\hat{i}} \end{bmatrix} \begin{bmatrix} w_{i,n} \\ w_{i,n}^* \end{bmatrix} \\ & + 2\tilde{\gamma} \sum_n \text{Re} \left(\begin{bmatrix} \sum_{i' \neq i, \hat{i}} G_{i'i}^H w_{i',n} \\ \sum_{i' \neq i, \hat{i}} G_{i'i}^H w_{i',n} \end{bmatrix}^H \begin{bmatrix} w_{i,n} \\ w_{i,n}^* \end{bmatrix} \right) + \text{const.}, \end{aligned} \quad (12)$$

where $\text{Re}(\cdot)$ denotes a real part of a complex number and const. denotes the terms independent of $w_{i,n}$. The index \hat{i} is the frequency index opposite of i across the Nyquist frequency, where $i + \hat{i} = I + 2$ ($i \neq 1$). It is difficult to optimize (12) analytically because the terms $w_{i,n}^T G_{ii} w_{i,n}$ and $w_{i,n}^H G_{i\hat{i}} w_{i,n}^*$ are included in the cost function. Note that this problem does not occur in the case $i = 1$ (corresponding to 0 Hz) or $i = I/2$, because $w_{i,n}$ ($i = 1, I/2$) is real-valued to satisfy the constraint that the separated signal is real-valued in the time domain. The demixing matrix related to $i = 1, I/2$ can be updated using VCD without auxiliary functions. We design the auxiliary function without the terms $w_{i,n}^T G_{ii} w_{i,n}$ and $w_{i,n}^H G_{i\hat{i}} w_{i,n}^*$ in the case $i \neq 1, I/2$.

In this paper, we follow the strategy used in [14] to derive the auxiliary function on the basis of Claim IV.1 (Claim 1 in [14]) for (12).

Claim IV.1 (Claim 1 in [14]). For any integer $q \geq 1$, $w \in \mathbb{C}^q$, $\alpha \in \mathbb{C}^q$, and positive semidefinite matrix $\hat{G} \in \mathbb{C}^{q \times q}$, we have

$$\begin{aligned} w^H \hat{G} w & \leq w^H \hat{G}^+ w \\ & - 2\text{Re} \left(\alpha^H (\hat{G}^+ - \hat{G}) w \right) + \alpha^H (\hat{G}^+ - \hat{G}) \alpha, \end{aligned} \quad (13)$$

where \hat{G}^+ is a positive semidefinite matrix satisfying $\hat{G}^+ - \hat{G} \succeq \mathbf{O}$. Equality holds when $\alpha = w$.

Let \hat{G}^+ , \hat{G} , w , and α in Claim IV.1 be set as

$$\hat{G}^+ = \begin{bmatrix} G_{ii}^+ & \mathbf{O} \\ \mathbf{O} & G_{\hat{i}\hat{i}}^+ \end{bmatrix}, \quad \hat{G} = \begin{bmatrix} G_{ii} & G_{i\hat{i}} \\ G_{\hat{i}i} & G_{\hat{i}\hat{i}} \end{bmatrix},$$

$w = [w_{i,n}^T, w_{i,n}^{*T}]^T$, and $\alpha = [\alpha_{i,n}^T, \alpha_{i,n}^{*T}]^T$, where G_{ii}^+ and $G_{\hat{i}\hat{i}}^+$ are chosen so that $\hat{G}^+ - \hat{G} \succeq \mathbf{O}$ is satisfied. By

using (13), we can derive an auxiliary function without the terms $w_{i,n}^T G_{ii} w_{i,n}$ and $w_{i,n}^H G_{i\hat{i}} w_{i,n}^*$ for the second term of (12).

In [14], since $G_{i\hat{i}}^{(\text{ind})}$ becomes the diagonal matrix due to no correlation with respect to m , the problem of finding the matrix \hat{G}^+ is simplified to that of finding 2×2 matrices, and the following relation is used in [14]:

$$\begin{bmatrix} g_{ii}^{(\text{ind})} + |g_{i\hat{i}}^{(\text{ind})}| & 0 \\ 0 & g_{\hat{i}\hat{i}}^{(\text{ind})} + |g_{i\hat{i}}^{(\text{ind})}| \end{bmatrix} - \begin{bmatrix} g_{ii}^{(\text{ind})} & g_{i\hat{i}}^{(\text{ind})} \\ g_{\hat{i}i}^{(\text{ind})} & g_{\hat{i}\hat{i}}^{(\text{ind})} \end{bmatrix} \succeq \mathbf{O}. \quad (14)$$

Since $G_{i\hat{i}}$ has nondiagonal terms due to the correlation regarding m , we cannot, however, simply follow this derivation in [14]. For this problem, we find the following relation as

$$\begin{bmatrix} G_{ii} + (G_{i\hat{i}} G_{\hat{i}i})^{\frac{1}{2}} & \mathbf{O}_{M \times M} \\ \mathbf{O}_{M \times M} & G_{\hat{i}\hat{i}} + (G_{\hat{i}i} G_{i\hat{i}})^{\frac{1}{2}} \end{bmatrix} - \begin{bmatrix} G_{ii} & G_{i\hat{i}} \\ G_{\hat{i}i} & G_{\hat{i}\hat{i}} \end{bmatrix} \succeq \mathbf{O}, \quad (15)$$

which is an extension of (14). This fact can be proved by singular value decomposition, but we omit the details for space limitation.

By using Claim IV.1 and (15), we design the total auxiliary function for (12) as

$$\begin{aligned} \mathcal{L}_{\text{Prop}}^+ = & J [w_{i,n}^H D_{i,n}^+ w_{i,n} + 2\text{Re} (b_{i,n}^{+H} w_{i,n}) - \log |\det \mathbf{W}_i|^2] \\ & + \text{const.}, \end{aligned} \quad (16)$$

where const. denotes the terms independent of $w_{i,n}$ and $D_{i,n}^+$ and $b_{i,n}^+$ are respectively defined as

$$\begin{aligned} D_{i,n}^+ & = \frac{1}{J} \sum_j \frac{x_{ij} x_{ij}^H}{\sum_k t_{ik,n} v_{kj,n}} + \frac{2\tilde{\gamma}}{J} \left(G_{ii} + (G_{i\hat{i}} G_{\hat{i}i})^{\frac{1}{2}} \right), \quad (17) \\ b_{i,n}^+ & = \frac{2\tilde{\gamma}}{J} \left[\sum_{i' \neq i, \hat{i}} G_{i'i}^H w_{i',n} - (G_{i\hat{i}} G_{\hat{i}i})^{\frac{1}{2}} \alpha_{i,n} + G_{\hat{i}\hat{i}}^H \alpha_{i,n}^* \right]. \end{aligned} \quad (18)$$

Here, $\alpha_{i,n} \in \mathbb{C}^M$ is an auxiliary variable and the equality for the auxiliary function holds when $\alpha_{i,n} = w_{i,n}$. Since (16) is the sum of a quadratic form, a linear term, and a log-determinant term, an update rule using VCD [17] can be obtained as follows:

$$\hat{\zeta}_{i,n} \leftarrow (\mathbf{W}_i D_{i,n}^+)^{-1} e_n, \quad (19)$$

$$\hat{\zeta}_{i,n} \leftarrow (D_{i,n}^+)^{-1} b_{i,n}^+, \quad (20)$$

$$\chi_{i,n} \leftarrow \zeta_{i,n}^H D_{i,n}^+ \zeta_{i,n}, \quad (21)$$

$$\hat{\chi}_{i,n} \leftarrow \zeta_{i,n}^H D_{i,n}^+ \hat{\zeta}_{i,n}, \quad (22)$$

$$w_{i,n} \leftarrow \begin{cases} \frac{\zeta_{i,n}}{\sqrt{\chi_{i,n}}} - \hat{\zeta}_{i,n} & (\hat{\chi}_{i,n} = 0), \\ \frac{\hat{\chi}_{i,n}}{2\chi_{i,n}} \left(1 - \sqrt{1 + \frac{4\chi_{i,n}}{|\hat{\chi}_{i,n}|^2}} \right) \zeta_{i,n} - \hat{\zeta}_{i,n} & (\text{otherwise}), \end{cases} \quad (23)$$

where $e_n \in \mathbb{C}^N$ is a vector whose n th component is one and the others are zero. This update rule guarantees the monotonic nonincrease in the cost function (11). All the update rules in

TABLE I: Averages of SDR improvements over 100 trials for each method

ILRMA	Consist. ILRMA	ILRMA w/ OSICR			Prop. w/o temp.			Prop. w/ temp.		
		$\beta = 10^{-6}$	$\beta = 10^{-5}$	$\beta = 10^{-4}$	$\gamma = 10^3$	$\gamma = 10^4$	$\gamma = 10^5$	$\gamma_0 = 10^3$	$\gamma_0 = 10^4$	$\gamma_0 = 10^5$
9.6	11.0	9.6	9.5	8.7	10.5	10.8	5.3	10.4	11.2	8.7

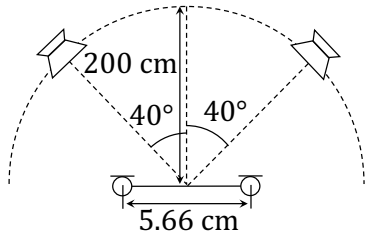


Fig. 2: Recording conditions of impulse response.

the proposed method including those for the NMF variables $\{t_{ik,n}\}$ and $\{v_{kj,n}\}$ guarantee the monotonic nonincrease in the cost function (11).

C. Weight of regularization

In this paper, to prevent the effect of the regularization term from varying with the power and length of the observed signals, we normalize the regularization weight $\tilde{\gamma}$ as

$$\tilde{\gamma} = \gamma \cdot \frac{J}{\sum_{i,j} \|\mathbf{x}_{ij}\|_2^2} \quad (24)$$

and design $\gamma \geq 0$ instead.

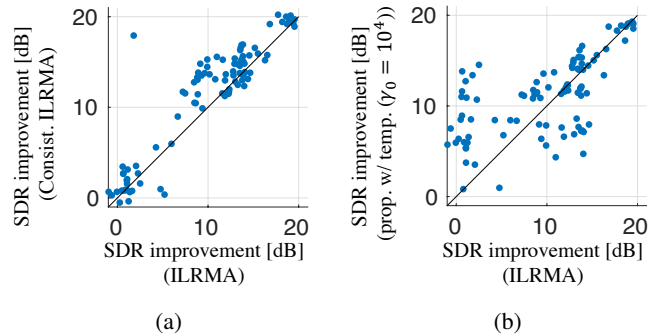
In some regularization-based methods, the scheduling of the weight parameter is often performed. We use ‘‘tempering,’’ in which the weight parameter is reduced with increasing iteration number, because the regularization to be consistent may conflict with the separation. In the tempering of γ , we set the the parameter γ at the c th iteration as $\max\{0, \gamma_0 - \gamma_0(c-1)/(C_0-1)\}$. Note that the parameter γ is γ_0 at the first iteration, decreases linearly thereafter, and is 0 after the C_0 th iteration.

V. SIMULATION EXPERIMENT

A. Experimental conditions

To confirm the performance of the proposed method, a simulation experiment was conducted using two sources and two microphones ($M = N = 2$). Ten pairs of music signals from SiSEC2011[19] were prepared as dry sources. They were convolved with the impulse responses E2A ($T_{60} = 300$ ms) in the RWCP database [20]. The sources and microphones were placed as shown in Fig. 2. Both the dry sources and the impulse responses were resampled at 16 kHz. The STFT was computed with a window length of 256 ms, a hopsize of 64 ms, and the Hamming window.

We compared ILRMA [7], Consist. ILRMA [13], ILRMA with OICR (ILRMA w/ OICR) [14], the proposed method with constant γ (prop. w/o temp.), and the proposed method with


 Fig. 3: Pairplots for SDR improvements of ILRMA and (a) Consist. ILRMA or (b) prop. w/ temp. ($\gamma_0 = 10^4$).

the tempering of γ (prop. w/ temp.). The number of bases for NMF K was 10. The demixing matrix was initialized to the identity matrix and the NMF variables were initialized with uniformly distributed random values over $(0, 1)$. Ten trials were performed using different random seeds, i.e., we conducted experiments under 100 conditions (10 pairs of sources \times 10 random seeds). The number of iterations was 400. We used $\beta = 10^{-6}, 10^{-5}, 10^{-4}$ in ILRMA w/ OICR, $\gamma = 10^3, 10^4, 10^5$ in prop. w/o temp., and $\gamma_0 = 10^3, 10^4, 10^5$ and $C_0 = 50$ in prop. w/ temp. The source-to-distortion ratio (SDR) [21] improvement was used as a measure of the separation performance.

B. Results and discussion

Table I shows the average of SDR improvements. Prop. w/o temp. and prop. w/ temp. improved the separation performance compared with ILRMA except for $\gamma = 10^5$ in prop. w/o temp. and $\gamma_0 = 10^5$ in prop. w/ temp., in which the initial value of γ was very large. In particular, prop. w/ temp. ($\gamma_0 = 10^4$) showed an average SDR improvement of 0.2 dB higher than Consist. ILRMA. In this experiment, ILRMA w/ OICR showed a lower separation performance than ILRMA, but this may be due to differences between the evaluation data used in [14] and in this study. Fig. 3 shows the pairplots for the SDR improvements of ILRMA with Consist. ILRMA and prop. w/ temp. ($\gamma_0 = 10^4$). Prop. w/ temp. ($\gamma_0 = 10^4$) showed higher SDR improvements than Consist. ILRMA in many cases when the SDR improvements of ILRMA were lower than 5 dB. The proposed regularization term seems to prevent the ILRMA from converging to poor solutions.

VI. CONCLUSION

Several methods have been proposed with the spectrogram consistency considered in ILRMA. One of them is a

regularization-based method that induces the separated signals to be consistent. Although the regularization term independent of the observed signals is used in the conventional method, we designed the regularization term that depends on the observed signals. For the optimization of the new cost function with the new regularization term, we designed the auxiliary function that is an extension of that in the conventional method. By applying VCD to the auxiliary function, we derived the update rule that guarantees the monotonic nonincrease in the cost function. A numerical experiment confirmed that the proposed method with the tempering of the weight parameter of the regularization showed a higher performance than the conventional methods.

ACKNOWLEDGMENT

The authors would like to thank Rintaro Ikeshita and Tomohiro Nakatani (NTT) for useful discussions.

REFERENCES

- [1] P. Comon, "Independent component analysis, a new concept?" *Signal Processing*, vol. 36, no. 3, pp. 287–314, 1994.
- [2] P. Smaragdis, "Blind separation of convolved mixtures in the frequency domain," *Neurocomputing*, vol. 22, no. 1, pp. 21–34, 1998.
- [3] S. Kurita, H. Saruwatari, S. Kajita, K. Takeda, and F. Itakura, "Evaluation of blind signal separation method using directivity pattern under reverberant conditions," in *Proc. ICASSP*, vol. 5, 2000, pp. 3140–3143.
- [4] H. Sawada, R. Mukai, S. Araki, and S. Makino, "A robust and precise method for solving the permutation problem of frequency-domain blind source separation," *IEEE Trans. SAP*, vol. 12, no. 5, pp. 530–538, 2004.
- [5] H. Sawada, S. Araki, and S. Makino, "Measuring dependence of bin-wise separated signals for permutation alignment in frequency-domain bss," in *Proc. ISCAS*, 2007, pp. 3247–3250.
- [6] T. Kim, T. Eltoft, and T.-W. Lee, "Independent vector analysis: An extension of ICA to multivariate components," in *Proc. ICA*, 2006, pp. 165–172.
- [7] 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. ASLP*, vol. 24, no. 9, pp. 1626–1641, 2016.
- [8] D. Kitamura, N. Ono, and H. Saruwatari, "Experimental analysis of optimal window length for independent low-rank matrix analysis," in *Proc. EUSIPCO*, 2017, pp. 1170–1174.
- [9] D. Griffin and J. Lim, "Signal estimation from modified short-time fourier transform," *IEEE Trans. ASSP*, vol. 32, no. 2, pp. 236–243, 1984.
- [10] J. Le Roux and E. Vincent, "Consistent Wiener filtering for audio source separation," *IEEE SP Letters*, vol. 20, no. 3, pp. 217–220, 2012.
- [11] N. Perraudin, P. Balazs, and P. L. Søndergaard, "A fast Griffin–Lim algorithm," in *Proc. WASPAA*, 2013, pp. 1–4.
- [12] K. Yatabe, "Consistent ICA: Determined BSS meets spectrogram consistency," *IEEE SP Letters*, vol. 27, pp. 870–874, 2020.
- [13] D. Kitamura and K. Yatabe, "Consistent independent low-rank matrix analysis for determined blind source separation," *EURASIP JASP*, vol. 2020, no. 46, pp. 1–35, 2020.
- [14] S. Misawa, N. Takamune, K. Yatabe, D. Kitamura, and H. Saruwatari, "Blind source separation using independent low-rank matrix analysis with spectrogram-consistency regularization," in *Proc. APSIPA*, 2023, pp. 1035–1042.
- [15] D. D. Lee and H. S. Seung, "Learning the parts of objects by non-negative matrix factorization," *Nature*, vol. 401, no. 6755, pp. 788–791, 1999.
- [16] D. R. Hunter and K. Lange, "Quantile regression via an MM algorithm," *JCGS*, vol. 9, no. 1, pp. 60–77, 2000.
- [17] "Independent deeply learned matrix analysis with automatic selection of stable microphone-wise update and fast sourcewise update of demixing matrix," *Signal Processing*, vol. 178, pp. 1–12, 2021.
- [18] N. Ono, "Stable and fast update rules for independent vector analysis based on auxiliary function technique," in *Proc. WASPAA*, 2011, pp. 189–192.
- [19] S. Araki, F. Nesta, E. Vincent, *et al.*, "The 2011 signal separation evaluation campaign (SiSEC2011): - Audio source separation -," in *Proc. LVA/ICA*, 2012, pp. 414–422.
- [20] 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 *Proc. LREC*, 2000, pp. 965–968.
- [21] E. Vincent, R. Gribonval, and C. Févotte, "Performance measurement in blind audio source separation," *IEEE Trans. ASLP*, vol. 14, no. 4, pp. 1462–1469, 2006.