

UNIFIED APPROACH FOR AUDIO SOURCE SEPARATION WITH MULTICHANNEL FACTORIAL HMM AND DOA MIXTURE MODEL

Takuya Higuchi¹⁾ and Hirokazu Kameoka^{1),2)}

¹⁾Graduate School of Information Science and Technology, The University of Tokyo
7-3-1 Hongo, Bunkyo-ku, Tokyo 113-8656, Japan

²⁾NTT Communication Science Laboratories, Nippon Telegraph and Telephone Corporation
3-1 Morinosatowakamiya, Atsugi, Kanagawa 243-0198, Japan
{higuchi,kameoka}@hil.t.u-tokyo.ac.jp

ABSTRACT

We deal with the problems of blind source separation, dereverberation, audio event detection and direction-of-arrival (DOA) estimation. We previously proposed a generative model of multichannel signals called the multichannel factorial hidden Markov model, which allows us to simultaneously solve these problems through a joint optimization problem formulation. In this approach, we modeled the spatial correlation matrix of each source as a weighted sum of the spatial correlation matrices corresponding to all possible DOAs. However, it became clear through real environment experiments that the estimate of the spatial correlation matrix tended to deviate from the actual correlation matrix since the plane wave assumption does not hold due to reverberation and noise components. To handle such deviations, we propose introducing a prior distribution over the spatial correlation matrices called the DOA mixture model instead of using the weighted sum model. The experiment showed that the proposed method provided 1.94 [dB] improvement compared with our previous method in terms of the signal-to-distortion ratios of separated signals.

Index Terms— Blind source separation, voice activity detection, dereverberation, DOA estimation

1. INTRODUCTION

Blind source separation (BSS) refers to a technique for separating out individual source signals from microphone array inputs when the transfer characteristics between the sources and microphones are unknown. To solve the BSS problem, it is generally necessary to make some assumptions about the sources, and formulate an appropriate optimization problem based on criteria designed according to those assumptions. For example, if the observed signals outnumber the sources, we can employ independent component analysis (ICA) [1] by assuming that the sources are statistically independent of each other. For monaural source separation, one successful approach involves applying non-negative matrix factorization (NMF) to the magnitude spectrogram of a mixture signal, interpreted as a non-negative matrix [2, 3]. Up to now, several attempts have been made to extend this approach to a multichannel case in order to allow us to use the spatial information as an additional clue to source separation [4, 5]. Moreover, we previously proposed the multichannel factorial hidden Markov model [6], where we used the information of audio events for source separation and simultaneously solved the problems of BSS and audio event detection. Furthermore,

we extended the model for dereverberation [7], by approximating the multichannel observed signal recorded in a reverberant condition as a form of a convolution of the frequency array response and the source signal in the time-frequency domain. Thus, we modeled the impulse response out of the frame of STFT by introducing a time sequence of frequency response arrays in the time-frequency domain. Furthermore, we used direction-of-arrivals (DOAs) of sources as a clue to estimation of spatial correlation matrices [8]. Generally, it is known that the spatial correlation matrix of the direct sound of a point source has a certain structure described by the DOA. In the model proposed in [8], we used a weighted sum of matrices parametrized by all possible DOAs, same as in [9], and estimates of the weight variables corresponded to estimates of spatial correlation matrices at all frequencies. However, in practical situations, a spatial correlation matrix often deviates from a theoretical structure because of reverberation and noises within the frame of the STFT. Therefore, in this paper, we apply the DOA mixture model [10] for modeling the generative process of a spatial correlation matrix, and design a prior distribution over a spatial correlation matrix based on DOAs. In the present model, we estimate a spatial correlation matrix at each frequency bin based on DOAs. Through the parameter inference of our new generative model, we can simultaneously perform source separation, source activity detection, dereverberation and DOA estimation based on a unified maximum likelihood criterion.

2. MULTICHANNEL FACTORIAL HIDDEN MARKOV MODEL

2.1. Mixing model with convolutive mixture approximation [7]

First we consider a situation where I source signals are recorded by M microphones. In a reverberant condition, the length of the impulse responses are relatively long and so an instantaneous mixture approximation is not always true. Therefore we approximately express the observed signals as a form of a convolution of the frequency array response and the source signal in the time-frequency domain.

$$\mathbf{y}(\omega_k, t_l) \approx \sum_{i=1}^I \sum_{\tau=0}^{\mathcal{T}} \mathbf{a}_i(\omega_k, t_\tau) s_i(\omega_k, t_l - t_\tau). \quad (1)$$

Here, let $y_m(\omega_k, t_l) \in \mathbb{C}$ be the short-time Fourier transform (STFT) component observed at the m -th microphone, and $s_i(\omega_k, t_l) \in \mathbb{C}$ be the STFT component of the i -th source signal. $1 \leq k \leq K$ and $1 \leq l \leq L$ are the frequency

This work was supported by JSPS KAKENHI Grant Number 26730100.

and time indices, respectively. $\mathbf{a}_i(\omega_k, t_\tau)$ denotes the frequency array response for source i at frequency ω_k and time t_τ . $0 \leq \tau \leq T$ is the time index of the frequency array response in the time-frequency domain. Note that $\mathbf{a}_i(\omega_k, t_{1:T})$ denote the frequency array responses which correspond to the impulse responses out of the frames of the STFT. For convenience of notation, we hereafter use subscripts k, l and τ to indicate ω_k, t_l and t_τ respectively.

2.2. Generative modeling of source signals based on audio events [6]

We now describe the generative process of the source signal $s_{i,k,l}$ based on its audio event. First, we assume each signal has a specific spectrum and utilize NMF for $\sigma_{i,k,l}^2$, which is an expected value of power of $s_{i,k,l}$. $\sigma_{i,k,l}^2$ is factorized as

$$\sigma_{i,k,l}^2 = w_{i,k} h_{i,l}, \quad (2)$$

where $w_{i,k}$ and $h_{i,l}$ are non-negative variables. The generative model of $s_{i,k,l}$ is written as

$$s_{i,k,l} | w_{i,k}, h_{i,l} \sim \mathcal{N}_{\mathbb{C}}(s_{i,k,l}; 0, w_{i,k} h_{i,l}), \quad (3)$$

conditioned on $w_{i,k}$ and $h_{i,l}$. Under the condition of Eq. (3), the generative model of the multichannel observed signal is equivalent to the method proposed in [4, 5].

In many cases, a sound signal consists of several spectra. Therefore, first we introduce latent variable $z_{i,l}^{(q)} \in \{1, \dots, Q\}$ to denote a state of i -th source's spectrum at time l , and the state sequence $z_{i,1}^{(q)}, \dots, z_{i,L}^{(q)}$ follows a Markov chain:

$$z_{i,l}^{(q)} | z_{i,l-1}^{(q)} \sim \text{Categorical}(z_{i,l}^{(q)}; \boldsymbol{\rho}_{i,z_{i,l-1}^{(q)}}^{(q)}), \quad (4)$$

where $\text{Categorical}(x; \mathbf{y}) = y_x$, $\boldsymbol{\rho}_{i,q}^{(q)} = (\rho_{i,q,1}^{(q)}, \dots, \rho_{i,q,Q}^{(q)})$ denotes the transition probability of state q to each state $1, \dots, Q$, and $\boldsymbol{\rho}_i^{(q)} = (\rho_{i,q,q'}^{(q)})_{Q \times Q}$ denotes the transition matrix. Then the power spectrum of the i -th source at time l is assumed to be determined according to $z_{i,l}^{(q)}$ and $w_{i,k,q}$ denotes a spectrum of i -th source at state q . Thus, the generative model of $s_{i,k,l}$ is rewritten as

$$s_{i,k,l} | w_{i,k,1:Q}, h_{i,l}, z_{i,l}^{(q)} \sim \mathcal{N}_{\mathbb{C}}(s_{i,k,l}; 0, w_{i,k,z_{i,l}^{(q)}} h_{i,l}). \quad (5)$$

Next, the source's power also tend to differ according to the source's state such as "active" or "inactive." Therefore we first introduce latent variable $z_{i,l}^{(j)} \in \{1, \dots, J\}$ to denote a state of i -th source's power at time l , and the state sequence $z_{i,1}^{(j)}, \dots, z_{i,L}^{(j)}$ follows a Markov chain:

$$z_{i,l}^{(j)} | z_{i,l-1}^{(j)} \sim \text{Categorical}(z_{i,l}^{(j)}; \boldsymbol{\rho}_{i,z_{i,l-1}^{(j)}}^{(j)}). \quad (6)$$

Then we assume $h_{i,l}$ follows a gamma distribution which has different parameters according to $z_{i,l}^{(j)}$,

$$h_{i,l} | z_{i,l}^{(j)} \sim \text{Gamma}(h_{i,l}; \alpha_{z_{i,l}^{(j)}} \beta_{z_{i,l}^{(j)}}), \quad (7)$$

where $\alpha_{1:J}$ and $\beta_{1:J}$ are the shape and scale parameters of a gamma distribution, and $\text{Gamma}(x; \alpha, \beta) = \frac{x^{\alpha-1} e^{-x/\beta}}{\Gamma(\alpha) \beta^\alpha}$.

As we want $h_{i,l}$ to take a small value when $z_{i,l}^{(j)}$ is the "inactive" (i.e., silent) state, we set the hyperparameters of the

gamma distribution of that state so that it becomes a sparsity-inducing distribution. As regards the gamma distributions of the remaining states, we consider setting the hyperparameters so that they are like uniform distributions.

2.3. Generative process of a spatial correlation matrix based on DOA mixture model

If we assume that a source is far from the microphones, the frequency array response has a certain structure depending on DOA. Therefore, we can express a spatial correlation matrix by using the DOAs. Since a spatial correlation matrix $\mathbf{C}_{i,k,\tau}$ is defined as $\mathbf{C}_{i,k,\tau} = \mathbf{a}_{i,k,\tau} \mathbf{a}_{i,k,\tau}^H$ by using $\mathbf{a}_{i,k,\tau}$, specifically, with $M = 2$ microphones, the spatial correlation matrix for a source at direction θ such that $0 \leq \theta \leq \pi$ is defined as a function of ω depending on θ

$$\mathbf{J}(\theta, \omega) = \begin{bmatrix} 1 \\ e^{j\omega B \cos \theta / C} \end{bmatrix} [1 \ e^{j\omega B \cos \theta / C}]^*, \quad (8)$$

where j is the imaginary unit, B [m] is the distance between the two microphones, and C [m/s] is the speed of sound. If the DOA θ_i of source i is known, the spatial correlation matrix for the direct wave should be equal to $\mathbf{J}(\theta_i, \omega_k)$. However, $\mathbf{C}_{i,k,0}$ would not be equal to $\mathbf{J}(\theta_i, \omega_k)$ because of reverberant components and noises within the frame of the STFT. One way is modeling a spatial correlation matrix by using a weighted sum of $\mathbf{J}(\theta_i, \omega_k)$ of all possible DOAs, same as in [8, 9]. In that model, estimates of the weight variables corresponded to estimates of spatial correlation matrices at all frequencies, and so it would not be enough for modeling the deviation of spatial correlation matrices. Therefore, in this paper, we incorporate the DOA mixture model [10] into the multichannel factorial hidden Markov model, and estimate a spatial correlation matrix at each frequency bin using the information of DOAs. First, we now introduce a discrete set of O possible directions, $\vartheta_1, \dots, \vartheta_O$, which are all assumed to be constants. We then assume that each source signal propagates from one of these directions. For each source i , an index $z_i^{(o)}$ of direction is drawn according to a categorical distribution $\boldsymbol{\rho}_i^{(o)} = (\rho_{i,1}, \dots, \rho_{i,O})$

$$z_i^{(o)} | \boldsymbol{\rho}_i^{(o)} \sim \text{Categorical}(z_i^{(o)}; \boldsymbol{\rho}_i^{(o)}), \quad (9)$$

$$\theta_i = \vartheta_{z_i^{(o)}}. \quad (10)$$

Then, a spatial correlation matrix follows a whishart distribution:

$$\mathbf{C}_{i,k,0} | z_i^{(o)} \sim \mathcal{W}_{\mathbb{C}}(\mathbf{C}_{i,k,0}; \gamma, \mathbf{J}_{\vartheta_{z_i^{(o)}}} + \epsilon \mathbf{I}), \quad (11)$$

where $\mathcal{W}_{\mathbb{C}}(\mathbf{X}; \gamma, \mathbf{Y}) \propto |\mathbf{X}|^{(\gamma-M)/2} \exp(-\frac{1}{2} \text{tr}(\mathbf{X} \mathbf{Y}^{-1}))$ and $\gamma (\geq M + 1)$ is a hyperparameter of the whishart distribution. Here \mathbf{I} is an identity matrix and ϵ is a small value, and they enable inverse operation.

Our overall generative model is given by Eqs. (4), Eqs. (6), Eqs. (7) Eqs. (9), Eqs. (11) and

$$\mathbf{y}_{k,l} | \mathbf{a}_{1:I,k,0:T}, w_{1:I,k,1:Q}, h_{1:I,l-T:l}, z_{1:I,l-T:l}^{(q)} \sim \mathcal{N}_{\mathbb{C}}(\mathbf{y}_{k,l}; 0, \sum_{i,\tau} \mathbf{C}_{i,k,\tau} w_{i,k,z_{i,l-\tau}^{(q)}} h_{i,l-\tau}), \quad (12)$$

conditioned on $\mathbf{a}_{1:I,k,0:T}$, $w_{1:I,k,1:Q}$, $h_{1:I,l-T:l}$, $z_{1:I,l-T:l}^{(q)}$. Parameter estimation of the generative model allows us to solve the problems of source separation, source activity detec-

tion, dereverberation and DOA estimation based on a unified maximum likelihood criterion.

3. PARAMETER ESTIMATION ALGORITHM BASED ON AUXILIARY FUNCTION METHOD

In this section, we describe a parameter estimation algorithm for our generative model based on an auxiliary function method. The random variables of interest in our model are $\mathbf{W} = w_{1:I,1:K,1:Q}$, $\mathbf{H} = h_{1:I,1:L}$ and $\mathbf{C} = \mathbf{C}_{1:I,1:K,0:\tau}$. We denote the entire set of the above parameters as Θ . $\mathbf{Z}^{(q)} = z_{1:I,1:L}$, $\mathbf{Z}^{(j)} = z_{1:I,1:L}$ and $\mathbf{Z}^{(o)} = z_{1:I}$ are latent variables. In the following, $\rho^{(q)}$, $\rho^{(j)}$, $\rho^{(o)}$, α and β are constants that are determined experimentally.

The objective function $L(\Theta)$ is defined as $L(\Theta) = \log p(\Theta|\mathbf{Y})$, where $\mathbf{Y} = \mathbf{y}_{1:K,1:L}$ is a set consisting of the time-frequency components of observed multichannel signals. By using the conditional distributions defined in Sec. 2, we can write down $\log p(\Theta|\mathbf{Y})$. However, $L(\Theta)$ is difficult to maximize analytically, and so we derive an efficient parameter estimation algorithm based on an auxiliary function method.

The optimization problem of maximizing $L(\Theta)$ with respect to Θ is difficult to solve analytically. However, we can invoke an auxiliary function approach to derive an iterative algorithm that searches for the estimate of Θ , as with [5]. To apply an auxiliary function approach to the current optimization problem, the first step is to construct an auxiliary function $L^+(\Theta, \Lambda)$ satisfying $L(\Theta) = \max_{\Lambda} L^+(\Theta, \Lambda)$. We refer to Λ as an auxiliary variable. It can then be shown that $L(\Theta)$ is non-decreasing under the updates $\Theta \leftarrow \arg\max_{\Theta} L^+(\Theta, \Lambda)$ and $\Lambda \leftarrow \arg\max_{\Lambda} L^+(\Theta, \Lambda)$. The proof of this shall be omitted owing to space limitations. Thus, $L^+(\Theta, \Lambda)$ should be designed as a function that can be maximized analytically with respect to Θ and Λ . Such a function can be constructed as follows.

$$\begin{aligned} L(\Theta) &\geq L^+(\Theta, \Lambda) \\ &= -\frac{1}{2} \sum_{k,l} \left[\sum_{i,q,\tau} \lambda_{q,i,l-\tau}^{(q)} \right. \\ &\quad \left(\frac{\text{tr}(\mathbf{y}_{k,l} \mathbf{y}_{k,l}^H \mathbf{R}_{i,k,l,q,\tau} \mathbf{C}_{i,k}^{-1} \mathbf{R}_{i,k,l,q,\tau})}{w_{i,k,q} h_{i,l-\tau}} \right. \\ &\quad \left. + \text{tr}(\mathbf{U}_{k,l}^{-1} \mathbf{C}_{i,k}) w_{i,k,q} h_{i,l-\tau} \right) + \log |\mathbf{U}_{k,l}| \Big] \\ &\quad + \sum_{j,i,l} \lambda_{j,i,l}^{(j)} \left((\alpha_j - 1) \log h_{i,l} - h_{i,l} / \beta_j \right. \\ &\quad \left. - \alpha_j \log \beta_j - \log \Gamma(\alpha_j) \right) \\ &\quad + \sum_{i,k,o} d_{i,o} \left[-\frac{1}{2} \text{tr}(\mathbf{C}_{i,k,0} (\mathbf{J}_{k,o} + \epsilon \mathbf{I})^{-1}) \right. \\ &\quad \left. + \frac{\gamma - M}{2} \left(\text{tr}(\tilde{\mathbf{U}}_{i,k}^{-1} \mathbf{C}_{i,k,0}) \right. \right. \\ &\quad \left. \left. + \log |\tilde{\mathbf{U}}_{i,k}| \right) \right] + \sum_{q,i,l} \lambda_{q,i,l}^{(q)} \log p(\mathbf{Z}^{(q)}) \end{aligned}$$

$$+ \sum_{j,i,l} \lambda_{j,i,l}^{(j)} \log p(\mathbf{Z}^{(j)}) + C_{L^+}, \quad (13)$$

where $\mathbf{R}_{i,k,l,\tau,q}$, $\mathbf{U}_{k,l}$, $\tilde{\mathbf{U}}_{i,k}$, $\lambda_{q,i,l}^{(q)}$, $\lambda_{j,i,l}^{(j)}$ and $d_{i,o}$ are auxiliary variables. $\mathbf{R}_{i,k,l,\tau,q}$, $\mathbf{U}_{k,l}$ and $\tilde{\mathbf{U}}_{i,k}$ satisfy Hermitian positive definiteness and $\sum_{i,\tau} \mathbf{R}_{i,k,l,\tau,q} = \mathbf{I}$. $\lambda_{q,i,l}^{(q)}$, $\lambda_{j,i,l}^{(j)}$ and $d_{i,o}$ satisfy $\sum_q \lambda_{q,i,l}^{(q)} = 1$, $\sum_j \lambda_{j,i,l}^{(j)} = 1$ and $\sum_o d_{i,o} = 1$. We denote the set of the auxiliary variables as Λ . C_{L^+} is the sum of constant terms. $\text{tr}(\cdot)$ is the trace of a matrix. The equality $L(\Theta) = L^+(\Theta, \Lambda)$ is satisfied when

$$\mathbf{R}_{i,k,l,\tau,q} = \mathbf{C}_{i,k,\tau} w_{i,k,q} h_{i,l-\tau} \hat{\mathbf{X}}_{k,l}^{-1}, \quad (14)$$

$$\mathbf{U}_{k,l} = \hat{\mathbf{X}}_{k,l}, \quad (15)$$

$$\tilde{\mathbf{U}}_{i,k} = \mathbf{C}_{i,k,0}, \quad (16)$$

with respect to $\mathbf{R}_{i,k,l,\tau,q}$, $\mathbf{U}_{k,l}$ and $\tilde{\mathbf{U}}_{i,k}$. With respect to $\lambda_{q,i,l}^{(q)}$,

$$\lambda_{q,i,l}^{(q)} = F_{q,i,l}^{(q)} B_{q,i,l}^{(q)} / \sum_q F_{q,i,l}^{(q)} B_{q,i,l}^{(q)}, \quad (17)$$

$$F_{q,i,l}^{(q)} = p(\Theta | z_{i,l}^{(q)} = q) \sum_{q'} F_{q,i,l-1}^{(q)} \rho_{i,q',q}^{(q)}, \quad (18)$$

$$B_{q,i,l}^{(q)} = \sum_{q'} B_{q',i,l+1}^{(q)} p(\Theta | z_{i,l+1}^{(q)} = q') \rho_{i,q,q'}^{(q)}, \quad (19)$$

where

$$\begin{aligned} &p(\Theta | z_{i,l}^{(q)} = q) \\ &\propto \exp \left[-\frac{1}{2} \sum_{i,k,\tau} \left(\text{tr}(\mathbf{y}_{k,l+\tau} \mathbf{y}_{k,l+\tau}^H \mathbf{R}_{i,k,l+\tau,\tau,q}) \right. \right. \\ &\quad \left. \left. \mathbf{C}_{i,k,\tau}^{-1} \mathbf{R}_{i,k,l+\tau,\tau,q} \right) / w_{i,k,q} h_{i,l} \right. \\ &\quad \left. + \text{tr}(\mathbf{U}_{k,l+\tau}^{-1} \mathbf{C}_{i,k,\tau}) w_{i,k,q} h_{i,l} \right]. \quad (20) \end{aligned}$$

With respect to $\lambda_{j,i,l}^{(j)}$,

$$\lambda_{j,i,l}^{(j)} = F_{j,i,l}^{(j)} B_{j,i,l}^{(j)} / \sum_j F_{j,i,l}^{(j)} B_{j,i,l}^{(j)}, \quad (21)$$

$$F_{j,i,l}^{(j)} = p(\Theta | z_{i,l}^{(j)} = j) \sum_{j'} F_{j,i,l-1}^{(j)} \rho_{i,j',j}^{(j)}, \quad (22)$$

$$B_{j,i,l}^{(j)} = \sum_{j'} B_{j',i,l+1}^{(j)} p(\Theta | z_{i,l+1}^{(j)} = j') \rho_{i,j,j'}^{(j)}, \quad (23)$$

where

$$\begin{aligned} &p(\Theta | z_{i,l}^{(j)} = j) \\ &\propto \exp \left((\alpha_j - 1) \log h_{i,l} - h_{i,l} / \beta_j \right. \\ &\quad \left. - \alpha_j \log \beta_j - \log \Gamma(\alpha_j) \right). \quad (24) \end{aligned}$$

With respect to $d_{i,o}$,

$$\begin{aligned} &d_{i,o} = p(\Theta | z_i^{(o)} = o) \\ &\propto \exp \left(-\frac{1}{2} \text{tr}(\mathbf{C}_{i,k,0} (\mathbf{J}_{\vartheta_o,k} + \epsilon \mathbf{I})^{-1}) \right), \quad (25) \end{aligned}$$

Table 1. The experimental condition of parameters

Q	O	α	β	γ	ϵ	ϵ'
15	30	1	$\beta_1 = 10^{-3}, \beta_2 = 10^{10}$	3	10^{-3}	10^{-3}

and $\sum_o d_{i,o} = 1$.

By setting the partial differential of L^+ at zero, we can obtain the update rule of each parameter in Θ . Specifically, the partial differential of L^+ with respect to $C_{i,k,0}$ is given by

$$\begin{aligned} & \frac{\partial L^+}{\partial C_{i,k,0}} \\ &= \frac{1}{2} \sum_{l,q} \lambda_{q,i,l} \left(\frac{C_{i,k,0}^{-1} \mathbf{R}_{i,k,l,0,q} \mathbf{y}_{k,l} \mathbf{y}_{k,l}^H \mathbf{R}_{i,k,l,0,q} C_{i,k,0}^{-1}}{w_{i,k,q} h_{i,l}} \right. \\ & \quad \left. - U_{k,l}^{-1} w_{i,k,q} h_{i,l} \right) + \sum_o d_{i,o} \left(-\frac{1}{2} (\mathbf{J}_{\theta_o,k} + \epsilon \mathbf{I})^{-1} \right) \\ & \quad + \frac{\gamma - M}{2} \tilde{U}_{i,k}^{-1}. \end{aligned} \quad (26)$$

By setting this at zero, we obtain an algebraic Riccati equation:

$$C_{i,k,0} \mathbf{A}_{i,k,0} C_{i,k,0} = \mathbf{B}_{i,k,0}, \quad (27)$$

where

$$\begin{aligned} \mathbf{A}_{i,k,0} &= \sum_{l,q} \lambda_{q,i,l} w_{i,k,q} h_{i,l} \hat{\mathbf{X}}_{k,l}^{-1} \\ & \quad + \sum_o d_{i,o} \left(-(\mathbf{J}_{k,o} + \epsilon \mathbf{I})^{-1} \right) + (\gamma - M) \tilde{U}_{i,k}^{-1}, \end{aligned} \quad (28)$$

$$B_{i,k,0} = C_{i,k,0} \left(\sum_{l,q} \lambda_{q,i,l} w_{i,k,q} h_{i,l} \hat{\mathbf{X}}_{k,l}^{-1} \mathbf{y}_{k,l} \mathbf{y}_{k,l}^H \hat{\mathbf{X}}_{k,l}^{-1} \right) C_{i,k,0}. \quad (29)$$

We can solve this equation with the same procedure proposed in [5] and obtain the update rule of $C_{i,k,0}$. For updating $C_{i,k,0}$, we must calculate the inverse of $\tilde{U}_{i,k}$. Even if $C_{i,k,0}$ equals $\mathbf{J}_{\theta_o,k} + \epsilon \mathbf{I}$, we can calculate the inverse of $\tilde{U}_{i,k}$ which is updated by Eq. (16). However, experimentally, $\tilde{U}_{i,k}$ would not always be a full-rank matrix in the iterations. Therefore, we update $\tilde{U}_{i,k}$ approximately as:

$$\tilde{U}_{i,k} \approx C_{i,k,0} + \epsilon' \mathbf{I}, \quad (30)$$

where ϵ' is a small value. On account of the space, we omit the update rules of the other parameters in Θ , however, they can be also obtained analytically by setting the partial differentials of L^+ at zeros.

4. EXPERIMENTAL EVALUATION

We evaluated the performance of the present method in terms of the abilities of source separation. We used 5 mixed stereo signals (therefore the number of the microphones M is 2) as the experimental data, each of which we obtained by mixing three speech signals (speech of a male and two females) from the ATR database [11] and was convolved with the measured room impulse response from the RWCP database [12] (in which the distance between the microphones was 11.48 cm and the reverberation time was 0 ms/380 ms). Thus, the

Table 2. SDR and SIR improvements obtained by the proposed method, the method in [7] and [8]

	RT=0 ms		RT=380 ms	
	SDR	SIR	SDR	SIR
method in [7]	-0.53	4.65	-1.55	3.73
method in [8]	1.90	7.88	-1.57	4.99
proposed	3.78	14.12	0.37	8.65

three signals were artificially located 30° , 90° and 130° from the microphones respectively. Fig. 1 shows an example of a spectrogram of the observed signal (the reverberation time was 380 ms). The sampling rate was 16 kHz. To compute the STFT components of the observed signal, the STFT frame length was set at 64 ms and a Hamming window was used with an overlap length of 48 ms. We set the hyperparameters as Table 1. We expected that $j = 1$ is an inactive state and $j = 2$ is an active state. We set \mathcal{T} as 0 for the anechoic mixed signals, as 3 for echoic ones. d was initially set as shown in Fig. 2 (a). $C_{i,k,0}$ was set to $1/\sqrt{M} \times \mathbf{I}$, $C_{i,k,1:\mathcal{T}}$ were set to $10^{-2}/\sqrt{M} \times \mathbf{I}$. \mathbf{W} was initially randomized. \mathbf{H} was set as 1 initially. $\rho_i^{(j)}$ was set as $\rho_{i,1}^{(j)} = (0.9, 0.1)$ and $\rho_{i,2}^{(j)} = (0.1, 0.9)$, all components of $\rho_i^{(q)}$ were set as $1/Q$. We set $\rho^{(o)}$ individually for each source so that directions of 3 sources did not overlap each other. The parameter estimation algorithm was run for 50 iterations. In order to avoid an undesirable local optima, we iterated the proposed algorithm 25 times with setting \mathcal{T} as 0, then gradually increased \mathcal{T} up to 3 according to the iteration. Moreover, in order to avoid huge numerical errors, we fixed \mathbf{W} when $\mathcal{T} > 0$. The numerical errors occurred probably because we fitted the product of $\mathbf{W}_{i,k}$ and $C_{i,k,1:\mathcal{T}}$ to observed signals, and so each of $\mathbf{W}_{i,k}$ and $C_{i,k,1:\mathcal{T}}$ was arbitrary scaled. We chose the methods proposed in [7] and [8] as comparisons. The method of [7] does not use the information of DOAs, and that of [8] models a spatial correlation matrix by using a weighted sum of \mathbf{J} . The separated signal was obtained by Wiener filtering. As evaluation measures, we used the signal-to-distortion ratio (SDR) and the signal-to-interference ratio (SIR) [13]. The SDR and SIR are expressed in decibels (dB), and a higher SDR/(SIR) indicates superior quality. The SDRs and SIRs were calculated with magnitude spectrograms. The averages of SDRs and SIRs of the mixed signals were -1.98 and -1.04[dB] when RT=0 ms, -4.77, -0.97[dB] when RT=380 ms respectively.

Table 2 shows the average SDR and SIR improvements obtained by our previous and present methods. The average SDRs and SIRs obtained with the present method were superior to those obtained with our previous methods. Fig. 2 (b) and (c) show examples of d obtained by the proposed method when RT=0 ms and RT=380 ms respectively. The arrows show the actual direction of the sources. We can see that the DOA of the sources were estimated roughly by the proposed method. Fig. 3 shows (a) a spectrogram of an anechoic source signal, (b) that of a separated and dereverberated signal with the proposed method and (c) the result of $\lambda^{(i)}$ which corresponds to the estimated activity of the source obtained with the mixed signal showed in Fig. 1. We expected $q = 1$ was an inactive state by setting the hyperparameters of the gamma distribution properly (see Table 1), and the result shows the voice activity was roughly detected.

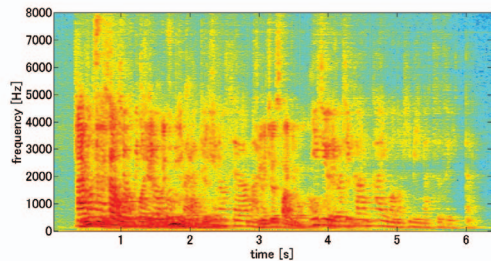


Fig. 1. An example of the spectrogram of the mixed signal (RT=380 ms).

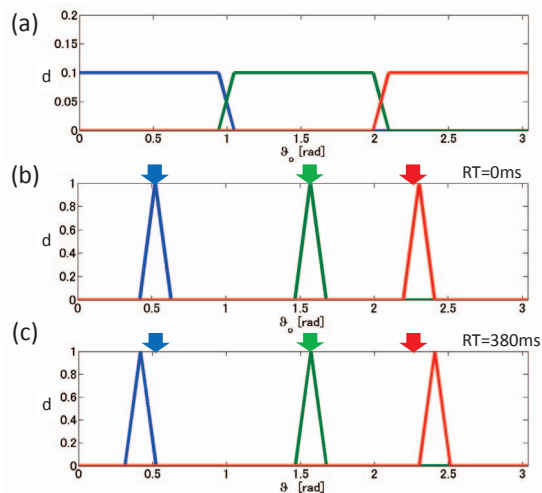


Fig. 2. (a) First condition of d , (b) an example of the estimation results of d obtained with the mixed signal (RT=0 ms) and (c) that obtained with the mixed signal (RT=380 ms). d of each source is colored differently. d corresponds to how likely the source is at the direction.

5. CONCLUSION

In this paper, we incorporate the DOA mixture model into the multichannel factorial hidden Markov model. Parameter estimation allows us to solve the problems of BSS, VAD, dereverberation and DOA estimation simultaneously. The average of SDRs obtained by the proposed method was 1.94 dB higher than that obtained by our previous method.

REFERENCES

- [1] A. Hyvärinen *et al.*, *Independent Component Analysis*, John Wiley & Sons, 2001.
- [2] D. D. Lee, and H. S. Seung, “Learning the parts of objects with nonnegative matrix factorization,” *Nature*, vol. 401, pp. 788–791, 1999.
- [3] P. Smaragdis, and J. C. Brown, “Non-negative matrix factorization for polyphonic music transcription,” *WASPAA 2003*, pp. 177–180, Oct. 2003.
- [4] A. Ozerov, and C. Févotte, “Multichannel nonnegative matrix factorization in convolutive mixtures for audio source separation,” *IEEE Trans. ASLP*, vol. 18, no. 3, pp. 550–563, 2010.
- [5] H. Sawada *et al.*, “Efficient algorithms for multichannel extensions of Itakura-Saito nonnegative matrix factorization,” *ICASSP 2012*, pp. 261–264, 2012.

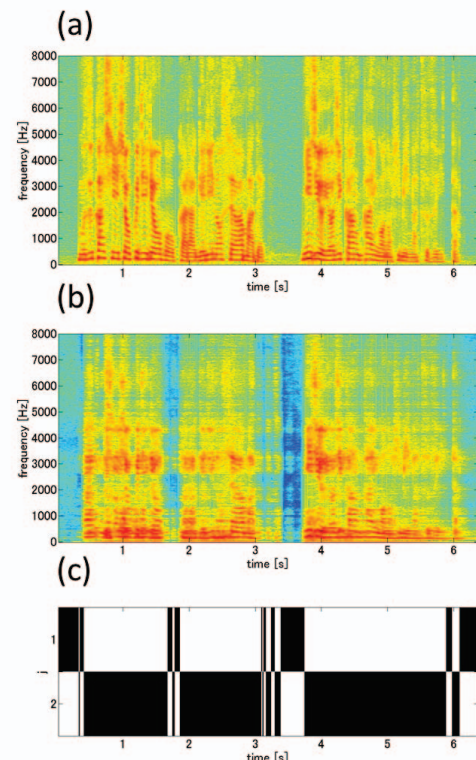


Fig. 3. (a) An example of a spectrogram of a source signal, (b) that of the separated and dereverberated signal obtained by the proposed method and (c) the estimation result of $\lambda^{(j)}$ which corresponds to source’s activity obtained with the mixed signal showed in Fig. 1. Black indicates higher value and the state is more likely to be assigned at that time in (c).

- [6] T. Higuchi, *et al.*, “A unified approach for underdetermined blind signal separation and source activity detection by multichannel factorial hidden Markov models,” *Interspeech 2014*, pp. 850–854, 2014.
- [7] T. Higuchi and H. Kameoka, “Joint audio source separation and dereverberation based on multichannel factorial hidden Markov model,” *MLSP 2014*.
- [8] T. Higuchi and H. Kameoka, “Unified approach for underdetermined BSS, VAD, dereverberation and DOA estimation with multichannel factorial HMM,” *GlobalSIP 2014*, pp. 725–729, 2014.
- [9] J. Nikunen and T. Virtanen, “Multichannel audio separation by direction of arrival based spatial covariance model and non-negative matrix factorization,” *ICASSP 2014*, pp. 6727–6731, 2014.
- [10] H. Kameoka *et al.*, “Blind separation of infinitely many sparse sources,” *IWAENC 2012*, pp. 1–4, 2012.
- [11] A. Kurematsu *et al.*, “ATR Japanese speech database as a tool of speech recognition and synthesis,” *Speech Communication*, pp. 357–363, 1990.
- [12] S. Nakamura *et al.*, “Acoustical sound database in real environments for sound scene understanding and hands-free speech recognition,” *LREC 2000*, pp. 965–968, 2000.
- [13] E. Vincent *et al.*, “Performance measurement in blind audio source separation,” *IEEE Trans. ASLP*, pp. 1462–1469, 2006.