

JOINT AUDIO SOURCE SEPARATION AND DEREVERBERATION BASED ON MULTICHANNEL FACTORIAL HIDDEN MARKOV 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

This paper proposes a unified approach for jointly solving underdetermined source separation, audio event detection and dereverberation of convolutive mixtures. 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. Several attempts have recently been made to extend this approach to a multichannel case in order to utilize the spatial correlation of the multichannel inputs as an additional clue for source separation. The multichannel NMF assumes that an observed signal is a mixture of a limited number of source signals each of which has a static power spectral density scaled by a time-varying amplitude. We have previously proposed an extension of this approach, in which the variations over time of the spectral density and the total power of each source is modeled by a hidden Markov model (HMM). This has allowed us to solve source activity detection and source separation simultaneously through model parameter inference. While this method was based on an anechoic mixing model, the aim of this paper is to further extend the above approach to deal with reverberation by incorporating an echoic mixing model into the generative model of observed signals. Through an experiment of underdetermined source separation under reverberant conditions, we confirmed that the proposed method provided a 9.61 dB improvement compared with the conventional method in terms of the signal-to-interference ratio.

Index Terms: source separation, dereverberation, audio event detection, non-negative matrix factorization, multichannel factorial hidden Markov model,

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. The best-known commercial application of BSS techniques is their use in teleconferencing systems. To solve BSS problems, 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. However, in an underdetermined case, the independence assumption is too weak to allow us to determine a unique solution and so directly applying ICA will not work well.

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]. With this approach, the spectrogram of a mixture signal is factorized into the product of a basis matrix consisting of basis spectra and an activation matrix consisting of time-varying amplitudes associated with the basis spectra. Several attempts have recently been made to extend this approach to a multichannel case in order to allow for the use of the spatial correlation of multichannel inputs as an additional clue for separation, which have opened a door to a new promising approach for underdetermined BSS [4][5]. This approach is based on the assumption that an observed signal is a mixture of a limited number of source signals each of which has a static power spectral density (i.e., the basis spectrum) scaled by a time-varying amplitude. However, many source signals in real world are non-stationary in nature and the variations of the spectral densities are much richer in time. Another important fact is that many sources including speech tend to stay inactive for some while until they switch to an active mode. This implies that the total power of a source may depend on its underlying state. To reasonably characterize such a non-stationary nature of source signals, we previously extended the multichannel NMF model by modeling the transition of the set consisting of the spectral densities and the total power of each source using a hidden Markov model (HMM). We call this model the “multichannel factorial hidden Markov model (MFHMM)” [6]. With this model, we are able to perform source separation and source activity detection simultaneously. However, in the models mentioned above, the length of the impulse response from a source to microphones is assumed to be sufficiently shorter than the frame length of the STFT so that the observed signal can be approximately modeled by an instantaneous mixture in the time-frequency domain. To deal with reverberation, this

paper aims to further extend the above approach by incorporating an echoic mixing model into MFHMM. Parameter inference of the present model allows us to simultaneously solve source separation, source activity detection and dereverberation based on a unified maximum likelihood criterion.

The remainder of this paper is organized as follows: Sec. 2 formulates a generative model of a multichannel observed signal under a reverberant condition, Sec. 3 presents the generative model of a source signal based on an HMM, Sec. 4 derives a parameter estimation algorithm for the present model and Sec. 5 presents some results of a source separation experiment.

2. ECHOIC MIXING MODEL

First we consider a situation where I source signals are recorded by M microphones. Here, let $y_m(t) \in \mathbb{R}$ be the observed signal at the m -th microphone, and $s_i(t) \in \mathbb{R}$ be the signal of the i -th source. The observed signal can be written in the time domain:

$$\mathbf{y}(t) = \sum_{i=1}^I \int_{-\infty}^{\infty} \mathbf{a}_i(\tau) s_i(t - \tau) d\tau, \quad (1)$$

where $\mathbf{y}(t) = (y_1(t), \dots, y_M(t))^T \in \mathbb{R}^M$ and $\mathbf{a}_i(t) = (a_{i,1}(t), \dots, a_{i,M}(t))^T \in \mathbb{R}^M$. $a_{i,m}(t)$ denotes the impulse response between source i and microphone m . If we assume that the length of the impulse response from a source to microphones is sufficiently shorter than the frame length of the STFT, the observed signal can be approximated fairly well by an instantaneous mixture in the time-frequency domain:

$$\mathbf{y}(\omega_k, t_l) = \sum_{i=1}^I \mathbf{a}_i(\omega_k) s_i(\omega_k, t_l), \quad (2)$$

where $\mathbf{y}(\omega_k, t_l) = (y_1(\omega_k, t_l), \dots, y_M(\omega_k, t_l))^T \in \mathbb{C}^M$ and $\mathbf{a}_i(\omega_k) = (a_{i,1}(\omega_k), \dots, a_{i,M}(\omega_k))^T \in \mathbb{C}^M$. 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. $1 \leq k \leq K$ and $1 \leq l \leq L$ are the frequency and time indices, respectively. $\mathbf{a}_i(\omega_k)$ denotes the frequency array response for source i at frequency ω_k . However, 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:

$$\mathbf{y}(\omega_k, t_l) \approx \sum_{i=1}^I \sum_{n=0}^N \mathbf{a}_i(\omega_k, t_n) s_i(\omega_k, t_l - t_n). \quad (3)$$

$0 \leq n \leq N$ is the time index of the frequency array response in the time-frequency domain. Note that $\mathbf{a}_i(\omega_k, t_1 : t_N)$ denote the frequency array responses which correspond to the impulse responses out of the frames of the STFT. This approximation is useful for dereverberation [7] and the validity of the approximation is experimentally shown. For convenience of notation, we hereafter use subscripts k, l and n to indicate ω_k, t_l and t_n respectively.

3. MULTICHANNEL FACTORIAL HMM

3.1. Generative process of observed signals

Here we describe the generative process of an observed signal based on Eq. (3). If we assume that each source signal follows a piecewise stationary Gaussian process, then $s_{i,k,l}$ follows a complex normal distribution with mean 0 and covariance $\sigma_{i,k,l}^2$,

$$s_{i,k,l} | \sigma_{i,k,l} \sim \mathcal{N}_{\mathbb{C}}(s_{i,k,l}; 0, \sigma_{i,k,l}^2), \quad (4)$$

where $\sigma_{i,k,l}^2$ denotes the power spectral density of i -th source at frequency k and time l and $\mathcal{N}_{\mathbb{C}}(\mathbf{x}; \boldsymbol{\mu}, \Sigma) \propto \exp(-(\mathbf{x} - \boldsymbol{\mu})^H \Sigma^{-1} (\mathbf{x} - \boldsymbol{\mu}))$. From Eq. (3) and Eq. (4), $\mathbf{y}_{k,l}$ is also normally distributed such that

$$\mathbf{y}_{k,l} | \mathbf{a}_{1:I,k,0:N}, \sigma_{i,k,l-N:l} \sim \mathcal{N}_{\mathbb{C}}(\mathbf{y}_{k,l}; 0, \sum_{i,n} \mathbf{C}_{i,k,n} \sigma_{i,k,l-n}^2), \quad (5)$$

conditioned on $\mathbf{a}_{1:I,k,0:N}$ and $\sigma_{i,k,l-N:l}$ where $\mathbf{C}_{i,k,n} = \mathbf{a}_{i,k,n} \mathbf{a}_{i,k,n}^H$. $\mathbf{C}_{i,k,n}$ represents n -th spatial correlation matrix of the i -th source at frequency k .

3.2. Generative model for multichannel NMF [4, 5]

In the regular NMF model (see [8]), the power spectra of a source signal is assumed to be static up to a scale factor. We can incorporate this assumption into the above model by setting

$$\sigma_{i,k,l}^2 = w_{i,k} h_{i,l}, \quad (6)$$

where $\sigma_{i,k,l}^2$ is assumed to be factorized into the product of the static power spectrum $w_{i,k}$ and the time-varying amplitude $h_{i,l}$. The generative model of $s_{i,k,l}$ is thus rewritten as

$$s_{i,k,l} | w_{i,k}, h_{i,l} \sim \mathcal{N}_{\mathbb{C}}(s_{i,k,l}; 0, w_{i,k} h_{i,l}), \quad (7)$$

conditioned on $w_{i,k}$ and $h_{i,l}$.

3.3. Generative modeling of source signals using HMMs

As described above, the multichannel NMF model roughly assumes that the power spectra of each sound source is static up to a scale factor. However, many sound sources exhibit different spectra according to underlying ‘‘states’’ of the sources. For example, the spectra of the sound of a piano note would be different in ‘‘attack,’’ ‘‘decay,’’ ‘‘sustain’’ and ‘‘release’’ segments. Another important fact is that the total power of a source also depends on its underlying state. To reasonably characterize such a non-stationary nature of source signals, here we model the sequence of the power spectra and the total powers of each source using an HMM.

Now we introduce a latent variable $z_{i,l} \in \{1, \dots, D\}$ to denote a state of the i -th source at time l . The state sequence $z_{i,1}, \dots, z_{i,L}$ is assumed to follow a Markov chain:

$$z_{i,l} | z_{i,l-1} \sim \text{Categorical}(z_{i,l}; \boldsymbol{\rho}_{z_{i,l-1}}), \quad (8)$$

where $\text{Categorical}(x; \mathbf{y}) = y_x$, $\boldsymbol{\rho}_d = (\rho_{d,1}, \dots, \rho_{d,D})$ denotes the transition probability of state d to each state $1, \dots, D$, and $\boldsymbol{\rho} = (\rho_{d,d'})_{D \times D}$ denotes the transition matrix. Here, we assume that $h_{i,l}$ follows a gamma distribution with hyperparameters determined according to $z_{i,l}$,

$$h_{i,l}|z_{i,l} \sim \text{Gamma}(h_{i,l}; \alpha_{z_{i,l}} \beta_{z_{i,l}}), \quad (9)$$

where $\alpha_{1:D}$ and $\beta_{1:D}$ 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 $z_{i,l}$ 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 become uniform distributions. We expect that this setting allows us to solve source separation and source activity detection in a cooperative manner. Let us use $w_{i,k,d}$ to denote the power spectrum of the i -th source at state d . The power spectrum of the i -th source at time l is also assumed to be determined according to $z_{i,l}$. Thus, the generative model of $s_{i,k,l}$ is eventually written as

$$s_{i,k,l}|w_{i,k,1:D}, h_{i,l}, z_{i,l} \sim \mathcal{N}_{\mathbb{C}}(s_{i,k,l}; 0, w_{i,k,z_{i,l}} h_{i,l}). \quad (10)$$

Since the generative model of $\mathbf{y}_{k,l}$ contains multiple HMMs associated with the underlying sources, the overall model can be viewed as a Factorial HMM. Our overall generative model is given by Eqs. (8), Eqs. (9) and

$$\mathbf{y}_{k,l} | \mathbf{a}_{1:I,k,0:N}, w_{1:I,k,1:D}, h_{1:I,l-N:l}, z_{1:I,l-N:l} \sim \mathcal{N}_{\mathbb{C}}(\mathbf{y}_{k,l}; 0, \sum_{i,n} \mathbf{C}_{i,k,n} w_{i,k,z_{i,l-n}} h_{i,l-n}), \quad (11)$$

conditioned on $\mathbf{a}_{1:I,k,0:N}$, $w_{1:I,k,1:D}$, $h_{1:I,l-N:l}$ and $z_{1:I,l-N:l}$.

3.4. Related work

The present model is equivalent to the model proposed in [6] when the number N is set at 0, which leads to an anechoic mixing model.

4. ALGORITHM FOR PARAMETER ESTIMATION

4.1. Objective function

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 $\mathbf{W} = w_{1:I,1:K,1:D}$, $\mathbf{H} = h_{1:I,1:L}$, $\mathbf{C} = \mathbf{C}_{1:I,1:K,0:N}$ and $\mathbf{Z} = z_{1:I,1:L}$. We denote the entire set of the above parameters as Θ . In the following, $\boldsymbol{\rho}$ is constants that is determined experimentally. Our goal is to compute the posterior

$$p(\Theta|\mathbf{Y}) = \frac{p(\mathbf{Y}, \Theta)}{p(\mathbf{Y})}, \quad (12)$$

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. 3, we can write the joint distribution $p(\mathbf{Y}, \Theta)$ as

$$p(\mathbf{Y}, \Theta) \propto p(\mathbf{Y}|\Theta)p(\mathbf{H}|\mathbf{Z})p(\mathbf{Z}). \quad (13)$$

The objective function is defined as $L(\Theta) = \log p(\Theta|\mathbf{Y})$. Our goal is to obtain $\hat{\Theta}$ such that

$$\hat{\Theta} = \underset{\Theta}{\text{argmax}} \log p(\Theta|\mathbf{Y}). \quad (14)$$

By using Eqs. (12), (13) and (14), the current optimization problem can be rewritten as

$$\hat{\Theta} = \underset{\Theta}{\text{argmax}} (\log p(\mathbf{Y}|\Theta) + \log p(\mathbf{H}|\mathbf{Z}) + \log p(\mathbf{Z})). \quad (15)$$

According to the generative model defined in Sec. 3, $\log p(\mathbf{Y}|\Theta)$ is written as

$$\begin{aligned} & \log p(\mathbf{Y}|\Theta) \\ &= -\frac{1}{2} \sum_{k,l} (M \log 2\pi + \log |\hat{\mathbf{X}}_{k,l}| + \mathbf{y}_{k,l}^H \hat{\mathbf{X}}_{k,l}^{-1} \mathbf{y}_{k,l}), \end{aligned} \quad (16)$$

where $\hat{\mathbf{X}}_{k,l} = \sum_{i,n} \mathbf{C}_{i,k,n} w_{i,k,z_{i,l-n}} h_{i,l-n}$.

4.2. Optimization 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 $\Theta \leftarrow \underset{\Theta}{\text{argmax}} L^+(\Theta, \Lambda)$ and $\Lambda \leftarrow \underset{\Lambda}{\text{argmax}} L^+(\Theta, \Lambda)$. 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.

$$\begin{aligned} & L(\Theta) \\ & \geq L^+(\Theta, \Lambda) \\ &= -\frac{1}{2} \sum_{k,l} \left\{ \sum_{i,n} \left(\frac{\text{tr}(\mathbf{y}_{k,l} \mathbf{y}_{k,l}^H \mathbf{R}_{i,k,l,n} \mathbf{C}_{i,k,n}^{-1} \mathbf{R}_{i,k,l,n})}{w_{i,k,z_{i,l-n}} h_{i,l-n}} \right) \right. \\ & \quad \left. + \text{tr}(\mathbf{U}_{k,l}^{-1} \mathbf{C}_{i,k,n}) w_{i,k,z_{i,l-n}} h_{i,l-n} \right) + \log |\mathbf{U}_{k,l}| - M \Big\} \\ & \quad + \sum_{i,l} \{ (\alpha_{z_{i,l}} - 1) \log h_{i,l} - h_{i,l} / \beta_{z_{i,l}} - \alpha_{z_{i,l}} \log \beta_{z_{i,l}} \} \\ & \quad + \log p(\mathbf{Z}), \end{aligned} \quad (17)$$

where $\mathbf{R}_{i,k,l,n}$ and $\mathbf{U}_{k,l}$ are auxiliary variables that satisfy Hermitian positive definiteness and $\sum_{i,n} \mathbf{R}_{i,k,l,n} = \mathbf{I}$. We

denote the set of the auxiliary variables as Λ . $\text{tr}(\cdot)$ is the trace of a matrix. The equality $L(\Theta) = L^+(\Theta, \Lambda)$ is satisfied when

$$\mathbf{R}_{i,k,l,n} = \mathbf{C}_{i,k,n} w_{i,k,z_{i,l-n}} h_{i,l-n} \hat{\mathbf{X}}_{k,l}^{-1}, \quad (18)$$

$$\mathbf{U}_{k,l} = \hat{\mathbf{X}}_{k,l}. \quad (19)$$

Therefore, we can monotonically increase L by repeating the following two steps.

1. Maximizing L^+ with respect to \mathbf{R} and \mathbf{U} .
2. Maximizing L^+ with respect to \mathbf{W} , \mathbf{H} , \mathbf{C} and \mathbf{Z} .

Step 1 consists in updating \mathbf{R} and \mathbf{U} using Eqs. (18) and (19). In step 2, we can obtain update rules of \mathbf{W} , \mathbf{H} , \mathbf{C} by setting the partial derivative of L^+ with respect to each of the parameters at zero. The partial derivatives of L^+ with respect to \mathbf{W} and \mathbf{H} are given by

$$\begin{aligned} & \frac{\partial L^+}{\partial w_{i,k,z_{i,l}}} \\ &= \frac{1}{2} \sum_{l,n} \left(\frac{\text{tr}(\mathbf{y}_{k,l+n} \mathbf{y}_{k,l+n}^H \mathbf{R}_{i,k,l+n,n} \mathbf{C}_{i,k,n}^{-1} \mathbf{R}_{i,k,l+n,n})}{w_{i,k,z_{i,l}}^2 h_{i,l}} \right. \\ & \left. - \text{tr}(\mathbf{U}_{k,l+n}^{-1} \mathbf{C}_{i,k,n}) h_{i,l} \right), \end{aligned} \quad (20)$$

$$\begin{aligned} & \frac{\partial L^+}{\partial h_{i,l}} \\ &= \frac{1}{2} \sum_{k,n} \left(\frac{\text{tr}(\mathbf{y}_{k,l+n} \mathbf{y}_{k,l+n}^H \mathbf{R}_{i,k,l+n,n} \mathbf{C}_{i,k,n}^{-1} \mathbf{R}_{i,k,l+n,n})}{w_{i,k,z_{i,l}} h_{i,l}^2} \right. \\ & \left. - \text{tr}(\mathbf{U}_{k,l+n}^{-1} \mathbf{C}_{i,k,n}) w_{i,k,z_{i,l}} \right) \\ & + (\alpha_{z_{i,l}} - 1)/h_{i,l} - 1/\beta_{z_{i,l}}, \end{aligned} \quad (21)$$

respectively. By setting them at zero, we obtain the following update rules:

$$w_{i,k,z_{i,l}} \leftarrow \sqrt{\frac{\sum_{l,n} \frac{\text{tr}(\mathbf{y}_{k,l+n} \mathbf{y}_{k,l+n}^H \mathbf{R}_{i,k,l+n,n} \mathbf{C}_{i,k,n}^{-1} \mathbf{R}_{i,k,l+n,n})}{h_{i,l}}}{\sum_{l,n} \text{tr}(\mathbf{U}_{k,l+n}^{-1} \mathbf{C}_{i,k,n}) h_{i,l}}}, \quad (22)$$

$$h_{i,l} \leftarrow \frac{(\alpha_{z_{i,l}} - 1) + \sqrt{(\alpha_{z_{i,l}} - 1)^2 + \mu_{i,l} \nu_{i,l}}}{\nu_{i,l}}, \quad (23)$$

where

$$\mu_{i,l} = \sum_{k,n} \frac{\text{tr}(\mathbf{y}_{k,l+n} \mathbf{y}_{k,l+n}^H \mathbf{R}_{i,k,l+n,n} \mathbf{C}_{i,k,n}^{-1} \mathbf{R}_{i,k,l+n,n})}{w_{i,k,z_{i,l}}}, \quad (24)$$

$$\nu_{i,l}$$

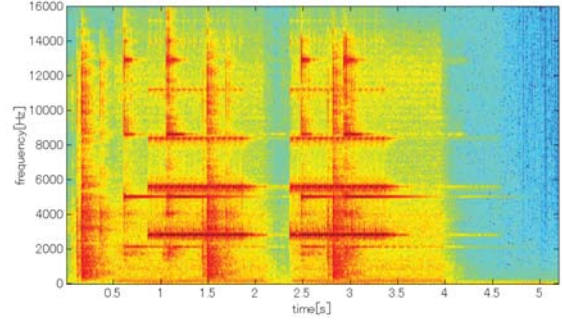


Fig. 1. Spectrogram of a mixture signal.

$$= \sum_{k,n} \text{tr}(\mathbf{U}_{k,l+n}^{-1} \mathbf{C}_{i,k,n}) w_{i,k,z_{i,l}} + 2/\beta_{z_{i,l}}. \quad (25)$$

The partial derivatives of L^+ with respect to \mathbf{C} is given by

$$\begin{aligned} \frac{\partial L^+}{\partial \mathbf{C}_{i,k,n}} &= \sum_l \left(\frac{\mathbf{C}_{i,k,n}^{-1} \mathbf{R}_{i,k,l,n} \mathbf{y}_{k,l} \mathbf{y}_{k,l}^H \mathbf{R}_{i,k,l,n} \mathbf{C}_{i,k,n}^{-1}}{w_{i,k,z_{i,l-n}} h_{i,l-n}} \right. \\ & \left. - \mathbf{U}_{k,l}^{-1} w_{i,k,z_{i,l-n}} h_{i,l-n} \right). \end{aligned} \quad (26)$$

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

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

where

$$\mathbf{A}_{i,k,n} = \sum_l w_{i,k,z_{i,l-n}} h_{i,l-n} \hat{\mathbf{X}}_{k,l}^{-1},$$

$$\mathbf{B}_{i,k,n} = \mathbf{C}_{i,k,n} \left(\sum_l w_{i,k,z_{i,l-n}} h_{i,l-n} \right.$$

$$\left. \hat{\mathbf{X}}_{k,l}^{-1} \mathbf{y}_{k,l} \mathbf{y}_{k,l}^H \hat{\mathbf{X}}_{k,l}^{-1} \right) \mathbf{C}_{i,k,n}. \quad (28)$$

We can solve this equation by using a method in [5]. We perform an eigenvalue decomposition of a $2M \times 2M$ matrix

$$\begin{bmatrix} 0 & -\mathbf{A}_{i,k,n} \\ -\mathbf{B}_{i,k,n} & 0 \end{bmatrix}, \quad (29)$$

and let $\mathbf{e}_{1,i,k,n} \dots \mathbf{e}_{M,i,k,n}$ be eigenvectors with negative eigenvalues. It is guaranteed that there are exactly M negative eigenvalues. Then, let us decompose the $2M$ -dimensional eigenvectors as

$$\mathbf{e}_{m,i,k,n} = \begin{bmatrix} \mathbf{f}_{m,i,k,n} \\ \mathbf{g}_{m,i,k,n} \end{bmatrix}, \quad (30)$$

for $m = 1 \dots M$ where $\mathbf{f}_{m,i,k,n}$ and $\mathbf{g}_{m,i,k,n}$ are M -dimensional vectors. We obtain the update rule for $\mathbf{C}_{i,k,n}$ as

$$\mathbf{C}_{i,k,n} \leftarrow \mathbf{G}_{i,k,n} \mathbf{F}_{i,k,n}^{-1}, \quad (31)$$

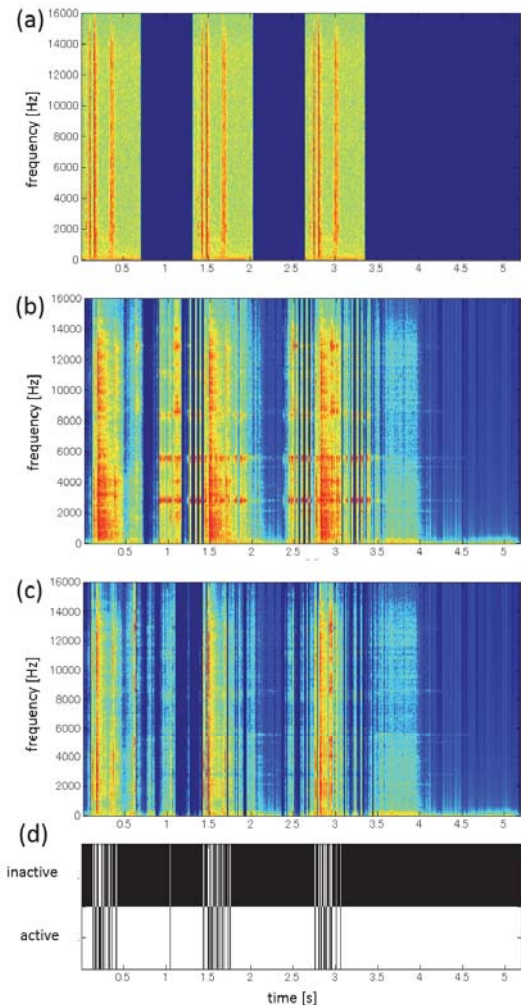


Fig. 2. (a) A spectrogram of the source signal of the stapler recorded in an anechoic condition, (b) that of the separated signal of the stapler obtained by the conventional method, (c) that of the separated and dereverberated signal of the stapler obtained by the proposed method, and (d) the acoustic event detection result obtained by the proposed method. Black indicates the state is assigned at the time.

where $F_{i,k,n} = [f_{1,i,k,n}, \dots, f_{M,i,k,n}]$ and $G_{i,k,n} = [g_{1,i,k,n}, \dots, g_{M,i,k,n}]$.

L^+ is equal up to constant terms to the sum of the log posteriors of HMMs, when viewed as a function of Z . Thus, we can invoke the Viterbi algorithm to search for the optimal path $z_{i,1}, \dots, z_{i,L}$ for each i individually. Note that updating W , H , C and Z corresponds to solving the problems of source separation, source activity detection and dereverberation based on a unified objective function.

5. EXPERIMENTAL EVALUATION

We evaluated the performance of the present method in terms of the abilities of source separation, source activity detec-

Table 1. The SIRs of the three sources obtained by the conventional and proposed methods.

SIR [dB]	bell	phone	stapler
Proposed	43.35	32.05	7.19
Conventional	32.37	30.15	-8.77

tion and dereverberation in a supervised case. We used a mixed stereo signal as the experimental data, each of which we obtained by mixing the non-speech signals (sounds of a cell phone, a bell and a stapler) from the RWCP database [9] and was convolved with the measured room impulse response from the RWCP database [9] (in which the distance between the microphones was 2.83 cm and the reverberation time was 600 ms). Thus, the three signals were artificially located 60° , 90° and 130° from the microphones respectively. Fig. 1 shows a spectrogram of the observed signal in the reverberant condition. The sampling rate was 32 kHz. To compute the STFT components of the observed signal, the STFT frame length was set at 16 ms and a Hamming window was used with an overlap length of 8 ms. We set the number of states of HMMs D as 2. We expected that $d = 1$ was an inactive state and $d = 2$ was an active state, by setting α_1 and β_1 as 1 and 10^{-2} respectively, and α_2 and β_2 as 1 and 10^{20} respectively. The diagonal elements of $C_{i,k,0}$ were initially all set to $1/\sqrt{M}$, and the off-diagonal elements were initially set to zero. For $n = 1, \dots, N$, the diagonal elements of $C_{i,k,n}$ were set to $10^{-3}/\sqrt{M}$, and the off-diagonal elements were also set to zero initially. We first learned W , H and ρ from the three clean source signals recorded in the anechoic condition with the proposed method starting from random initial matrices W and H (therefore this experiment was a supervised source separation). The parameter estimation algorithm was run for 150 iterations. In order to avoid an undesirable local optima, we iterated the proposed algorithm 100 times with setting N as 0, then gradually increased N up to 20 according to the iteration. We chose the method proposed in [6] as a comparison. This method equals to be our proposed method when we set $N = 0$. The separated signal $\hat{y}_{i,k,l}$ was obtained by Wiener filtering

$$\hat{y}_{i,k,l} = w_{i,k,z_{i,l}} h_{i,l} C_{i,k,0} \hat{X}_{k,l}^{-1} y_{k,l}. \quad (32)$$

As evaluation measures, we used the signal-to-interference ratio (SIR) [10]. The SIR is expressed in decibels (dB), and a higher SIR indicates superior quality. The input SIR of the sounds of the bell, the cell phone and the stapler were -16.61, 16.58 and -39.17 [dB], respectively.

Table 1 shows the SIRs obtained by the conventional and proposed methods. The average of the SIRs obtained with the proposed method was 9.61 [dB] more than that obtained with the conventional approach. These results show the effectiveness of the proposed method for source separation in a reverberant condition. Fig. 2 shows (a) a spectrogram of the source signal of the stapler recorded in an anechoic condition, (b) that of the separated signal of the stapler obtained by the conventional method, (c) that of the separated and dereverberated signal of the stapler obtained by the proposed method,

and (d) the acoustic event detection result obtained by the proposed method respectively. Black indicates the state is assigned at the time. We can see that the reverberant components were relatively removed by the proposed method and the acoustic event was roughly detected.

6. CONCLUSION

This paper proposed a unified approach for jointly solving underdetermined source separation, audio event detection and dereverberation of convolutive mixtures based on MFHMM with an echoic mixing model. Through an experiment of supervised source separation under reverberant conditions, we confirmed that the proposed method provided a 9.61 dB improvement compared with the conventional method in terms of the signal-to-interference ratio.

7. REFERENCES

- [1] A. Hyvärinen, J. Karhunen and E. Oja, *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. Smaragdīs and J. C. Brown, “Non-negative matrix factorization for polyphonic music transcription,” in *Proc. 2003 IEEE Workshop on Applications of Signal Processing to Audio and Acoustics (WASPAA 2003)*, pp. 177–180, 2003.
- [4] A. Ozerov and C. Févotte, “Multichannel nonnegative matrix factorization in convolutive mixtures for audio source separation,” *IEEE Trans. Audio, Speech and Language Processing*, vol. 18, no. 3, pp. 550–563, Mar.2010.
- [5] H. Sawada, H. Kameoka, S. Araki and N. Ueda, “Efficient algorithms for multichannel extensions of Itakura-Saito nonnegative matrix factorization,” in *Proc. 2012 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP 2012)*, pp. 261–264, 2012.
- [6] T. Higuchi, H. Takeda, T. Nakamura and H. Kameoka, “A unified approach for underdetermined blind signal separation and source activity detection by multichannel factorial hidden Markov models,” in *Proc. 15th Annual Conference of the International Speech Communication Association (Interspeech 2014)*, to appear.
- [7] T. Nakatani, T. Yoshioka, K. Kinoshita, M. Miyoshi and B.-H. Juang, “Blind speech dereverberation with multichannel linear prediction based on short time Fourier transform representation,” in *Proc. 2008 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP 2008)*, pp. 85–88, 2008.
- [8] C. Févotte, N. Bertin and J. L. Durrieu, “Nonnegative matrix factorization with the Itakura-Saito divergence: With application to music analysis,” *Neural computation*, vol. 21, no. 3, pp. 793–830, 2009.
- [9] 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. 2nd International Conference on Language Resources & Evaluation (LREC 2000)*, pp. 965–968, 2000.
- [10] E. Vincent, R. Gribonval and C. Févotte, “Performance measurement in blind audio source separation,” *IEEE Transactions on Audio, Speech, and Language Processing*, pp. 1462–1469, 2006.