

A unified approach for underdetermined blind signal separation and source activity detection by multichannel factorial hidden Markov models

Takuya Higuchi¹, Hirofumi Takeda¹, Tomohiko Nakamura¹ and Hirokazu Kameoka^{1,2}

¹Graduate School of Information Science and Technology, The University of Tokyo

²NTT Communication Science Laboratories, NTT Corporation

{higuchi, h-takeda, nakamura, kameoka}@hil.t.u-tokyo.ac.jp

Abstract

This paper proposes to introduce a new model called “the multichannel factorial hidden Markov Model (MFHMM)” for underdetermined blind signal separation (BSS). 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. Up to now, multichannel extensions of NMF, which allow for the use of spatial information as an additional clue for source separation, have been proposed by several authors and proven to be an effective approach for underdetermined BSS. 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 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. Moreover, many sources including speech tend to stay inactive for some while until they switch to an active mode, implying that the total power of a source may depend on its underlying state. To reasonably characterize such a non-stationary nature of source signals, this paper proposes to extend 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). By letting each HMM contain states corresponding to active and inactive modes, we will show that voice activity detection and source separation can be solved simultaneously through parameter inference of the present model. The experiment showed that the proposed algorithm provided a 7.65 dB improvement compared with the conventional multichannel NMF in terms of the signal-to-distortion ratio.

Index Terms: blind signal separation, source activity detection, a hidden Markov model, non-negative matrix factorization

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. An important feature of this approach is that it is capable of finding a finite set of basis spectra that are considered to be the dominant elements constituting the observed spectrogram, in an unsupervised manner. Up to now, several attempts have been made to extend this approach to a multichannel case in order to allow for the use of spatial information 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. For example, the sound of a piano note would be more accurately characterized by a succession of several basis spectra corresponding to “attack,” “decay,” “sustain” and “release” segments than only by a single basis spectrum. 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, this paper proposes to extend 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). With this model, we would be able to flexibly reflect the time-varying nature of each source by appropriately specifying or training the transition probabilities prior to analysis. We formulate an entire generative model of a multichannel mixture signal by incorporating a generative model of the spectral densities of each source, described using an HMM. Each hidden state of the HMM represents the state of the corresponding source. We call this model “the multichannel factorial hidden Markov model (MFHMM)”.

In general, simultaneous estimation is preferable when several estimation problems are interdependent. If we knew when each source is active and inactive, source separation would become a relatively simple matter. On the other hand, if all the sources were already separated, voice activity detection would become a relatively simple matter. This simply implies that the problems of source separation and voice/source activity detection are interdependent of each other. It is important to note that through parameter inference of the present model, we would be able to simultaneously perform voice activity detection and source separation 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 and source signals based on the factorial HMM, Sec. 3 describes a parameter inference algorithm for the present model and Sec. 4 presents some experimental results.

2. Multichannel factorial HMM

2.1. Mixing model

First we consider a situation where I source signals are recorded by M microphones. 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. $1 \leq k \leq K$ and $1 \leq l \leq L$ are the frequency and time indices, respectively. 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 frequency domain:

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

where $\mathbf{y}(\omega_k, t_l) = (y_1(\omega_k, t_l), \dots, y_M(\omega_k, t_l))^T \in \mathbb{C}^M$. $\mathbf{a}_i(\omega_k)$ denotes the frequency array response for source i at frequency ω_k . For convenience of notation, we hereafter use subscripts k and l to indicate ω_k and t_l respectively.

2.2. Generative process of observed signals

Here we describe the generative process of an observed signal based on Eq. (1). 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 $\sigma_{i,k,l}^2$,

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

where $\sigma_{i,k,l}^2$ denotes the power spectral density of i -th source at frequency k and time l . From Eq. (1) and Eq. (2), $\mathbf{y}_{k,l}$ is also normally distributed such that

$$\mathbf{y}_{k,l} | \mathbf{a}_{1:I,k}, \sigma_{1:I,k,l} \sim \mathcal{N}_{\mathbb{C}}(\mathbf{y}_{k,l}; 0, \sum_i \mathbf{C}_{i,k} \sigma_{i,k,l}^2), \quad (3)$$

conditioned on $\mathbf{a}_{1:I,k} = \{\mathbf{a}_{1,k}, \dots, \mathbf{a}_{I,k}\}$, $\sigma_{1:I,k,l}^2 = \{\sigma_{1,k,l}^2, \dots, \sigma_{I,k,l}^2\}$ where $\mathbf{C}_{i,k} = \mathbf{a}_{i,k} \mathbf{a}_{i,k}^H$ and $\mathcal{N}_{\mathbb{C}}(\mathbf{x}; \boldsymbol{\mu}, \boldsymbol{\Sigma}) \propto \exp(-(\mathbf{x} - \boldsymbol{\mu})^H \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}))$. $\mathbf{C}_{i,k}$ represents a spatial correlation matrix of the i -th source at frequency k .

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

In the regular NMF model (see [6]), 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

$$\sigma_{i,k,l}^2 = w_{i,k} h_{i,l}, \quad (4)$$

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

$$s_{i,k,l} | w_{i,k}, h_{i,l} \sim \mathcal{N}_{\mathbb{C}}(s_{i,k,l}; 0, w_{i,k} h_{i,l}), \quad (5)$$

conditioned on $w_{i,k}$ and $h_{i,l}$. Since the generative model of $\mathbf{Y} = \{\mathbf{y}_{k,l}\}_{k,l}$ under the assumption Eq. (5) can be viewed as a natural extension of NMF to a multichannel case, a BSS approach based on this model is called the multichannel NMF [4, 5].

2.4. 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 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 also depends on its underlying state. To reasonably characterize such a non-stationary nature of source signals, here we propose to model the sequence of the power spectra and the total powers of each source using an HMM.

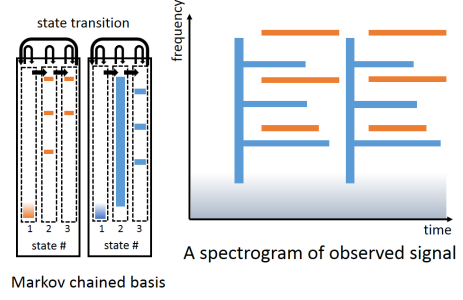


Figure 1: Illustration of the concept of the FHMM.

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:

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

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 $z_{i,l}$,

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

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. In particular, we consider setting α and β at 1 and 10^{-2} respectively for the inactive state, and at 1 and 10^{20} respectively for the active states. We expect that this setting allows us to solve source separation and voice 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

$$s_{i,k,l} | w_{i,k,1:D}, h_{i,l}, z_{i,l} \sim \mathcal{N}_{\mathbb{C}}(s_{i,k,l}; 0, w_{i,k,z_{i,l}} h_{i,l}). \quad (8)$$

Since the generative of $\mathbf{y}_{k,l}$ contains multiple HMMs associated with the underlying sources, the overall model can be viewed as a Factorial HMM. Fig. 1 shows an illustration of the proposed FHMM. Note that we can integrate an assumption that the signal tend to stay non-active for some time by setting the transition probabilities of self-transitions at reasonably large values in a specific state of the HMM. If we assume the emission probabilities of the HMMs as uniform distributions, our overall generative model is given by Eqs. (6) and

$$\mathbf{y}_{k,l} | \mathbf{a}_{1:I,k}, w_{1:I,k,1:D}, h_{1:I,l}, z_{1:I,l} \sim \mathcal{N}_{\mathbb{C}}(\mathbf{y}_{k,l}; 0, \sum_i \mathbf{C}_{i,k} w_{i,k,z_{i,l}} h_{i,l}), \quad (9)$$

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

3. Algorithm for parameter estimation

3.1. 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 $\mathbf{W} = w_{1:I,1:K,1:D}$, $\mathbf{H} = h_{1:I,1:L}$, $\mathbf{C} = \mathbf{C}_{1:I,1:K}$ and $\mathbf{Z} =$

$z_{1:L,1:L}$. We denote the entire set of the above parameters as Θ . In the following, ρ is constants that is determined experimentally. Our goal is to compute the posterior

$$p(\Theta|\mathbf{Y}) = \frac{p(\mathbf{Y}, \Theta)}{p(\mathbf{Y})}, \quad (10)$$

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 the joint distribution $p(\mathbf{Y}, \Theta)$ as

$$p(\mathbf{Y}, \Theta) \propto p(\mathbf{Y}|\Theta)p(\mathbf{H}|\mathbf{Z})p(\mathbf{Z}). \quad (11)$$

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

$$\hat{\Theta} = \underset{\Theta}{\operatorname{argmax}} \log p(\Theta|\mathbf{Y}). \quad (12)$$

By using Eqs. (10), (11) and (12), the current optimization problem can be rewritten as

$$\hat{\Theta} = \underset{\Theta}{\operatorname{argmax}} (\log p(\mathbf{Y}|\Theta) + \log p(\mathbf{H}|\mathbf{Z}) + \log p(\mathbf{Z})). \quad (13)$$

According to the generative model defined in Sec. 2, $\log p(\mathbf{Y}|\Theta)$ is written as

$$\begin{aligned} \log p(\mathbf{Y}|\Theta) &= \sum_{k,l} \left(-\frac{M}{2} \log 2\pi - \frac{1}{2} \log |\hat{\mathbf{X}}_{k,l}| - \frac{1}{2} \mathbf{y}_{k,l}^H \hat{\mathbf{X}}_{k,l}^{-1} \mathbf{y}_{k,l} \right), \end{aligned} \quad (14)$$

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

3.2. 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-increasing under the updates $\Theta \leftarrow \underset{\Theta}{\operatorname{argmax}} L^+(\Theta, \Lambda)$ and $\Lambda \leftarrow \underset{\Lambda}{\operatorname{argmax}} L^+(\Theta, \Lambda)$. 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.

$$\begin{aligned} &L(\Theta) \\ &\geq L^+(\Theta, \Lambda) \\ &= -\frac{1}{2} \sum_{k,l} \left(\sum_i \left(\frac{\operatorname{tr}(\mathbf{y}_{k,l}^H \mathbf{y}_{k,l} \mathbf{R}_{i,k,l} \mathbf{C}_{i,k}^{-1} \mathbf{R}_{i,k,l})}{w_{i,k,z_{i,l}} h_{i,l}} \right) \right. \\ &\quad \left. + \operatorname{tr}(\mathbf{U}_{k,l}^{-1} \mathbf{C}_{i,k}) w_{i,k,z_{i,l}} h_{i,l} \right) + \log |\mathbf{U}_{k,l}| - M \\ &\quad + \sum_{i,l} \left((\alpha_{z_{i,l}} - 1) \log h_{i,l} - h_{i,l} / \beta_{z_{i,l}} - \alpha_{z_{i,l}} \log \beta_{z_{i,l}} \right) \\ &\quad + \log p(\mathbf{Z}), \end{aligned} \quad (15)$$

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

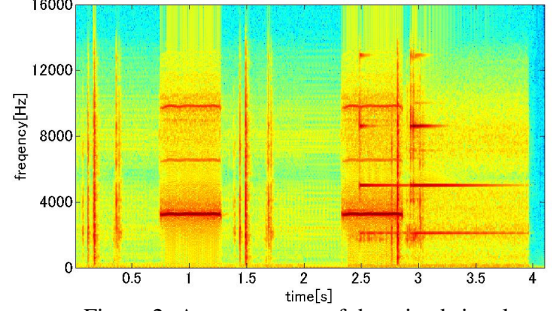


Figure 2: A spectrogram of the mixed signal.

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

$$\mathbf{R}_{i,k,l} = \mathbf{C}_{i,k} w_{i,k,z_{i,l}} h_{i,l} \hat{\mathbf{X}}_{k,l}^{-1}, \quad (16)$$

$$\mathbf{U}_{k,l} = \hat{\mathbf{X}}_{k,l}. \quad (17)$$

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. (16) and (17). 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

$$\frac{\partial L^+}{\partial w_{i,k,z_{i,l}}} = \sum_l \left(\frac{\operatorname{tr}(\mathbf{y}_{k,l}^H \mathbf{y}_{k,l} \mathbf{R}_{i,k,l} \mathbf{C}_{i,k}^{-1} \mathbf{R}_{i,k,l})}{w_{i,k,z_{i,l}}^2 h_{i,l}} - \operatorname{tr}(\mathbf{U}_{k,l}^{-1} \mathbf{C}_{i,k}) h_{i,l} \right), \quad (18)$$

$$\begin{aligned} \frac{\partial L^+}{\partial h_{i,l}} &= \sum_k \left(\frac{\operatorname{tr}(\mathbf{y}_{k,l}^H \mathbf{y}_{k,l} \mathbf{R}_{i,k,l} \mathbf{C}_{i,k}^{-1} \mathbf{R}_{i,k,l})}{w_{i,k,z_{i,l}} h_{i,l}^2} \right. \\ &\quad \left. - \operatorname{tr}(\mathbf{U}_{k,l}^{-1} \mathbf{C}_{i,k}) w_{i,k,z_{i,l}} \right) \\ &\quad + (\alpha_{z_{i,l}} - 1) / h_{i,l} - 1 / \beta_{z_{i,l}}, \end{aligned} \quad (19)$$

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

$$w_{i,k,z_{i,l}} \leftarrow \sqrt{\frac{\sum_l \operatorname{tr}(\mathbf{y}_{k,l}^H \mathbf{y}_{k,l} \mathbf{R}_{i,k,l} \mathbf{C}_{i,k}^{-1} \mathbf{R}_{i,k,l})}{\sum_l \operatorname{tr}(\mathbf{U}_{k,l}^{-1} \mathbf{C}_{i,k}) h_{i,l}}}, \quad (20)$$

$$h_{i,l} \leftarrow \frac{(\alpha_{z_{i,l}} - 1) + \sqrt{(\alpha_{z_{i,l}} - 1)^2 + 4\mu_{i,l}\nu_{i,l}}}{2\mu_{i,l}}, \quad (21)$$

where

$$\mu_{i,l} = \sum_k \frac{\operatorname{tr}(\mathbf{y}_{k,l}^H \mathbf{y}_{k,l} \mathbf{R}_{i,k,l} \mathbf{C}_{i,k}^{-1} \mathbf{R}_{i,k,l})}{w_{i,k,z_{i,l}}}, \quad (22)$$

$$\nu_{i,l} = \sum_k \operatorname{tr}(\mathbf{U}_{k,l}^{-1} \mathbf{C}_{i,k}) w_{i,k,z_{i,l}} + 1 / \beta_{z_{i,l}}. \quad (23)$$

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

$$\frac{\partial L^+}{\partial \mathbf{C}_{i,k}} = \sum_l \left(\frac{\mathbf{C}_{i,k}^{-1} \mathbf{R}_{i,k,l} \mathbf{y}_{k,l}^H \mathbf{y}_{k,l} \mathbf{R}_{i,k,l} \mathbf{C}_{i,k}^{-1}}{w_{i,k,z_{i,l}} h_{i,l}} \right)$$

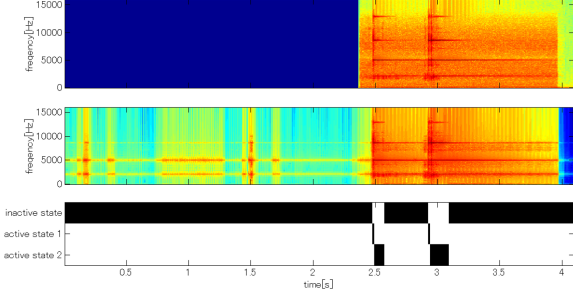


Figure 3: Examples of spectrogram of signal of a bell (top), that of separated signal (middle), and its activity detection result (bottom). Black indicates the state is assigned at the time.

Table 1: The average output SDRs and standard deviations of the three sources by the conventional and proposed methods.

SDR(\pm SD) [dB]	bell	whistle	stapler
Proposed	19.92(\pm 11.50)	31.21(\pm 6.69)	16.81(\pm 10.09)
Conventional	13.33(\pm 8.22)	23.37(\pm 4.57)	8.28(\pm 5.98)

$$- \mathbf{U}_{k,i}^{-1} w_{i,k,z_{i,l}} h_{i,l} \Big). \quad (24)$$

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

$$\mathbf{C}_{i,k} \mathbf{A}_{i,k} \mathbf{C}_{i,k} = \mathbf{B}_{i,k}, \quad (25)$$

where

$$\mathbf{A}_{i,k} = \sum_l w_{i,k,z_{i,l}} h_{i,l} \hat{\mathbf{X}}_{k,l}^{-1}, \mathbf{B}_{i,k} = \mathbf{C}_{i,k} \left(\sum_l w_{i,k,z_{i,l}} h_{i,l} \hat{\mathbf{X}}_{k,l}^{-1} \mathbf{y}_{k,l} \mathbf{y}_{k,l}^H \right) \mathbf{C}_{i,k}. \quad (26)$$

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

$$\begin{bmatrix} 0 & -\mathbf{A}_{i,k} \\ -\mathbf{B}_{i,k} & 0 \end{bmatrix}, \quad (27)$$

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

$$\mathbf{e}_{m,i,k} = \begin{bmatrix} \mathbf{f}_{m,i,k} \\ \mathbf{g}_{m,i,k} \end{bmatrix} \quad (28)$$

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

$$\mathbf{C}_{i,k} \leftarrow \mathbf{G}_{i,k} \mathbf{F}_{i,k}^{-1}, \quad (29)$$

where $\mathbf{F}_{i,k} = [\mathbf{f}_{1,i,k}, \dots, \mathbf{f}_{M,i,k}]$ and $\mathbf{G}_{i,k} = [\mathbf{g}_{1,i,k}, \dots, \mathbf{g}_{M,i,k}]$.

L^+ is equal up to constant terms to the sum of the log posteriors of HMMs, when viewed as a function of \mathbf{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 \mathbf{W} , \mathbf{H} , \mathbf{C} and \mathbf{Z} corresponds to solving the problems of blind signal separation and source activity detection based on a unified objective function.

4. Experimental evaluation

We evaluated the performance of the present method in terms of the ability to separate sources and detect their activity. We used a mixed stereo signal as the experimental data, each of which we obtained by mixing the non-speech signals (sounds of a whistle, a bell and a stapler) from the RWCP database[8]

and was convolved with the measured room impulse response from the RWCP database [8] (in which the distance between the microphones was 5.85 cm and the reverberation time was 0 ms). Thus, the three signals were artificially located 20° , 60° and 100° from the microphones respectively. 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. Fig. 2 shows a spectrogram of the mixed signal. We set the number of states of HMMs D as 3. We expected that $d = 1$ is an inactive state and $d = 2$ and 3 are active states, by setting α_1 and β_1 as 1 and 10^{-2} respectively, and $\alpha_{2:3}$ and $\beta_{2:3}$ 1 and 10^{20} respectively. the transition probability ρ as $\rho_1 = (0.9, 0.1, 0)$, $\rho_2 = (0, 0.5, 0.5)$ and $\rho_3 = (0.5, 0, 0.5)$ based on the assumption that a sound source tends to stay inactive. The diagonal elements of \mathbf{C} were initially all set to $1/\sqrt{M}$, and the off-diagonal elements were initially set to zero. We first performed the single channel NMF based on IS-divergence starting from different random initial matrices \mathbf{W} and \mathbf{H} , then we set the results of \mathbf{W} , \mathbf{H} as the initial parameters of $\hat{\mathbf{W}}$ and $\hat{\mathbf{H}}$. The parameter estimation algorithm was run for 100 iterations. We chose the method proposed in [5] as a comparison. The separated signal $\hat{\mathbf{y}}_{i,k,l}$ was obtained by Wiener filtering

$$\hat{\mathbf{y}}_{i,k,l} = w_{i,k,z_{i,l}} h_{i,l} \mathbf{C}_{i,k} \hat{\mathbf{X}}_{k,l}^{-1} \mathbf{y}_{k,l}. \quad (30)$$

As an evaluation measure, we used the signal-to-distortion ratio (SDR) [9]. The SDR is expressed in decibels (dB), and a higher SDR indicates superior quality. The input SDRs [dB] of the bell, the whistle and the stapler were -8.36, 6.64 and -12.78 [dB], respectively.

Table 1 shows the average SDRs and standard deviations for the ten trials obtained by the conventional and proposed methods. The average SDRs of the proposed method were superior to those of the conventional method for each signal. The total average of the SDRs obtained with the proposed method was 7.65 dB more than that obtained with the conventional approach. These results show the effectiveness of the proposed method for BSS. Examples of separated signals are available at http://www.hil.t.u-tokyo.ac.jp/~higuchi/demo/Examples_interspeech.htm. Fig. 3 shows examples of spectrogram of signal of a bell (top), that of separated signal (middle), and its activity detection result (bottom). The activities of the bell were almost all estimated correctly. However, our results also showed that abilities of signal separation and source activity detection by our method depended on the first condition of \mathbf{W} and \mathbf{H} to some extent. As future work, we plan to solve this kind of problem by learning \mathbf{W} and ρ from a clean signal and apply our proposed method to acoustic event detection.

5. Conclusion

In this paper we extend our previous approach to a multichannel observed signal and proposes a comprehensive approach to deal with an underdetermined BSS problem and source activity detection. Specifically, we describe a generative model of mixture signals by incorporating a generative model of spectra for each source, using MFHMMs. Each hidden state of the present HMM represents states of each source such as its activities. Through the estimation of the parameters of the overall generative model, we can simultaneously the underdetermined blind signal separation and source activity detection. The experiment showed that the proposed algorithm provided a 7.65 dB improvement compared with the conventional method in terms of the signal-to-distortion ratio.

6. Acknowledgement

This work was supported by JSPS KAKENHI Grant Number 26730100.

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. Smaragdis, and J. C. Brown, “Non-negative matrix factorization for polyphonic music transcription,” *Proc. WASPAA 2003*, Oct. 2003, pp. 177–180.
- [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,” *IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP) 2012*, pp. 261-264, 2012.
- [6] 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.
- [7] Masahiro Nakano, Le Roux Jonathan, Hirokazu Kameoka, Yu Kitano, Nobutaka Ono, and Shigeki Sagayama, “Nonnegative matrix factorization with Markov-chained bases for modeling time-varying patterns in music spectrograms,” *Latent Variable Analysis and Signal Separation*, vol. 6365, pp. 149–156, 2010.
- [8] 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.
- [9] 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.