

Bayesian Nonparametric Approach to Blind Separation of Infinitely Many Sparse Sources

Hirokazu KAMEOKA^{†,††a)}, Member, Misa SATO^{†††*}, Takuma ONO[†], Nobutaka ONO^{††††}, and Shigeki SAGAYAMA[†], Nonmembers

SUMMARY This paper deals with the problem of underdetermined blind source separation (BSS) where the number of sources is unknown. We propose a BSS approach that simultaneously estimates the number of sources, separates the sources based on the sparseness of speech, estimates the direction of arrival of each source, and performs permutation alignment. We confirmed experimentally that reasonably good separation was obtained with the present method without specifying the number of sources.

key words: underdetermined blind signal separation, speech, sparseness, permutation alignment, Bayesian nonparametrics, direction of arrival, Dirichlet process, stick-breaking construction, variational inference

1. Introduction

Blind Source Separation (BSS) is a technique for separating out individual source signals from microphone inputs when the transfer characteristics between sources and microphones are unknown. BSS is potentially useful for the development of such applications as hands-free teleconference systems and automatic meeting transcription systems. In meeting situations, for instance, it is likely that the number of participants (speakers) will change during the meeting or a loud, unexpected noise such as a door slamming will occur in the room. Thus, it is often difficult to pre-specify the exact number of all possible sources present in real environments. Many conventional BSS algorithms are designed to use the number of sources as the input when performing separation, and most of these algorithms do not work well if the assumed and actual numbers of sources are not the same. This paper presents a novel BSS algorithm that allows the number of sources to be inferred along with the separation.

To estimate the unknown mixing matrix and source signals solely from observed signals, we must 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], which estimates the separation matrix (the inverse of the mixing matrix) such that the independence of the source estimates is maximized. However, to achieve a BSS algorithm that works without assuming the number of sources, we must always consider an underdetermined case where there are fewer observations than sources. In an underdetermined case, there are an infinite number of solutions for source signals even if the mixing process is known. The independence assumption is too weak to allow us to determine a unique solution and so directly applying ICA will not work in this case. In underdetermined situations, we typically need a stronger assumption than independence.

One successful approach for underdetermined BSS involves utilizing the sparse nature of speech [2]–[7], which relies on the fact that the time-frequency components of speech are near zero across most of the time-frequency grid points. Since the time-frequency components of speech rarely overlap even when multiple speakers are speaking simultaneously, the main focus of this approach is how to design a time-frequency mask with which we can extract only the components of target speech from the mixture.

The signals observed at each microphone can be modeled as a convolutive mixture of source signals. To exploit the sparse nature of speech, we must convert it to a time-frequency representation. If we assume the use of a short-time Fourier transform (STFT) to obtain a time-frequency representation with a frame length sufficiently longer than the length of the impulse response from a source to a microphone, an observed signal can be approximated fairly well by an instantaneous mixture in the frequency domain. BSS based on this observation model is called frequency domain BSS. While frequency domain BSS allows for a fast implementation compared with BSS that uses a time domain convolutive mixture model, it requires us to solve an additional problem called the permutation alignment problem. That is, we must group together the separated components of different frequency bins that are considered to originate from the same source in order to construct a separated signal. Some methods such as [7] are designed to perform frequency-bin-wise source separation followed by permutation alignment. However, permutation alignment and source separation problems should be solved in a cooperative manner since the clues used for permutation alignment can also be helpful

Manuscript received January 22, 2013.

Manuscript revised May 23, 2013.

[†]The authors are with the Graduate School of Information Science and Technology, The University of Tokyo, Tokyo, 113-8656 Japan.

^{††}The author is with NTT Communication Science Laboratories, NTT Corporation, Atsugi-shi, 243-0198 Japan.

^{†††}The author is with School of Engineering, The University of Tokyo, Tokyo, 113-8656 Japan.

^{††††}The author is with National Institute of Informatics, Tokyo, 101-8430 Japan.

^{*}Presently, with Graduate School of Engineering, The University of Tokyo.

a) E-mail: kameoka@hil.t.u-tokyo.ac.jp

DOI: 10.1587/transfun.E96.A.1928

for source separation. Thus, we consider it important to develop a method that can simultaneously perform frequency-bin-wise source separation and permutation alignment.

Motivated by the above, this paper proposes a novel BSS approach that simultaneously performs (1) an estimation of the number of sources, (2) source separation based on the sparseness of speech, and (3) permutation alignment.

2. Mixing Model

We first consider a situation where N source signals are captured 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. ω_k is the angular frequency corresponding to the k -th frequency bin and t_l is the time of the l -th frame, respectively. If we assume that the length of the impulse response from a source to a microphone 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:

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

where $\mathbf{y}(\omega_k, t_l) = (y_1(\omega_k, t_l), \dots, y_M(\omega_k, t_l))^T$. $\mathbf{a}_i(\omega_k) = (a_{i,1}(\omega_k), \dots, a_{i,M}(\omega_k))^T$ is the frequency array response for source i , which is assumed to be time-invariant. $\mathbf{n}(\omega_k, t_l)$ is assumed to comprise all kinds of components that cannot be expressed by the instantaneous mixture representation (e.g., background noise and reverberant components).

We now exploit the sparseness of speech and assume that only one source is active at each time-frequency point, as with [2]–[7]. By using $z_{k,l} \in \{1, \dots, N\}$ to denote the (unknown) active source index at time-frequency point (ω_k, t_l) , (1) can be rewritten as

$$\mathbf{y}(\omega_k, t_l) = \mathbf{a}_{z_{k,l}}(\omega_k) s(\omega_k, t_l) + \mathbf{n}(\omega_k, t_l). \quad (2)$$

Notice that the superscript i is dropped from $s_i(\omega_k, t_l)$ in (2) as it is no longer necessary since we are assuming $s_i(\omega_k, t_l) = 0$ for $i \neq z_{k,l}$. Namely, $s(\omega_k, t_l)$ signifies the component of an active source at time-frequency point (ω_k, t_l) . For convenience of notation, we hereafter use subscripts k and l to indicate ω_k and t_l .

3. Generative Model

3.1 Generative Process of Observed Signals

Here we describe the generative process of an observed signal on the basis of (2). Let us assume that the noise component $\mathbf{n}_{k,l}$ follows a complex normal distribution with mean $\mathbf{0}$ and covariance $\Sigma_k^{(n)}$. Then, from (2), $\mathbf{y}_{k,l}$ is also normally distributed such that

$$\mathbf{y}_{k,l} | \mathbf{a}_{1:N,k}, s_{k,l}, z_{k,l} \sim \mathcal{N}_{\mathbb{C}}(\mathbf{y}_{k,l}; \mathbf{a}_{z_{k,l},k} s_{k,l}, \Sigma_k^{(n)}), \quad (3)$$

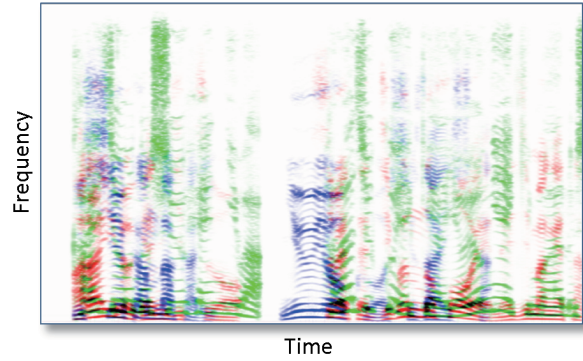


Fig. 1 Spectrograms of speech signals uttered by three speakers, each of which is indicated by a different color. Typically, the time-frequency components of speech rarely overlap even when multiple speakers are speaking simultaneously.

conditioned on $\mathbf{a}_{1:N,k} = \{\mathbf{a}_{1,k}, \dots, \mathbf{a}_{N,k}\}$, $s_{k,l}$ and $z_{k,l}$, where $\mathcal{N}_{\mathbb{C}}(\mathbf{x}; \boldsymbol{\mu}, \Sigma) \propto \exp(-(\mathbf{x} - \boldsymbol{\mu})^H \Sigma^{-1} (\mathbf{x} - \boldsymbol{\mu}))$. Using this model, Izumi et al. formulated an Expectation-Maximization (EM) algorithm for the maximum likelihood estimation of the unknown parameters by treating $z_{k,l}$ as a latent variable [6], [7]. To flexibly incorporate hierarchical prior distributions over the unknown parameters as well as the number of sources, we employ a Bayesian approach.

3.2 Mixture of Infinitely Many Sparse Sources

We do not usually obtain any information about which source is active at each time-frequency point. Thus, we regard $z_{k,l}$ as a latent variable and consider describing its generative process. If the number of sources is N , it would be natural to assume that the probability of choosing an index $z_{k,l}$ from the set of source indices, $\{1, \dots, N\}$, can be described as a categorical distribution

$$z_{k,l} | \boldsymbol{\pi} \sim \text{Categorical}(z_{k,l}; \boldsymbol{\pi}), \quad (4)$$

where $\text{Categorical}(x; \mathbf{y}) = y_x$ (with $\mathbf{y} = (y_1, \dots, y_N)$ and $\sum_{i=1}^N y_i = 1$). The i -th element of $\boldsymbol{\pi}$ defines how likely the source index i is to be chosen. Since we also have no information about $\boldsymbol{\pi}$, we consider describing its generative process using a “symmetric” distribution. For the convenience of the following analysis, we assume that $\boldsymbol{\pi}$ has been generated from a symmetric Dirichlet distribution

$$\boldsymbol{\pi} \sim \text{Dirichlet}(\boldsymbol{\pi}; \alpha_0/N, \dots, \alpha_0/N), \quad (5)$$

where $\text{Dirichlet}(\mathbf{y}; z_1, \dots, z_N) \propto \prod_{i=1}^N y_i^{z_i-1}$. The shape of the Dirichlet distribution is governed by a concentration hyperparameter α_0 .

Thus far, we have considered the case of a finite number N of sources. It can be shown that with a size-biased permutation of (5) followed by taking the limit $N \rightarrow \infty$, the above generative processes (3), (4) and (5) approach

$$\boldsymbol{\pi} \sim \text{GEM}(\boldsymbol{\pi}; \alpha_0), \quad (6)$$

$$z_{k,l} | \boldsymbol{\pi} \sim \text{Categorical}(z_{k,l}; \boldsymbol{\pi}), \quad (7)$$

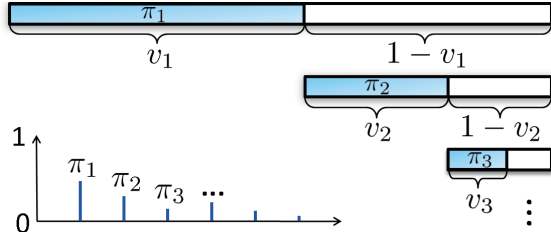


Fig. 2 Illustration of the stick-breaking construction. The construction of $\pi = (\pi_1, \pi_2, \dots, \pi_\infty)$ can be understood metaphorically as follows. Starting with a stick of length 1, we break it at v_1 , assigning π_1 as the length of the stick we just broke off. We now recursively break the other portion to obtain π_2, π_3 and so forth.

$$y_{k,l} | a_{1:\infty,k} s_{k,l}, z_{k,l} \sim \mathcal{N}_{\mathbb{C}}(y_{k,l}; a_{z_{k,l},k} s_{k,l}, \Sigma_k^{(n)}). \quad (8)$$

$\text{GEM}(\pi; \alpha_0)$ is called the stick-breaking process [8], where the letters stand for Griffiths, Engen and McCloskey. For $i = 1, 2, \dots, \infty$, the i -th mixture weight π_i is generated via

$$v_i \sim \text{Beta}(v_i; 1, \alpha_0), \quad (9)$$

$$\pi_i = v_i \prod_{j=1}^{i-1} (1 - v_j), \quad (10)$$

where $\text{Beta}(y; z_1, z_2) \propto y^{z_1-1} (1-y)^{z_2-1}$. This is known as a constructive definition of the Dirichlet process [9], which can be thought of as an infinite dimensional Dirichlet distribution. The construction of $\pi = (\pi_1, \pi_2, \dots, \pi_\infty)$ can be understood metaphorically as follows. Starting with a stick of length 1, we break it at v_1 , assigning π_1 as the length of the piece of stick that we have just broken off. Now recursively break the other portion to obtain π_2, π_3 and so forth (see Fig. 2). $\pi \sim \text{GEM}(\pi; \alpha_0)$ thus produces exponentially decaying weights in expectation. This means that the source with a larger index will be less likely to be active and thus simple models with fewer sources are favored, given observed signals.

3.3 Mixture of Direction-of-Arrivals (DOAs)

Here we describe the generative process of the frequency response $a_{i,k}$ of the mixing system.

So far we have treated $a_{i,k}$ as an independent parameter across k . If the index i indicates an identical source across $\omega_1, \dots, \omega_K$, $a_{i,k}$ will have a certain structure that can be described using the property of acoustic wave propagation. We thus expect that the incorporation of an appropriate constraint into $a_{i,k}$ would help solve both the permutation alignment problem and the frequency-wise source separation problem simultaneously through parameter inference. If each source is assumed to be located far from the microphones so that the signal can be treated approximately as a plane wave, the interchannel time difference between the microphones depends only on the direction of arrival (DOA) of the source. Since the time delay between two microphones corresponds to the phase difference of the frequency response of the microphone array, the complex array response can be expressed explicitly by using the DOAs

of the source. Specifically, with $M = 2$ microphones, the complex array response for a source at direction θ such that $0 \leq \theta < 2\pi$ is defined as a function of ω depending on θ

$$h(\theta, \omega) = \begin{bmatrix} 1 \\ e^{j\omega B \cos \theta / C} \end{bmatrix}, \quad (11)$$

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 frequency response $a_{i,k}$ should be equal to $h(\theta_i, \omega_k)$. However, due to such factors as the plane wave assumption and the narrowband instantaneous mixture approximation, the actual frequency response $a_{i,k}$ may diverge from the “ideal frequency response” $h(\theta_i, \omega_k)$ to some extent. One way to simplify the process of this kind of deviation is to assume a probability distribution on $a_{i,k}$ with the expected value of $h(\theta_i, \omega_k)$. Here, we assume for convenience that $a_{i,k}$ is generated from a complex normal distribution with mean $h(\theta_i, \omega_k)$. Note that we do not usually obtain any information about the direction from which each source is emanating. Thus, we regard the DOA of each source as a latent variable and further consider describing its generative process. As explained in detail below, the entire generative process of $a_{i,k}$ can then be described as a “mixture of DOAs”.

Let us firstly introduce a discrete set of D possible directions, $\vartheta_1, \dots, \vartheta_D$, which are all assumed to be constants. For instance, consider defining ϑ_d as $\vartheta_d = (d-1)\pi/D$, ($d = 1, \dots, D$), which means dividing π into D equal angles. We then assume that each source is emanating from one of these directions. First, we consider the generative process of the DOA θ_i of source i . For each source i , an index c_i of direction is drawn according to a categorical distribution $\rho = (\rho_1, \dots, \rho_D)$

$$c_i | \rho \sim \text{Categorical}(c_i; \rho). \quad (12)$$

By using c_i , θ_i is then given as

$$\theta_i = \vartheta_{c_i}. \quad (13)$$

As with 3.2, we assume that ρ has been generated from a symmetric Dirichlet distribution

$$\rho \sim \text{Dirichlet}(\rho; \beta_0/D, \dots, \beta_0/D). \quad (14)$$

As mentioned above, the frequency response $a_{i,k}$ is assumed to be generated from a complex normal distribution with mean $h(\vartheta_{c_i}, \omega_k)$, given c_i ,

$$a_{i,k} | c_i \sim \mathcal{N}_{\mathbb{C}}(a_{i,k}; h(\vartheta_{c_i}, \omega_k), \Sigma_k^{(a)}), \quad (15)$$

where $\Sigma_k^{(a)}$ denotes the covariance of the complex normal distribution, which is assumed to be a constant. An illustration of the generative process of $a_{i,k}$ described above is shown in Fig. 3.

Overall, the entire generative model is described in plate notation in Fig. 4.

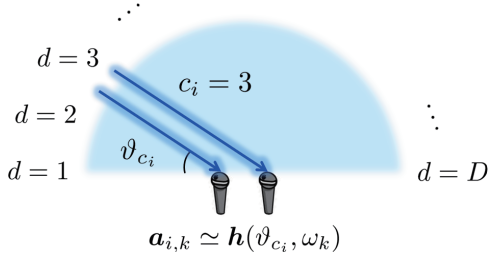


Fig. 3 Illustration of the generative process of the complex array response $\mathbf{a}_{i,k}$ based on the DOA mixture model.

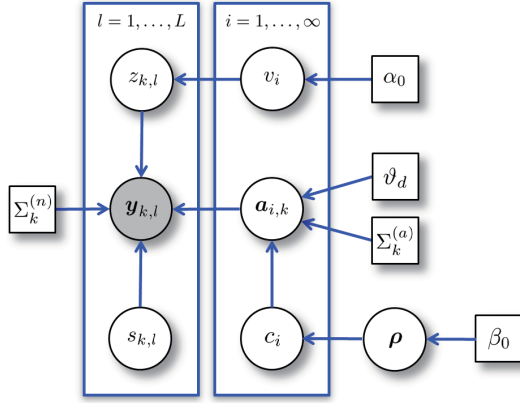


Fig. 4 Plate notation of present overall generative model.

4. Approximate Posterior Inference

4.1 Variational Bayesian Approach

In this section, we describe an approximate posterior inference algorithm for our generative model based on variational inference [10]. The random variables of interest in our model are

- $A = \{\mathbf{a}_{i,k}\}_{i,k}$: complex array response for source i ,
- $S = \{s_{k,l}\}_{k,l}$: component of active source at (ω_k, t_l) ,
- $Z = \{z_{k,l}\}_{k,l}$: index of active source at (ω_k, t_l) ,
- $V = \{v_i\}_i$: stick breaking proportion,
- $C = \{c_i\}_i$: index of direction for source i ,
- $\rho = (\rho_1, \dots, \rho_D)$: mixture weight for each DOA.

We denote the entire set of the above parameters as Θ . Our goal is to compute the posterior $p(\Theta|Y)$ where $Y = \{\mathbf{y}_{k,l}\}_{k,l}$ is a set consisting of the time-frequency components of observed multichannel signals. By using the conditional distributions defined in 3.2 and 3.3, we can write the joint distribution $p(Y, \Theta)$ as

$$p(Y, A, S, Z, V, C, \rho) = p(Y|A, S, Z)p(Z|V)p(V)p(A|C)p(C|\rho)p(\rho), \quad (16)$$

but to obtain the exact posterior $p(\Theta|Y)$, we must compute $p(Y)$, which involves many intractable integrals.

We can express this posterior variationally as the solution to an optimization problem:

$$\operatorname{argmin}_{q \in \mathcal{Q}} \text{KL}(q(\Theta) \| p(\Theta|Y)), \quad (17)$$

where $\text{KL}(\cdot \| \cdot)$ denotes the Kullback-Leibler (KL) divergence between its two arguments, i.e.,

$$\text{KL}(q(\Theta) \| p(\Theta|Y)) = \int q(\Theta) \log \frac{q(\Theta)}{p(\Theta|Y)} d\Theta. \quad (18)$$

Indeed, if we let \mathcal{Q} be the family of all distributions over Θ , the solution to the optimization problem is the exact posterior $p(\Theta|Y)$, since KL divergence is minimized when its two arguments are exactly equal. Of course, solving this optimization problem is just as intractable as directly computing the posterior. Although it may appear that no progress has been made, restricting $q(\Theta)$ to belong to a family of distributions of simpler form than $p(\Theta|Y)$ allows us to obtain principled approximate solutions.

For our model, we define the set of approximate distributions \mathcal{Q} as those that factor as follows:

$$\mathcal{Q} = \{q : q(A)q(S)q(Z)q(V)q(C)q(\rho)\}. \quad (19)$$

This approximation is often called a naive mean-field approximation. To define $q(A)$, $q(V)$ and $q(C)$, we need to construct distributions on the infinite sets $\{v_1, v_2, \dots\}$, $\{\mathbf{a}_{1,k}, \mathbf{a}_{2,k}, \dots\}$ and $\{c_1, c_2, \dots\}$. For this approach to be tractable, we truncate the variational distribution at some value N^* by setting $q(v_{N^*+1} = 1) = 1$. The mixture proportions π_i for $i > N^*$ will thus be zero, and we can ignore $\mathbf{a}_{i,k}$ and c_i for $i > N^*$. In practice, we set N^* at a sufficiently large integer. It is important to emphasize that truncating the variational distribution does not mean that the true posterior itself is truncated. As the truncation level N^* becomes larger, the approximations to the true posterior become more accurate.

4.2 Coordinate Ascent

We now present an algorithm for solving the optimization problem described in (17) and (19). Unfortunately, the optimization problem is non-convex, and it is intractable to find the global optimum. However, we can use a simple coordinate ascent algorithm to find a local optimum. Notice that (18) can be written as

$$\begin{aligned} \text{KL}(q(\Theta) \| p(\Theta|Y)) &= \int q(\Theta) \log \frac{q(\Theta)}{p(Y, \Theta)} d\Theta + \log p(Y). \end{aligned} \quad (20)$$

As the log evidence $\log p(Y)$ is fixed with respect to $q(\Theta)$, minimizing the first term, which is known as the (negative) variational free energy, minimizes the KL divergence of $p(\Theta|Y)$ from $q(\Theta)$. The algorithm optimizes one factor in the mean-field approximation of the posterior at a time while fixing all other factors. It can be shown using the calculus of variations that the “optimal” distribution for each of the factors can be expressed as:

$$\hat{q}(X) \propto \exp \mathbb{E}_{\Theta \setminus X} [\log p(Y, \Theta)], \quad (21)$$

where X indicates one of the factors and $\mathbb{E}_{\Theta \setminus X}[\log p(Y, \Theta)]$ is the expectation of the joint probability of the data and latent variables, taken over all variables except X . The mean-field update equations for the variational distributions are given in the following form:

$$\hat{q}(A) \leftarrow \prod_{i,k} \mathcal{N}_{\mathbb{C}}(\mathbf{a}_{i,k}; \mathbf{m}_{i,k}, \Gamma_{i,k}), \quad (22)$$

$$\hat{q}(S) \leftarrow \prod_{k,l} \mathcal{N}_{\mathbb{C}}(s_{k,l}; \mu_{k,l}, \sigma_{k,l}^2), \quad (23)$$

$$\hat{q}(Z) \leftarrow \prod_{k,l} \text{Categorical}(z_{k,l}; \boldsymbol{\phi}_{k,l}), \quad (24)$$

$$\hat{q}(V) \leftarrow \prod_i \text{Beta}(v_i; \gamma_{i,0}, \gamma_{i,1}), \quad (25)$$

$$\hat{q}(C) \leftarrow \prod_i \text{Categorical}(c_i; \boldsymbol{\psi}_i), \quad (26)$$

$$\hat{q}(\boldsymbol{\rho}) \leftarrow \text{Dirichlet}(\boldsymbol{\rho}; \zeta_1, \dots, \zeta_D), \quad (27)$$

where

$$\Gamma_{i,k}^{-1} = \left(\sum_l [\boldsymbol{\phi}_{k,l}]_i (|\mu_{k,l}|^2 + \sigma_{k,l}^2) \right) \Sigma_k^{(n)-1} + \Sigma_k^{(a)-1}, \quad (28)$$

$$\mathbf{m}_{i,k} = \Gamma_{i,k} \left(\Sigma_k^{(n)-1} \sum_l [\boldsymbol{\phi}_{k,l}]_i \mu_{k,l}^* \mathbf{y}_{k,l} + \Sigma_k^{(a)-1} \sum_d [\boldsymbol{\psi}_d]_i \mathbf{h}(\vartheta_d, \omega_k) \right), \quad (29)$$

$$\frac{1}{\sigma_{k,l}^2} = \sum_i [\boldsymbol{\phi}_{k,l}]_i \text{tr} \left[(\mathbf{m}_{i,k} \mathbf{m}_{i,k}^H + \Gamma_{i,k}) \Sigma_k^{(n)-1} \right], \quad (30)$$

$$\mu_{k,l} = \sigma_{k,l}^2 \left(\sum_i [\boldsymbol{\phi}_{k,l}]_i \mathbf{m}_{i,k}^H \right) \Sigma_k^{(n)-1} \mathbf{y}_{k,l}, \quad (31)$$

$$\begin{aligned} \varphi_{i,k,l} = \exp & \left[2\text{Re}(\mu_{k,l} \mathbf{y}_{k,l}^H \Sigma_k^{(n)-1} \mathbf{m}_{i,k}) \right. \\ & - (|\mu_{k,l}|^2 + \sigma_{k,l}^2) \text{tr} \left[(\mathbf{m}_{i,k} \mathbf{m}_{i,k}^H + \Gamma_{i,k}) \Sigma_k^{(n)-1} \right] \\ & + \Psi(\gamma_{i,0}) - \Psi(\gamma_{i,0} + \gamma_{i,1}) \\ & \left. + \sum_{j=1}^{i-1} (\Psi(\gamma_{j,1}) - \Psi(\gamma_{j,0} + \gamma_{j,1})) \right], \end{aligned} \quad (32)$$

$$[\boldsymbol{\phi}_{k,l}]_i = \frac{\varphi_{i,k,l}}{\sum_{i'=1}^{N^*} \varphi_{i',k,l}}, \quad (32)$$

$$\gamma_{i,0} = 1 + \sum_{k,l} [\boldsymbol{\phi}_{k,l}]_i, \quad (33)$$

$$\gamma_{i,1} = \alpha_0 + \sum_{j=i+1}^{N^*} \sum_{k,l} [\boldsymbol{\phi}_{k,l}]_i, \quad (34)$$

$$\begin{aligned} \varrho_{i,d} = \exp & \left[\sum_k \left\{ 2\text{Re}(\mathbf{h}(\vartheta_d, \omega_k)^H \Sigma_k^{(a)-1} \mathbf{m}_{i,k}) \right. \right. \\ & \left. \left. - \text{tr}[\mathbf{h}(\vartheta_d, \omega_k) \mathbf{h}(\vartheta_d, \omega_k)^H \Sigma_k^{(a)-1}] \right\} + \Psi(\zeta_d) \right], \end{aligned} \quad (35)$$

$$[\boldsymbol{\psi}_i]_d = \frac{\varrho_{i,d}}{\sum_{d'=1}^D \varrho_{i,d'}}, \quad (35)$$

$$\zeta_d = \frac{\beta_0}{D} + \sum_{i=1}^{N^*} \psi_{i,d}. \quad (36)$$

$[\cdot]_i$ denotes the i -th element of a vector, $\text{tr}[\cdot]$ denotes the trace of a matrix, and $\Psi(\cdot)$ denotes the digamma function.

Finally, the STFT components of the i -th separated signal can be obtained by $\phi_{i,k,l} \mu_{k,l}$. Since $q(\Theta)$ is an approximation to the true posterior $p(\Theta|Y)$, $\phi_{i,k,l} \mu_{k,l}$ corresponds to an approximation of the minimum mean square error estimator of the i -th source signal, i.e., $\hat{s}_{i,k,l} = \mathbb{E}[\mathbf{1}[z_{k,l} = i] s_{k,l} | Y] \simeq \mathbb{E}[z_{k,l} = i | Y] \mathbb{E}[s_{k,l} | Y] = \phi_{i,k,l} \mu_{k,l}$, where $\mathbf{1}[\cdot]$ denotes the indicator function that takes the value 1 if its argument is true and 0 otherwise.

5. Related Work

It has been brought to our attention that several papers related to this work have been presented independently by different authors [11], [12], after the publication of our conference paper [17].

Taghia et al. proposed in [11] a variational Bayesian approach for the approximate inference of the parameters of Sawada's model [7]. This method requires a model selection procedure to estimate the number of sources. The variational lower bound (an approximation of the log marginal likelihood) is used as the criterion for the model selection. By contrast, our method avoids selecting the number of sources explicitly and instead pushes this task down into the inference algorithm, incorporating automated Bayesian model selection in the inference procedure.

Otsuka et al. proposed in [12] a combined source separation and source localization method based on a Bayesian approach. Similar to the idea described in 3.3, they described the prior distribution of the spatial covariance matrix of each source using a mixture model of Wishart-distributed spatial covariance matrices, each of which is associated with a latent DOA variable. By contrast, the present model describes the prior distribution of the frequency array response of each source using a mixture model of Gaussian-distributed frequency array responses, each of which is associated with a latent DOA variable. Although these models appear slightly different, the concepts are very similar. It would be interesting to compare the different ways in which the use of these models affect the source separation ability. This should be investigated in the future.

6. Experiment

We evaluated the performance of the proposed method in terms of source separation ability.

We used stereo speech signals with a sampling rate of 16 kHz as test signals, which we obtained by mixing three speech signals [13] (one male and two female speakers) using a measured room impulse response [14] (in which the distance between the microphones was 5 cm and the reverberation time was 0 ms). The three sources were spaced 30 degrees apart. The geometry setting for the test mixtures

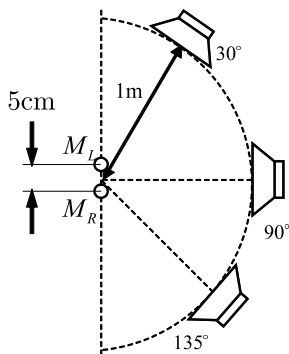


Fig. 5 Geometry setting for text mixtures.

is shown in Fig. 5. 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 16 ms.

Our preliminary experiments revealed that ζ_d was likely to be trapped in local optima due to the spatial aliasing that occurs at high frequencies. To avoid this, we adopted the following procedure: We first ran the variational inference algorithm using only the low-frequency region of the observed signals, after which we gradually increased the frequency range up to the Nyquist frequency during the iteration. $\Sigma_k^{(n)}$ and $\Sigma_k^{(a)}$ were set respectively at \mathbf{I} and $10^{-1.5} \times \mathbf{I}$. D was set at 180. All the variational parameters were initialized randomly. After convergence, each separated signal was obtained by multiplying $\mu_{k,l}$ by $\phi_{i,k,l}$. Figure 6 shows some examples of the spectrograms of separated source signals obtained with the present method. Owing to the property of the stick-breaking process, sources that are dominant at many time-frequency points are likely to be assigned small indices. Thus, we considered the separated signals of indices $i = 1, 2, 3$ to be the source estimates, which we used for the evaluation. We chose Sawada's method described in [7] as a comparison. The source code was kindly provided to us by its author. In this method, the number of sources must be specified manually. The following results report the performance in terms of the Signal-to-Distortion Ratio (SDR) [15]. The SDR is expressed in decibels (dB), and a higher SDR indicates superior quality.

The present method was tested with various settings of the truncation level N^* . As for Sawada's method, it was tested with various settings of the assumed number N of sources. Figure 7 shows the average SDRs obtained with Sawada's and the present methods with various N and N^* settings. As expected, the performance of the present method improves with increasing N^* , while that of Sawada's method deteriorates significantly when the assumed number of sources departs from the actual number. It is worth emphasizing that Sawada's method requires prior knowledge of the exact number of sources to achieve good separation, while the present approach does not. With Sawada's method, each source signal tends to be further separated into multiple components when the specified number of sources is greater than the actual number. Unlike with Sawada's

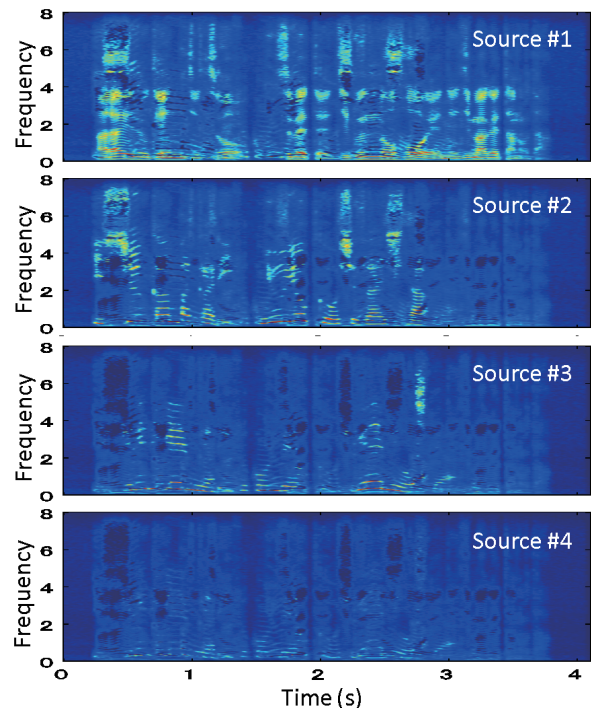


Fig. 6 Separated source signals obtained with present method.

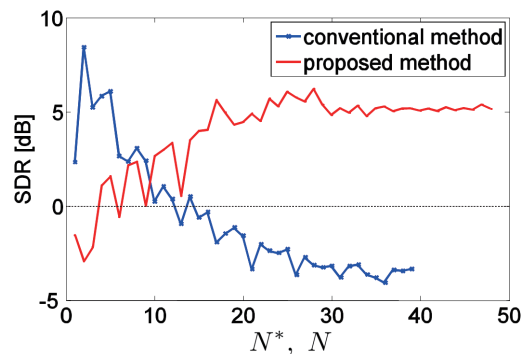


Fig. 7 SDRs obtained with conventional and present methods with different N and N^* settings.

method, such a phenomenon did not occur with the present method. This is because the present method searches for an explanation for an observed signal using a model with as low a complexity as possible.

Figure 8 shows an example of the estimate of the inter-channel phase difference of each source. The slope of $\arg([\mathbf{m}_{i,k}]_2/[\mathbf{m}_{i,k}]_1)$ along the frequency direction k corresponds to the DOA estimate of source i . Hence, if the DOAs of the sources were estimated correctly and the permutation problem was successfully solved, $\arg([\mathbf{m}_{i,k}]_2/[\mathbf{m}_{i,k}]_1)$ with the same index i should lie on a straight line along k . Figure 8 thus reveals that the DOA estimation and permutation alignment both worked successfully with the present method.

It is important to note that the DOA mixture model described in 3.3 assumes a completely anechoic environment. Moreover, the assumption of the sparseness of speech

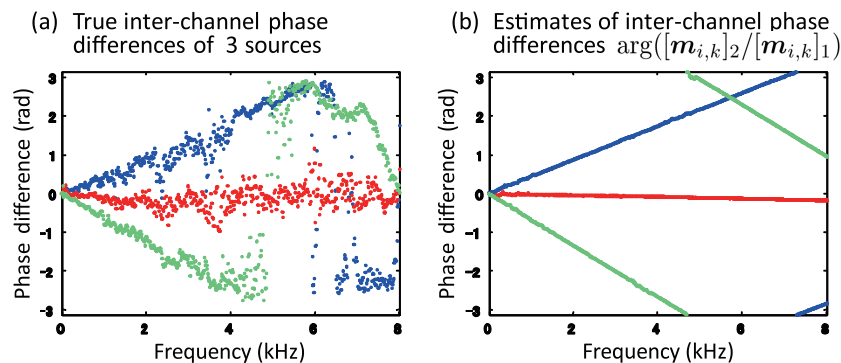


Fig. 8 An example of the estimate of the inter-channel phase difference of each source plotted in a different color. The slope of $\arg([m_{i,k}]_2/[m_{i,k}]_1)$ along the frequency direction k corresponds to the DOA estimate of source i . For each i , $\arg([m_{i,k}]_2/[m_{i,k}]_1)$ lies on a straight line along k . This means that both the DOA estimation and permutation alignment worked successfully with the present method.

described in Sect. 2 does not hold in reverberant environments. Indeed, the performance of the present method was poor when using real data recorded in reverberant environments. This implies the need for a combined generative model of the present model and a reverberation process, e.g., by adopting the idea described in [16].

7. Conclusion

This paper aimed at developing a BSS algorithm that works well even when the number of sources is unknown and proposed a novel BSS approach that simultaneously performs an estimation of the number of sources, source separation based on the sparseness of speech, an estimation of the direction of arrival of each source, and permutation alignment. We confirmed experimentally that reasonably good separations were obtained with the present method without specifying the number of sources.

Blind separation of moving sound sources is an important task to be solved for real applications. One possible extension to the present method involves representing the DOA indicator variable c_i as a time sequence, i.e., $\{c_{i,l}\}_{1 \leq l \leq L}$, and modeling its generative process using a hidden Markov model. This is currently under development.

It should be noted that this paper is an extended journal version of our conference papers [17], [18].

Acknowledgments

We thank Dr. Hiroshi Sawada (NTT) for providing us with the source code of his method. We also thank Mr. Yosuke Izumi (previously with The University of Tokyo) for providing us with Fig. 1.

References

- [1] A. Hyvärinen, J. Karhunen, and E. Oja, *Independent Component Analysis*, John Wiley & Sons, 2001.
- [2] Ö. Yılmaz and S. Rickard, "Blind separation of speech mixtures via time-frequency masking," *IEEE Trans. Signal Process.*, vol.52, no.7, pp.1830–1847, 2004.
- [3] M.I. Mandel, D.P.W. Ellis, and T. Jebara, "An EM algorithm for localizing multiple sound sources in reverberant environments," *Adv. Neural Info. Process. Syst.*, pp.953–960, 2006.
- [4] S. Araki, H. Sawada, R. Mukai, and S. Makino, "Underdetermined blind sparse source separation for arbitrarily arranged multiple sensors," *Signal Process.*, vol.87, no.8, pp.1833–1847, 2007.
- [5] Y. Mori, H. Saruwatari, T. Takatani, S. Ukai, K. Shikano, T. Hiekata, and T. Morita, "Real-time implementation of two-stage blind source separation combining SIMO-ICA and binary masking," *Proc. 9th International Workshop on Acoustic Echo and Noise Control (IWAENC 2005)*, pp.229–232, 2005.
- [6] Y. Izumi, N. Ono, and S. Sagayama, "Sparseness-based 2ch BSS using the EM algorithm in reverberant environment," *Proc. 2007 IEEE Workshop on Applications of Signal Processing to Audio and Acoustics (WASPAA 2007)*, pp.147–150, 2007.
- [7] H. Sawada, S. Araki, and S. Makino, "Underdetermined convolutive blind source separation via frequency bin-wise clustering and permutation alignment," *IEEE Trans. Audio Speech Language Process.*, vol.19, no.3, pp.516–527, 2010.
- [8] J. Sethuraman, "A constructive definition of Dirichlet priors," *Statistica Sinica*, vol.4, pp.639–650, 1994.
- [9] T.S. Ferguson, "A Bayesian analysis of some nonparametric problems," *Annals of Statistics*, vol.1, no.2, pp.209–230, 1973.
- [10] D.M. Blei and M.I. Jordan, "Variational inference for Dirichlet process mixtures," *Journal of Bayesian Analysis*, vol.1, no.1, pp.121–144, 2006.
- [11] J. Taghia, N. Mohammadiha, and A. Leijon, "A variational Bayes approach to the underdetermined blind source separation with automatic determination of the number of sources," *Proc. 2012 IEEE International Conferences on Acoustics, Speech and Signal Processing (ICASSP2012)*, pp.253–256, March 2012.
- [12] T. Otsuka, K. Ishiguro, H. Sawada, and H.G. Okuno, "Bayesian unification of sound source localization and separation with permutation resolution," *Proc. Twenty-Sixth AAAI Conference on Artificial Intelligence (AAAI-12)*, pp.2038–2045, July 2012.
- [13] A. Kurematsu, K. Takeda, Y. Sagisaka, S. Katagiri, H. Kuwabara, and K. Shikano, "ATR Japanese speech database as a tool of speech recognition and synthesis," *Transactions on Speech Communication*, pp.357–363, 1990.
- [14] 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," *Proc. 2nd International Conference on Language Resources & Evaluation (LREC 2000)*, pp.965–968, 2000.
- [15] E. Vincent, R. Gribonval, and C. Févotte, "Performance measurement in blind audio source separation," *IEEE Trans. Audio Speech Language Process.*, pp.1462–1469, 2006.

- [16] T. Yoshioka, T. Nakatani, M. Miyoshi, and H.G. Okuno, "Blind Separation and Dereverberation of Speech Mixtures by Joint Optimization," IEEE Trans. Audio Speech Language Process., vol.19, no.1, pp.69–84, 2011.
- [17] H. Kameoka, M. Sato, T. Ono, N. Ono, and S. Sagayama, "Bayesian nonparametric approach to underdetermined sparse BSS," Proc. ASJ meeting, 1-1-19, pp.713–716, March 2012 (in Japanese).
- [18] H. Kameoka, M. Sato, T. Ono, N. Ono, and S. Sagayama, "Blind separation of infinitely many sparse sources," Proc. 13th International Workshop on Acoustic Signal Enhancement (IWAENC 2012), Sept. 2012.

Appendix A: Probability Density Functions

The circular complex normal distribution with mean μ and covariance P is defined as

$$\mathcal{N}_{\mathbb{C}}(\mathbf{x}; \mu, P) = \frac{1}{|\pi P|} \exp\left((\mathbf{x} - \mu)^H P^{-1} (\mathbf{x} - \mu)\right) \quad (\text{A} \cdot 1)$$

$$= \exp\left[-\mathbf{x}^H P^{-1} \mathbf{x} + 2\text{Re}(\mu^H P^{-1} \mathbf{x}) - \mu^H P^{-1} \mu - \log |\pi P|\right], \quad (\text{A} \cdot 2)$$

and so

$$\mathbb{E}[\mathbf{x}] = \mu, \quad (\text{A} \cdot 3)$$

$$\mathbb{E}[\mathbf{x}\mathbf{x}^H] = P + \mu\mu^H, \quad (\text{A} \cdot 4)$$

where $\mathbb{E}[\cdot]$ denotes expectation.

The Dirichlet distribution of order I with parameters $\alpha_1, \dots, \alpha_I > 0$ is defined as

$$\text{Dirichlet}(\mathbf{x}; \alpha) = \frac{\prod_{i=1}^I \Gamma(\alpha_i)}{\Gamma(\sum_{i=1}^I \alpha_i)} \prod_{i=1}^I x_i^{\alpha_i-1}, \quad (\text{A} \cdot 5)$$

for all $x_1, \dots, x_I > 0$ satisfying $\sum_{i=1}^I x_i = 1$. The density is zero outside this open $(I - 1)$ -dimensional simplex. $\Gamma(\cdot)$ is the gamma function. The mean and entropy of the Dirichlet-distributed variable are given by

$$\mathbb{E}[x_i] = \frac{\alpha_i}{\sum_{i'=1}^I \alpha_{i'}}, \quad (\text{A} \cdot 6)$$

$$\mathbb{E}[\log x_i] = \Psi(\alpha_i) - \Psi\left(\sum_i \alpha_i\right). \quad (\text{A} \cdot 7)$$

Note that when $I = 2$ the Dirichlet distribution reduces to a Beta distribution.

Appendix B: Derivation of Update Equations

Here we show how we obtained the variational update equations. Recall that each variational update equation can be obtained via the formula given by (21).

B.1 Update Equation for $q(A)$

The term in $\log p(Y, \Theta)$ depending on A is given as

$$-\sum_{k,l} (\mathbf{y}_{k,l} - \mathbf{a}_{z_{k,l},k} s_{k,l})^H \Sigma_k^{(n)-1} (\mathbf{y}_{k,l} - \mathbf{a}_{z_{k,l},k} s_{k,l})$$

$$\begin{aligned} & - \sum_{i,k} (\mathbf{a}_{i,k} - \mathbf{h}(\vartheta_{c_i}, \omega_k))^H \Sigma_k^{(a)-1} (\mathbf{a}_{i,k} - \mathbf{h}(\vartheta_{c_i}, \omega_k)) \\ & \stackrel{A}{=} \sum_{i,k} \left[\sum_l \mathbf{1}[z_{k,l} = i] (2\text{Re}[s_{k,l}^* \mathbf{a}_{i,k}^H \Sigma_k^{(n)-1} \mathbf{y}_{k,l}] \right. \\ & \quad \left. - |s_{k,l}|^2 \mathbf{a}_{i,k}^H \Sigma_k^{(n)-1} \mathbf{a}_{i,k}) \right. \\ & \quad \left. + \sum_d \mathbf{1}[c_i = d] (2\text{Re}[\mathbf{a}_{i,k}^H \Sigma_k^{(a)-1} \mathbf{h}(\vartheta_d, \omega_k)] - \mathbf{a}_{i,k}^H \Sigma_k^{(a)-1} \mathbf{a}_{i,k}) \right], \end{aligned} \quad (\text{A} \cdot 8)$$

where $\stackrel{x}{=}$ denotes equality up to a term independent of x and $\mathbf{1}[\cdot]$ denotes the indicator function that takes the value 1 if its argument is true and 0 otherwise. The expectation of (A·8) taken over all the variables except A is given as

$$\begin{aligned} & \sum_{i,k} \left[\sum_l \mathbb{E}[z_{k,l} = i] (2\text{Re}[\mathbb{E}[s_{k,l}^*] \mathbf{a}_{i,k}^H \Sigma_k^{(n)-1} \mathbf{y}_{k,l}] \right. \\ & \quad \left. - \mathbb{E}[|s_{k,l}|^2] \mathbf{a}_{i,k}^H \Sigma_k^{(n)-1} \mathbf{a}_{i,k}) \right. \\ & \quad \left. + \sum_d \mathbb{E}[c_i = d] (2\text{Re}[\mathbf{a}_{i,k}^H \Sigma_k^{(a)-1} \mathbf{h}(\vartheta_d, \omega_k)] - \mathbf{a}_{i,k}^H \Sigma_k^{(a)-1} \mathbf{a}_{i,k}) \right], \end{aligned} \quad (\text{A} \cdot 9)$$

where

$$\mathbb{E}[z_{k,l} = i] = q(z_{k,l} = i) = [\boldsymbol{\phi}_{k,l}]_i \quad (\text{A} \cdot 10)$$

$$\mathbb{E}[c_i = d] = q(c_i = d) = [\boldsymbol{\psi}_i]_d \quad (\text{A} \cdot 11)$$

$$\mathbb{E}[s_{k,l}] = \mu_{k,l} \quad (\text{A} \cdot 12)$$

$$\mathbb{E}[|s_{k,l}|^2] = |\mu_{k,l}|^2 + \sigma_{k,l}^2. \quad (\text{A} \cdot 13)$$

We can thus confirm that (A·9) is equal up to a constant term to a quadratic form $-\sum_{i,k} (\mathbf{a}_{i,k} - \mathbf{m}_{i,k})^H \Gamma_{i,k}^{-1} (\mathbf{a}_{i,k} - \mathbf{m}_{i,k})$ where $\mathbf{m}_{i,k}$ and $\Gamma_{i,k}$ are given by (29) and (28), respectively. Hence, $q(A)$ is shown to be the product of multivariate complex normal distributions with mean $\mathbf{m}_{i,k}$ and covariance $\Gamma_{i,k}$

$$q(A) = \prod_{i,k} \mathcal{N}_{\mathbb{C}}(\mathbf{a}_{i,k}; \mathbf{m}_{i,k}, \Gamma_{i,k}). \quad (\text{A} \cdot 14)$$

B.2 Update Equation for $q(S)$

The term in $\log p(Y, \Theta)$ depending on S is given as

$$\begin{aligned} & - \sum_{k,l} (\mathbf{y}_{k,l} - \mathbf{a}_{z_{k,l},k} s_{k,l})^H \Sigma_k^{(n)-1} (\mathbf{y}_{k,l} - \mathbf{a}_{z_{k,l},k} s_{k,l}) \\ & \stackrel{S}{=} - \sum_{k,l} \sum_i \mathbf{1}[z_{k,l} = i] (|s_{k,l}|^2 \text{tr}[\mathbf{a}_{i,k} \mathbf{a}_{i,k}^H \Sigma_k^{(n)-1}] \\ & \quad - 2\text{Re}[s_{k,l}^* \mathbf{a}_{i,k}^H \Sigma_k^{(n)-1} \mathbf{y}_{k,l}]). \end{aligned} \quad (\text{A} \cdot 15)$$

The expectation of (A·15) taken over all the variables except S is given as

$$\begin{aligned} & - \sum_{k,l} \sum_i \mathbb{E}[z_{k,l} = i] (|s_{k,l}|^2 \text{tr}[\mathbb{E}[\mathbf{a}_{i,k} \mathbf{a}_{i,k}^H] \Sigma_k^{(n)-1}] \\ & \quad - 2\text{Re}[s_{k,l}^* \mathbb{E}[\mathbf{a}_{i,k}^H] \Sigma_k^{(n)-1} \mathbf{y}_{k,l}]), \end{aligned} \quad (\text{A} \cdot 16)$$

where

$$\mathbb{E}[z_{k,l} = i] = [\boldsymbol{\phi}_{k,l}]_i, \quad (\text{A} \cdot 17)$$

$$\mathbb{E}[\mathbf{a}_{i,k}^H] = \mathbf{m}_{i,k}^H, \quad (\text{A} \cdot 18)$$

$$\mathbb{E}[\mathbf{a}_{i,k} \mathbf{a}_{i,k}^H] = \mathbf{m}_{i,k} \mathbf{m}_{i,k}^H + \Gamma_{i,k}. \quad (\text{A} \cdot 19)$$

We can thus confirm that (A·16) is equal up to a constant term to a quadratic function $-\sum_{k,l} |s_{k,l} - \mu_{k,l}|^2 / \sigma_{k,l}^2$ where $\mu_{k,l}$ and $\sigma_{k,l}^2$ are given by (31) and (30), respectively. Hence, $q(S)$ is shown to be the product of complex normal distributions with mean $\mu_{k,l}$ and variance $\sigma_{k,l}^2$

$$q(S) = \prod_{k,l} \mathcal{N}_{\mathbb{C}}(s_{k,l}; \mu_{k,l}, \sigma_{k,l}^2). \quad (\text{A} \cdot 20)$$

B.3 Update Equation for $q(Z)$

The term in $\log p(Y, \Theta)$ depending on Z is given as

$$\begin{aligned} & \sum_{k,l} \left\{ -(\mathbf{y}_{k,l} - \mathbf{a}_{z_{k,l},k} s_{k,l})^H \Sigma_k^{(n)-1} (\mathbf{y}_{k,l} - \mathbf{a}_{z_{k,l},k} s_{k,l}) \right. \\ & \quad \left. + \log v_{z_{k,l}} \prod_{j=1}^{z_{k,l}-1} (1 - v_j) \right\} \\ & \stackrel{=}{=} \sum_{k,l} \sum_i \mathbf{1}[z_{k,l} = i] \left\{ 2\text{Re}[s_{k,l} \mathbf{y}_{k,l}^H \Sigma_k^{(n)-1} \mathbf{a}_{i,k}] \right. \\ & \quad \left. - |s_{k,l}|^2 \text{tr}[\mathbf{a}_{i,k} \mathbf{a}_{i,k}^H \Sigma_k^{(n)-1}] + \log v_i + \sum_{j=1}^{i-1} \log(1 - v_j) \right\}. \end{aligned} \quad (\text{A} \cdot 21)$$

The expectation of (A·21) taken over all the variables except Z is given as

$$\begin{aligned} & \sum_{k,l} \sum_i \mathbf{1}[z_{k,l} = i] \left\{ 2\text{Re}[\mathbb{E}[s_{k,l}] \mathbf{y}_{k,l}^H \Sigma_k^{(n)-1} \mathbb{E}[\mathbf{a}_{i,k}]] \right. \\ & \quad - \mathbb{E}[|s_{k,l}|^2] \text{tr}[\mathbb{E}[\mathbf{a}_{i,k} \mathbf{a}_{i,k}^H] \Sigma_k^{(n)-1}] \\ & \quad \left. + \mathbb{E}[\log v_i] + \sum_{j=1}^{i-1} \mathbb{E}[\log(1 - v_j)] \right\}, \end{aligned} \quad (\text{A} \cdot 22)$$

where

$$\mathbb{E}[\log v_i] = \Psi(\gamma_{i,0}) - \Psi(\gamma_{i,0} + \gamma_{i,1}), \quad (\text{A} \cdot 23)$$

$$\mathbb{E}[\log(1 - v_j)] = \Psi(\gamma_{j,1}) - \Psi(\gamma_{j,0} + \gamma_{j,1}). \quad (\text{A} \cdot 24)$$

Hence,

$$\begin{aligned} q(Z) &= \prod_{k,l} q(z_{k,l}) \\ q(z_{k,l}) &\propto \exp \left\{ 2\text{Re}[\mu_{k,l} \mathbf{y}_{k,l}^H \Sigma_k^{(n)-1} \mathbf{m}_{z_{k,l},k}] \right. \\ & \quad - (|\mu_{k,l}|^2 + \sigma_{k,l}^2) \text{tr}[(\mathbf{m}_{z_{k,l},k} \mathbf{m}_{z_{k,l},k}^H + \Gamma_{z_{k,l},k}) \Sigma_k^{(n)-1}] \\ & \quad + \Psi(\gamma_{z_{k,l},0}) - \Psi(\gamma_{z_{k,l},0} + \gamma_{z_{k,l},1}) \\ & \quad \left. + \sum_{j=1}^{z_{k,l}-1} (\Psi(\gamma_{j,1}) - \Psi(\gamma_{j,0} + \gamma_{j,1})) \right\}. \end{aligned} \quad (\text{A} \cdot 25)$$

B.4 Update Equation for $q(V)$

The term in $\log p(Y, \Theta)$ depending on V is given as

$$\begin{aligned} & \sum_{k,l} \log v_{z_{k,l}} \prod_{j=1}^{z_{k,l}-1} (1 - v_j) + \sum_{i=1}^{\infty} (\alpha_0 - 1) \log(1 - v_i) \\ &= \sum_i \left(\sum_{k,l} \mathbf{1}[z_{k,l} = i] \right) \log v_i \\ & \quad + \sum_i \left(\sum_{j=i+1}^{\infty} \sum_{k,l} \mathbf{1}[z_{k,l} = j] + \alpha_0 - 1 \right) \log(1 - v_i). \end{aligned} \quad (\text{A} \cdot 26)$$

The expectation of (A·26) taken over all the variables except V is given as

$$\begin{aligned} & \sum_i \left\{ \left(\sum_{k,l} \mathbb{E}[z_{k,l} = i] \right) \log v_i \right. \\ & \quad \left. + \left(\sum_{j=i+1}^{N^*} \sum_{k,l} \mathbb{E}[z_{k,l} = j] + \alpha_0 - 1 \right) \log(1 - v_i) \right\}. \end{aligned} \quad (\text{A} \cdot 27)$$

We can thus confirm that the exponential of (A·27) is equal up to a constant factor to

$$q(V) = \prod_i \text{Beta}(v_i; \gamma_{i,0}, \gamma_{i,1}), \quad (\text{A} \cdot 28)$$

where $\gamma_{i,0}$ and $\gamma_{i,1}$ are given by (33) and (34), respectively.

B.5 Update Equation for $q(C)$

The term in $\log p(Y, \Theta)$ depending on C is given as

$$\begin{aligned} & -\sum_{i,k} (\mathbf{a}_{i,k} - \mathbf{h}(\vartheta_{c_i}, \omega_k))^H \Sigma_k^{(a)-1} (\mathbf{a}_{i,k} - \mathbf{h}(\vartheta_{c_i}, \omega_k)) + \sum_i \log \rho_{c_i} \\ & \stackrel{=}{=} \sum_i \left\{ \sum_d \mathbf{1}[c_i = d] \left(\sum_k (2\text{Re}[\mathbf{a}_{i,k}^H \Sigma_k^{(a)-1} \mathbf{h}(\vartheta_d, \omega_k)] \right. \right. \\ & \quad \left. \left. - \mathbf{h}(\vartheta_d, \omega_k)^H \Sigma_k^{(a)-1} \mathbf{h}(\vartheta_d, \omega_k)) + \log \rho_d \right) \right\}. \end{aligned} \quad (\text{A} \cdot 29)$$

The expectation of (A·29) taken over all the variables except C is given as

$$\begin{aligned} & \sum_i \left\{ \sum_d \mathbf{1}[c_i = d] \left(\sum_k (2\text{Re}[\mathbf{m}_{i,k}^H \Sigma_k^{(a)-1} \mathbf{h}(\vartheta_d, \omega_k)] \right. \right. \\ & \quad \left. \left. - \mathbf{h}(\vartheta_d, \omega_k)^H \Sigma_k^{(a)-1} \mathbf{h}(\vartheta_d, \omega_k)) + \mathbb{E}[\log \rho_d] \right) \right\}, \end{aligned} \quad (\text{A} \cdot 30)$$

where

$$\mathbb{E}[\log \rho_d] = \Psi(\zeta_d) - \Psi\left(\sum_{d'} \zeta_{d'}\right). \quad (\text{A} \cdot 31)$$

We can thus confirm that the exponential of (A·30) is equal up to a constant factor to

$$q(C) = \prod_i \text{Categorical}(c_i; \boldsymbol{\psi}_i), \quad (\text{A} \cdot 32)$$

where $[\psi_i]_d$ is given by (35).

B.6 Update Equation for $q(\rho)$

The term in $\log p(Y, \Theta)$ depending on ρ is given as

$$\begin{aligned} & \sum_i \log \rho_{c_i} + (\beta_0/D - 1) \sum_d \log \rho_d \\ &= \sum_i \sum_d \mathbf{1}[c_i = d] \log \rho_d + (\beta_0/D - 1) \sum_d \log \rho_d \end{aligned} \quad (\text{A} \cdot 33)$$

The expectation of (A·33) taken over all the variables except ρ is given as

$$\sum_d \left(\sum_i [\psi_i]_d + \beta_0/D - 1 \right) \log \rho_d. \quad (\text{A} \cdot 34)$$

We can thus confirm that the exponential of (A·34) is equal up to a constant factor to

$$q(\rho) = \text{Dirichlet}(\rho; \zeta_1, \dots, \zeta_D), \quad (\text{A} \cdot 35)$$

where ζ_d is given by (36).



Hirokazu Kameoka received B.E., M.S. and Ph.D. degrees all from the University of Tokyo, Japan, in 2002, 2004 and 2007, respectively. He is currently a research scientist at NTT Communication Science Laboratories and a Visiting Associate Professor at the University of Tokyo. His research interests include computational auditory scene analysis, statistical signal processing, speech and music processing, and machine learning. He is a member of IEEE, the Information Processing Society of Japan (IPSJ)

and the Acoustical Society of Japan (ASJ). He received 13 awards over the past 10 years, including the Yamashita Memorial Research Award in 2005 from IPSJ, the Itakura Prize Innovative Young Researcher Award in 2007 and the Awaya Prize Young Researcher Award in 2008 from ASJ, IEEE Signal Processing Society 2008 SPS Young Author Best Paper Award in 2009, and IEICE ISS Young Researcher's Award in Speech Field in 2011.



Misa Sato received her B.E. degree from the University of Tokyo, Japan, in 2012. She is currently an MSc student at the Graduate School of Engineering, the University of Tokyo, Japan. Her research interests include natural language processing and machine learning.



Takuma Ono received the B.E. and M.S. degrees from the University of Tokyo, Japan, in 2010 and 2012, respectively. Since then, he has been with Zenkyoren (National Mutual Insurance Federation of Agricultural Cooperatives).



Nobutaka Ono received the B.E., M.S., and Ph.D. degrees in Mathematical Engineering and Information Physics from the University of Tokyo, Japan, in 1996, 1998, 2001, respectively. He joined the Graduate School of Information Science and Technology, the University of Tokyo, Japan, in 2001 as a Research Associate and became a Lecturer in 2005. He moved to the National Institute of Informatics, Japan, as an Associate Professor in 2011. His research interests include acoustic signal processing, specifically,

microphone array processing, source localization and separation, music signal processing, audio coding and watermarking, and optimization algorithms for them. He is the author or co-author of more than 90 articles in international journals and conference proceedings. He was a Tutorial speaker at ISMIR2010 and will organize a special session in EU-SIPCO2013. Dr. Ono has been an Associate Editor of the IEEE Transactions on Audio, Speech and Language Processing since 2012 and an Associate Editor of Acoustic Science and Technology since 2012. He is a chair of SiSEC (Signal Separation Evaluation Campaign) committee in 2013. He is a senior member of the IEEE Signal Processing Society and a member of the Acoustical Society of Japan (ASJ), the Information Processing Society of Japan (IPSJ), the Institute of Electrical Engineers of Japan (IEEJ), and the Society of Instrument and Control Engineers (SICE). He received the Sato Paper Award and the Awaya Award from ASJ in 2000 and 2007, respectively, and received the Igarashi Award at the Sensor Symposium on Sensors, Micromachines, and Applied Systems from IEEJ in 2004.



Shigeki Sagayama received the B.E., M.S., and Ph.D. degrees from the University of Tokyo, Tokyo, Japan, in 1972, 1974, and 1998, respectively, all in mathematical engineering and information physics. He joined Nippon Telegraph and Telephone Public Corporation (currently, NTT) in 1974 and started his career in speech analysis, synthesis, and recognition at NTT Labs in Musashino, Japan. From 1990, he was Head of the Speech Processing Department, ATR Interpreting Telephony Laboratories,

Kyoto, Japan where he was in charge of an automatic speech translation project. In 1993, he was responsible for speech recognition, synthesis, and dialog systems at NTT Human Interface Laboratories, Yokosuka, Japan. In 1998, he became a Professor of the Graduate School of Information Science, Japan Advanced Institute of Science and Technology (JAIST), Ishikawa. In 2000, he was appointed Professor at the Graduate School of Information Science and Technology (formerly, Graduate School of Engineering), the University of Tokyo. On his retirement from the University of Tokyo in 2013, he became a Project Professor at the National Institute of Informatics (NII). His major research interests include the processing and recognition of speech, music, acoustic signals, handwriting, and images. He was the leader of anthropomorphic spoken dialog agent project (Galatea Project) from 2000 to 2003. Prof. Sagayama received the National Invention Award from the Institute of Invention of Japan in 1991, the Chief Official's Award for Research Achievement from the Science and Technology Agency of Japan in 1996, and other academic awards including Paper Awards from the Institute of Electronics, Information and Communications Engineers, Japan (IEICEJ) in 1996 and from the Information Processing Society of Japan (IPSJ) in 1995. He is a member of the Acoustical Society of Japan, IEICEJ, and IPSJ.