

An Approach to Solving a Permutation Problem of Frequency Domain Independent Component Analysis for Blind Source Separation of Speech Signals

Masaru Fujieda, Takahiro Murakami, and Yoshihisa Ishida

Abstract—Independent component analysis (ICA) in the frequency domain is used for solving the problem of blind source separation (BSS). However, this method has some problems. For example, a general ICA algorithm cannot determine the permutation of signals which is important in the frequency domain ICA. In this paper, we propose an approach to the solution for a permutation problem. The idea is to effectively combine two conventional approaches. This approach improves the signal separation performance by exploiting features of the conventional approaches. We show the simulation results using artificial data.

Keywords—Blind source separation, Independent component analysis, Frequency domain, Permutation ambiguity.

I. INTRODUCTION

BLIND source separation (BSS) [1], [4] is an approach to estimating original source signals by using only the information of the mixed signals observed at a sensor array. This technique is applicable to the realization of the noise-robust speech recognition, the hearing aid which can enhance the specific sound that a user wants to listen, and so on. Independent component analysis (ICA) [4] is one of the statistical analysis methods and identifies the independent components in the random variables. ICA can be used for BSS of linear (instantaneous) mixtures. To achieve BSS of convolutive mixtures, several methods have been proposed [1], [2], [3], [5]. Although there are a number of applications for BSS to mixed speech signals in realistic acoustical environments, the separation performance is still not good enough [6].

In this paper, we propose an approach to solving a permutation problem of frequency domain ICA (FDICA) for BSS. In FDICA, it is necessary to solve the problem of ambiguity of scaling factors and permutation. Several methods have already been proposed for such problems [1], [2] but each of them has imperfection. For example, the method [1] is impossible to estimate the permutation when the envelopes of

source signals are similar. On the other hand, the method [2] is not dependent on the envelope of source signals. The method [1] performs better than the method [2] when the envelopes of source signals are different from each other. To improve the performance, we propose a new method for combining the features of these methods.

The following chapter explains the formulation of the general ICA problems and FDICA. In Section III, the algorithm to solve scaling and permutation problems is explained. In Section IV, the source separation experiments are performed and the simulation results are illustrated. Finally, Section V shows the conclusion of this paper.

II. INDEPENDENT COMPONENT ANALYSIS

A. Statement of the Problem

In general, the linear mixture model of ICA is given by the following equation:

$$\mathbf{x}(t) = \mathbf{A} \cdot \mathbf{s}(t), \quad (1)$$

where $\mathbf{x}(t) = [x_1(t), x_2(t), \dots, x_n(t)]^T$ is a set of observed signals, \mathbf{A} is an unknown mixing matrix, and $\mathbf{s}(t) = [s_1(t), s_2(t), \dots, s_m(t)]^T$ is a set of source signals which is assumed to be mutually independent. The purpose of ICA is to find a separation matrix \mathbf{W} so that the output signals

$$\mathbf{y}(t) = \mathbf{W} \cdot \mathbf{x}(t), \quad (2)$$

where $\mathbf{y}(t) = [y_1(t), y_2(t), \dots, y_m(t)]^T$ is mutually independent. However, there are still ambiguity of scaling factors and permutation of output signals because of lack of information about the amplitude and permutation of the source signals, that is, \mathbf{W} is allowed to satisfy the following equation:

$$\begin{aligned} \mathbf{y}(t) &= \mathbf{W} \cdot \mathbf{x}(t) \\ &= \mathbf{Q} \cdot \mathbf{D} \cdot \mathbf{s}(t), \end{aligned} \quad (3)$$

Manuscript received September 27, 2006.

The authors are with the School of Science and Technology, Meiji University, Kanagawa, Japan (corresponding author to provide phone: +81-44-934-7307; e-mail: ce66050@isc.meiji.ac.jp).

where \mathbf{D} is a diagonal matrix which represents scaling factors and \mathbf{Q} is a permutation matrix.

In the case of acoustical signals, the mixing model is expressed by the convolutive mixture as

$$\begin{aligned} \mathbf{x}(t) &= \sum_{\tau} \mathbf{A}(\tau) \cdot \mathbf{s}(t - \tau) \\ &= \mathbf{A}(\tau) * \mathbf{s}(t), \end{aligned} \quad (4)$$

where operator $*$ denotes convolution and $\mathbf{A}(\tau)$ is an unknown mixing filter. In the convolutive mixture case, (2) is rewritten by using a separation filter $\mathbf{W}(\tau)$ as

$$\begin{aligned} \mathbf{y}(t) &= \mathbf{W}(\tau) * \mathbf{x}(t) \\ &= \mathbf{Q} \cdot \mathbf{D} \cdot \mathbf{s}(t), \end{aligned} \quad (5)$$

B. Frequency Domain Independent Component Analysis (FDICA)

By using the short time Fourier transform (STFT), we can transpose the convolutive model to the instantaneous model with complex values as in the form

$$\hat{\mathbf{x}}(\omega, t_s) = \hat{\mathbf{A}}(\omega) \cdot \hat{\mathbf{s}}(\omega, t_s), \quad (6)$$

where $\hat{\mathbf{x}}(\omega, t_s)$ and $\hat{\mathbf{s}}(\omega, t_s)$ denote STFT of $\mathbf{x}(t)$ and $\mathbf{s}(t)$, respectively and $\hat{\mathbf{A}}(\omega)$ denotes the discrete Fourier transform (DFT) of the mixing filter $\mathbf{A}(\tau)$. $\hat{\mathbf{x}}(\omega, t_s)$ and $\hat{\mathbf{s}}(\omega, t_s)$ are often referred as spectrograms. For any fixed frequency ω in (6), $\hat{\mathbf{x}}(\omega, t_s)$ can be regarded as the linear mixture model. Therefore, we can apply ICA for the linear mixture model at each frequency independently.

In FDICA, ambiguity of scaling factors and the permutation of output signals become problematic when time domain signals are reconstructed.

We propose an approach to solving the permutation problem. For the scaling problem, the method presented in [1] is applied.

III. ALGORITHM

A. Solving Scaling Problem

When a separating matrix $\hat{\mathbf{W}}(\omega)$ is obtained from the spectrogram $\hat{\mathbf{x}}(\omega, t_s)$, the corresponding independent components are given by:

$$\hat{\mathbf{u}}(\omega, t_s) = \hat{\mathbf{W}}(\omega) \hat{\mathbf{x}}(\omega, t_s). \quad (7)$$

Then we can obtain the scaling factor as

$$\hat{\mathbf{v}}(\omega, t_s; i) = \hat{\mathbf{W}}^{-1}(\omega) \begin{pmatrix} 0 \\ \vdots \\ \hat{u}_i(\omega, t_s) \\ \vdots \\ 0 \end{pmatrix}, \quad (8)$$

where $\hat{u}_i(\omega, t_s)$ denotes the i -th independent component of $\hat{\mathbf{u}}(\omega, t_s)$ and $\hat{\mathbf{v}}(\omega, t_s; i)$ denotes the i -th independent components observed at each sensor [1].

B. Solving Permutation Problem

The proposed method is to combine the features of two conventional techniques presented in [1], [2]. In this section, we review these methods briefly.

1) Method based on temporal structure of speech signals (MTS)

This algorithm was proposed by Murata et al. [1], utilizes a property of speech signals, that is, signals are intrinsically non-stationary in a long range mainly because of amplitude modulation.

Let us rewrite $\hat{s}_i(\omega, t_s)$, i -th element of $\hat{\mathbf{s}}(\omega, t_s)$ in (6):

$$\hat{s}_i(\omega, t_s) = |\hat{s}_i(\omega, t_s)| e^{j\phi_i(\omega, t_s)}, \quad (9)$$

where $|\hat{s}_i(\omega, t_s)|$ and $\phi_i(\omega, t_s)$ are an absolute value, a phase of source signals in time-frequency domain, respectively. Since source signals are mutually independent, they satisfy the following equations:

$$\text{corr}(|\hat{s}_i(\omega, t_s)|, |\hat{s}_j(\omega, t_s)|) = 0, \quad i \neq j, \quad (10)$$

$$\text{corr}(|\hat{s}_i(\omega, t_s)|, |\hat{s}_j(\omega', t_s)|) = 0, \quad i \neq j, \quad \omega \neq \omega', \quad (11)$$

where

$$\begin{aligned} \text{corr}(|\hat{s}_i(\omega, t_s)|, |\hat{s}_j(\omega, t_s)|) \\ = \left(\overline{|\hat{s}_i(\omega, t_s)| \cdot |\hat{s}_j(\omega, t_s)|} \right) + \overline{|\hat{s}_i(\omega, t_s)|} \cdot \overline{|\hat{s}_j(\omega, t_s)|}, \end{aligned} \quad (12)$$

$$\overline{|\hat{s}_i(\omega, t_s)|} = \frac{1}{T} \sum_{s=1}^T |\hat{s}_i(\omega, t_s)|, \quad (13)$$

We assume that T is sufficiently large. On the other hand, for different frequency components corresponding to the same source signal, we can assume

$$\text{corr}(|\hat{s}_i(\omega, t_s)|, |\hat{s}_i(\omega', t_s)|) \neq 0, \quad \omega \neq \omega'. \quad (14)$$

Therefore, the correlation coefficient of their envelopes shown as

$$= \frac{r(|\hat{s}_i(\omega, t_s)|, |\hat{s}_j(\omega', t_s)|)}{\sqrt{\text{corr}(|\hat{s}_i(\omega, t_s)|, |\hat{s}_i(\omega, t_s)|) \text{corr}(|\hat{s}_j(\omega', t_s)|, |\hat{s}_j(\omega', t_s)|)}}, \quad (15)$$

may be available for estimating an appropriate combination of frequency elements.

Then, let us define a moving average operator ε for estimating the envelope of time series as

$$\varepsilon \hat{v}(\omega, t_s; i) = \frac{1}{2M+1} \sum_{t'_s=t_s-M}^{t_s+M} \sum_{j=1}^n |\hat{v}_j(\omega, t_s; i)|, \quad (16)$$

where M is a positive constant and $\hat{v}_j(\omega, t_s; i)$ denotes the j -th element of $\hat{v}(\omega, t_s; i)$.

We solve the permutation problem by the sorting based on the correlation of envelopes as follows:

- 1) Sort ω in order of independent components with the low correlation. This is done by sorting them in ascending order of similarity defined as

$$\text{sim}(\omega) = \sum_{i \neq j} r(\varepsilon \hat{v}(\omega, t_s; i), \varepsilon \hat{v}(\omega, t_s; j)), \quad (17)$$

$$\text{sim}(\omega_1) \leq \text{sim}(\omega_2) \leq \dots \leq \text{sim}(\omega_N). \quad (18)$$

- 2) For ω_1 , assign $\hat{v}(\omega_1, t_s; i)$ to $\hat{y}(\omega_1, t_s; i)$ as follows:

$$\hat{y}(\omega_1, t_s; i) = \hat{v}(\omega_1, t_s; i), \quad i = 1, \dots, n. \quad (19)$$

- 3) For ω_k , find a permutation $\sigma(i)$ which maximizes the correlation between the envelope of ω_k and the summed envelope from ω_1 through ω_{k-1} . This is achieved by maximizing the sum of correlation coefficients

$$\sum_{i=1}^m r\left(\varepsilon \hat{v}(\omega_k, t_s; \sigma(i)), \sum_{j=1}^{k-1} \varepsilon \hat{y}(\omega_j, t_s; i)\right), \quad (20)$$

for all the possible permutations of $i = 1, \dots, m$.

- 4) Assign the appropriate permutation to $\hat{y}(\omega_k, t_s; i)$:

$$\hat{y}(\omega_k, t_s; i) = \hat{v}(\omega_k, t_s; \sigma(i)), \quad i = 1, \dots, m. \quad (21)$$

- 5) Go to 3) until $k = N$.

As a result, we obtain separated spectrograms

$$\hat{y}(\omega_k, t_s; i), \quad i = 1, \dots, m. \quad (22)$$

- 2) Method based on relation of the mixing matrix at the adjacent frequencies (MRM)

Alternative approach, which was proposed by Asano et al. [2], utilizes a relation of the location vectors in the mixing matrix at adjacent frequencies. The relation between the mixing matrices at frequencies ω_k and ω_{k-1} are written as

$$\hat{A}(\omega_k) = T(\omega_k, \omega_{k-1}) \cdot \hat{A}(\omega_{k-1}) \quad (23)$$

where $T(\omega_k, \omega_{k-1})$ is a rotation matrix. Absolute values of mixing matrices $\hat{A}(\omega_k)$ and $\hat{A}(\omega_{k-1})$ are assumed to be unity for the sake of simplicity. If the difference between two frequencies $\Delta\omega = \omega_k - \omega_{k-1}$ is sufficiently small, the relation can be written as

$$\hat{A}(\omega_k) \approx \hat{A}(\omega_{k-1}), \quad T(\omega_k, \omega_{k-1}) \approx \mathbf{I}, \quad (24)$$

where \mathbf{I} is a unit matrix. Then, an angle between the location vectors is small. The location vector $\hat{a}_i(\omega_k)$ is obtained by the column vector in the mixing matrix $\hat{A}(\omega)$ as follows:

$$\hat{A}(\omega) = [\hat{a}_1(\omega) \quad \hat{a}_2(\omega) \quad \dots \quad \hat{a}_m(\omega)]. \quad (25)$$

$\hat{a}_i(\omega)$ is defined as

$$\hat{a}_i(\omega) = \begin{bmatrix} e^{-j\omega\tau_{i1}} \\ \vdots \\ e^{-j\omega\tau_{im}} \end{bmatrix}, \quad (26)$$

where $\tau_{j,i}$ denotes the propagation time from the i -th source signal to the j -th sensor. A symbol θ_i is defined as the angle between $\hat{a}_i(\omega_k)$ and $\hat{a}_i(\omega_{k-1})$. Then, θ_i is expected to be the smallest for the correct permutation.

Based on the above discussion, the permutation problem is solved by minimizing the sum of the angles $\{\theta_1, \dots, \theta_m\}$ between the location vectors in the adjacent frequencies. An estimate of the mixing matrix can be obtained by the inverse of the separation matrix $\hat{W}(\omega)$ as

$$\tilde{A}(\omega) = \hat{W}^{-1}(\omega). \quad (27)$$

Let us define the mixing matrix multiplied by the arbitrary permutation matrix \mathbf{P} as

$$\bar{A}^T(\omega) = \mathbf{P} \cdot \tilde{A}^T(\omega). \quad (28)$$

The permutation $\mathbf{P}\tilde{A}^T(\omega)$ exchanges the row vectors of

$\tilde{A}^T(\omega)$. The column vectors of $\tilde{A}(\omega)$ are denoted as $\tilde{A}(\omega) = [\tilde{a}_1(\omega), \dots, \tilde{a}_m(\omega)]$. The cosine of the angle θ_i between the two vectors, $\tilde{a}_i(\omega_k)$ and $\tilde{a}_i(\omega_{k-1})$, is defined as

$$\cos \theta_i = \frac{\tilde{a}_i^H(\omega_k) \tilde{a}_i(\omega_{k-1})}{\|\tilde{a}_i(\omega_k)\| \cdot \|\tilde{a}_i^H(\omega_{k-1})\|}. \quad (29)$$

By using this, the permutation matrix is determined as

$$\hat{P} = \arg \max_P F(P), \quad (30)$$

where the cost function $F(P)$ is defined as

$$F(P) = \frac{1}{m} \sum_{i=1}^m \cos \theta_i. \quad (31)$$

The above method assumes that the estimate of the mixing matrix $\tilde{A}(\omega)$ is a good approximation of the true mixing matrix $A(\omega)$. Therefore, if $\tilde{A}(\omega_{k-1})$ is a bad approximation, we may fail to estimate the correct permutation at the frequency ω_k .

To prevent this, the reference frequency is extended to the following frequency range:

$$\omega_{k-l} = \omega_k - l \cdot \Delta\omega_k, \text{ for } l = 1, \dots, L. \quad (32)$$

The cost function (31) is calculated at all L frequencies within this range. Let us define the value of the cost function at $\omega_{k-l} = \omega_k - l \cdot \Delta\omega_k$ as $F(P, l)$. Next, a confidence measure for $F(P, l)$ is considered. When the largest value of $\max F(P, l)$ closes to other cost functions $F(P, l)$, it may be difficult to determine the correct permutation, and then the value of $F(P, l)$ is not reliable. Based on this assumption, the following confidence measure is defined:

$$C(l) = \max_{P \in \Omega} [F(P, l)] - \max_{P \in \Omega'} [F(P, l)] \quad (33)$$

where Ω denotes a set of all possible P and Ω' denotes Ω with $\hat{P} = \arg \max_{P \in \Omega} [F(P, l)]$. An appropriate reference frequency is determined as

$$\hat{l} = \max_l C(l). \quad (34)$$

The permutation is then solved by using the information at this reference frequency as

$$\hat{P} = \arg \max_P F(P, \hat{l}). \quad (35)$$

3) Proposed method which combines two conventional methods (MCC)

In this paper, we propose a method which combines MTS and MRM. MTS has a disadvantage, that is, it is impossible to estimate the permutation when the envelopes of source signals are similar. On the other hand, MRM is not dependent on the envelope of source signals. However, MTS provides better performance than MRM when the envelopes of source signals are different each other. Considering these facts, we propose to combine two methods in order to cover each disadvantage.

We consider two situations. In the case that the envelopes of separated signals are similar, MTS is employed. In the case that the envelopes of separated signals are different, MRM is utilized. This gives better performance than using only MTS or MRM.

The proposed procedure is summarized as follows:

- 1) Set an appropriate threshold α ($|\alpha| \leq 1$).
- 2) For ω_1 , assign $\hat{v}(\omega_1, t_s; i)$ to $\hat{y}(\omega_1, t_s; i)$ as

$$\hat{y}(\omega_1, t_s; i) = \hat{v}(\omega_1, t_s; i), \quad i = 1, \dots, n. \quad (36)$$

- 3) For ω_k , use MTS methods (3) and (4).
- 4) $k = k + 1$, and go to 3) while $\text{sim}(\omega_k) < \alpha$.
- 5) In the ascending frequency order, apply MRM within the reference frequency range $[-L, L]$.

IV. EXPERIMENTAL RESULTS

In this section, we compare the performance of the conventional method with the proposed method.

A. Conditions for Experiments

For the source signals, three Japanese speech sounds recorded separately are used.

Intonation of Japanese is given by the change of the fundamental frequency. Thus, the envelopes of the Japanese speech can be often similar to each other. The sets are prepared as follows:

- Set 1: Envelopes are mutually different,
- Set 2: Envelopes are mutually similar.

The waveforms are shown in Fig. 1.

The mixing process is assumed to be

$$\begin{aligned} x_1(t) &= 0.4249s_1(t) + 0.3902s_2(t - \tau) \\ x_2(t) &= 0.3322s_1(t - \tau) + 0.3401s_2(t), \end{aligned} \quad (37)$$

where $s_1(t)$ is "aoiie," which means a blue house in English. $s_2(t)$ is "sakuragasaita," that means Japanese cherries blossomed in English, or "aiueo," five vowels in Japanese. τ is a delay factor. The other parameters are shown in Table I.

B. Evaluation

Since the true sources and mixing process are available, we can evaluate the performance using the noise-reduction rate (NRR) [3] that is defined by output signal-to-noise ratio (SNR) minus input SNR as

$$\begin{aligned} \text{SNR}_{out} &= 10 \log_{10} \left(\frac{\text{var}[s(t)]}{\text{var}[y(t) - s(t)]} \right) \\ \text{SNR}_{in} &= 10 \log_{10} \left(\frac{\text{var}[s(t)]}{\text{var}[x(t) - s(t)]} \right), \\ \text{NRR} &= \text{SNR}_{out} - \text{SNR}_{in} \\ &= 10 \log_{10} \left(\frac{\text{var}[x(t) - s(t)]}{\text{var}[y(t) - s(t)]} \right) \end{aligned} \quad (38)$$

where $\text{var}[\bullet]$ denotes the variance of the signal. In practice, we use the following equation:

$$\text{NRR} = \frac{1}{nm} \sum_{i=1}^n \sum_{j=1}^m 10 \log_{10} \left(\frac{\text{var}[x_i(t) - s_j(t)]}{\text{var}[y_{ij}(t) - s_j(t)]} \right), \quad (39)$$

$i = 1, \dots, n, \quad j = 1, \dots, m$

where $n = m = 2$. Table II shows the experimental results.

The performance of MCC was the best. When the envelope of source signals is mutually different, MTS performs better. However, in the frequency bin where the separation by ICA does not succeed, the estimate of permutation is incorrect because the envelope of spectrogram is mutually similar. In this case, the reliability of MRM is higher than that of MTS. On the other hand, when the envelope of source signals is mutually similar, MRM performs better. In the case that $\text{sim}(\omega_k)$ in (17) is sufficiently large, the reliability of MTS is higher than that of MRM. MCC can select the most appropriate approach for every frequency bin by setting the optimum threshold. As the result, MCC shows the best performance.

V. CONCLUSION

In this paper, a method to combine two conventional methods for solving the permutation problem has been proposed. Our approach is to overcome the drawbacks in the conventional algorithms and the effect of a new method is well-recognized in the experiments. The proposed method works well and gives better performance than the conventional methods.

REFERENCES

- [1] N. Murata, S. Ikeda, A. Ziehe, "An approach to blind source separation based on temporal structure of speech signals," in *Neurocomputing*, 41, 2001, pp. 1–24.
- [2] F. Asano, S. Ikeda, M. Ogawa, H. Asoh, N. Kitawaki, "Combined approach of array processing and independent component analysis for blind separation of acoustic signals," in *IEEE Trans. Speech and Audio Processing*, 11, No. 3, 2003, pp. 204–214.
- [3] S. Kurita, H. Saruwatari, S. Kajita, K. Takeda, F. Itakura, "Evaluation of blind signal separation method using directivity pattern under reverberant conditions," in *Proceedings of IEEE International Conference on Acoustics, Speech, and Signal Processing*, SAM-P2-5, Jun. 2000, pp. 3140–3143.
- [4] A. Hyvärinen, E. Oja, "Independent component analysis: a tutorial," in *Helsinki University of Technology, Laboratory of computer and Information Science*, Apr. 1999.
- [5] N. Murata, S. Ikeda, "An on-line algorithm for blind source separation on speech signals," in *Proceedings NOLTA'98*, Sep. 1998, pp. 923–926.
- [6] S. Araki, R. Mukai, S. Makino, T. Nishikawa, H. Saruwatari, "The fundamental limitation of frequency domain blind source separation for convolutive mixtures of speech," in *IEEE Trans. Speech and Audio Processing*, 11, No. 2, 2003, pp. 109–116.

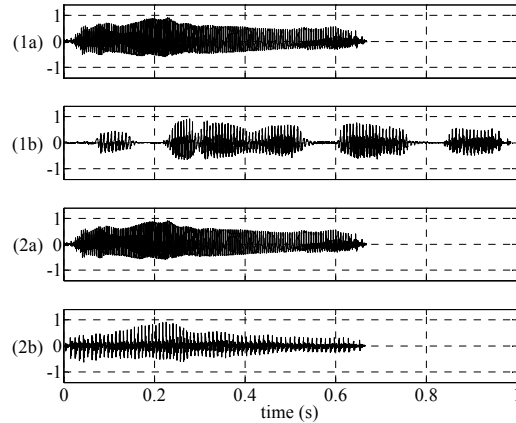


Fig. 1 Waveforms of source signals. (1a) is "aoiie," which means a blue house in English. (1b) is "sakuragasaita," that means Japanese cherries blossomed in English. (2a) is the same (1a). (2b) is "aiueo," five vowels in Japanese. (1a) and (1b) are a set of different envelopes. (2a) and (2b) are a set of similar envelopes

TABLE I
PARAMETERS OF EXPERIMENT

Parameter Name	Value
Sampling frequency	16 kHz
Delay τ	0.3125 ms (5 points)
Window length	8 ms (128 points)
Window function	Hamming window
Shifting time of analysis frames	1.25 ms (20 points)
Threshold α (for Set 1)	0.70
Threshold α (for Set 2)	-0.06
Window length of moving average M (16)	15
Number of reference frequency range L	3

TABLE II
NOISE-REDUCTION RATES FOR CONVOLUTIVE MIXTURE

Method	Envelope is mutually different	Envelope is mutually similar
MTS	18.22 dB	0.93 dB
MRM	12.86 dB	7.07 dB
MCC	18.74 dB	8.84 dB