

Real-time Moving Blind Source Extraction based on Constant Separating Vector and Auxiliary Function Technique

Sihan Yuan, Tetsuya Ueda, Shoji Makino

Waseda University

2-7 Hibikino, Wakamatsu-ku, Kitakyushu, Fukuoka 808-0135, Japan

E-mail: {yuansihan@fuji, t.ueda@akane, s.makino@}.waseda.jp

Abstract

In this paper, we propose a novel online algorithm for moving blind source extraction (BSE). The original algorithm is an offline algorithm based on the recently proposed constant separating vector (CSV) mixing model with auxiliary-function-based independent vector extraction (AuxIVE). The offline CSV-AuxIVE is not suitable for some devices that require real-time processing. In this case, we propose the online CSV-AuxIVE which only needs to know a part of the mixed signal and then sequentially processes the observed signal. Then, we verified the effectiveness of the proposed method on a specific task of the extraction of moving two-speaker signals. The experiment compares the online CSV-AuxIVE with the online AuxIVA. The result shows that under the same conditions of the source of interest (SOI), the average SDR of online CSV-AuxIVE is approximately 1.5 dB higher than that of online AuxIVA, regardless of the change in the range and speed of the interference (IR).

1. Introduction

As technology continues to advance, there are more and more speech devices used in a variety of practical scenarios, such as hearing aids, speech recognition modules in smart homes, and speech enhancement modules in online conferencing. All of these applications need to extract a specific source signal from the mixed signals recorded by the sensors, provided that the information about source signals is unknown. This technique is blind source extraction (BSE), or generally, is blind source separation (BSS). It has become one of the research hotspots in signal processing.

Many effective solutions have been proposed by many researchers. Among them, the most classical method is independent component analysis (ICA) [1, 2], which assumes that the source signal is statistically independent and source separation is accomplished by maximizing non-Gaussianity. Alternatively, among the frequency domain-based methods, the most commonly used is independent vector analysis (IVA) [3, 4], as an extension of independent component analysis (ICA), which uses a joint statistical source model to treat all frequencies simultaneously. For faster convergence, inde-

pendent vector analysis based on auxiliary functions (AuxIVA) [5, 6] has been proposed, which eliminates the effect of tuning parameters on the convergence speed by introducing auxiliary functions that guarantee the monotonic decrease of the objective function. However, the aforementioned methods perform better with static systems (meaning all sources are static). In real situation, the speakers will move back and forth. Thus, in such dynamic situations, these methods can be employed in an adaptive way within small intervals during which the mixture is approximately stationary. In [7], researchers extend the batch-processing AuxIVA to an online process for the autoregressive approximation of auxiliary variables. Although this online method can adapt the speaker's movement, the accuracy is not always so good because the movement range and the speed of the target source change. Especially, the performance will degrade when the source moves so quickly.

In recent years, the constant separating vector (CSV) mixing model [8] has been proposed for moving BSE. It involves both time-varying mixing parameters and time-invariant separation parameters, which effectively maintain the extraction performance regardless of the speaker's movement. Meanwhile, This model combined with the auxiliary function-based independent vector extraction (AuxIVE) to accelerate the convergence. However, CSV-AuxIVE [9] is an offline algorithm, which is not suitable for some devices that require real-time processing, such as hearing aids and intelligent speech assistants. In this paper, we propose an online version of the CSV-AuxIVE model. Unlike the offline algorithm, the proposed algorithm only needs to know part of the mixed speech before the beginning of extraction, and then process the input observed signals in a sequential manner. Finally, we evaluated the performance of the proposed algorithm in the experiment on dynamic two-speaker signals, comparing with online AuxIVA [7].

2. Conventional method: CSV-AuxIVE

2.1 Blind Source Extraction

Let us first explain the BSS scheme assuming that there are N sources recorded by M microphones. After the short-

time Fourier transformer (STFT), let $s_{k,l}^n$ and $x_{k,l}^m$ denote the STFT coefficients of the n -th source and the m -th microphone, where $k = 1, \dots, K$ and $l = 1, \dots, L$ represent the frequency bin and time frame, respectively. In the frequency-domain instantaneous mixture model, the relationship between the observed signal $\mathbf{x}_{k,l} = [x_{k,l}^1, \dots, x_{k,l}^M]^\top \in \mathbb{C}^M$ and the source signal $\mathbf{s}_{k,l} = [s_{k,l}^1, \dots, s_{k,l}^N]^\top \in \mathbb{C}^N$ can be written as:

$$\mathbf{x}_{k,l} = \mathbf{A}_{k,l} \mathbf{s}_{k,l}, \quad (1)$$

where $\mathbf{A}_{k,l}$ stands for $M \times N$ mixing matrix and $(\cdot)^\top$ denotes the transpose. If M is equal to N , the inverse matrix of $\mathbf{A}_{k,l}$ exists and is called separation matrix ($\mathbf{W}_{k,l}$), so (1) can be written as:

$$\mathbf{s}_{k,l} = \mathbf{W}_{k,l} \mathbf{x}_{k,l}, \quad (2)$$

For BSE, the task changes from separating each source to extracting a source of interest (SOI), so (1) can be rewritten as:

$$\mathbf{x}_{k,l} = \mathbf{a}_{k,l} s_{k,l} + \mathbf{y}_{k,l}, \quad (3)$$

where $s_{k,l}$ represents the SOI and $\mathbf{y}_{k,l}$ represents the remaining part of the mixture, which are independent of each other. In general, we assume that the SOI corresponds to the first source in $\mathbf{s}_{k,l} = [s_{k,l}^1, \dots, s_{k,l}^N]^\top \in \mathbb{C}^N$, implying that the aboved source signal can be divided into $\mathbf{s}_{k,l} = [s_{k,l}, \mathbf{z}_{k,l}^\top]^\top$. Consequently, the mixing vector $\mathbf{a}_{k,l}$ is the first column vector in the $\mathbf{A}_{k,l}$. Correspondingly, the separation vector $\mathbf{w}_{k,l}$ associated with the SOI is the first row in the $\mathbf{W}_{k,l}$, which ensures $s_{k,l} = \mathbf{w}_{k,l}^\top \mathbf{x}_{k,l}$. $(\cdot)^\text{H}$ denotes the Hermitian transpose. In this case, we limit the dimension of $\mathbf{y}_{k,l}$ to a subspace of dimension $M - 1$, so $\mathbf{y}_{k,l}$ cannot be separated individually.

2.2 CSV Mixing Modeling

In the real situation, the mixing system is time-varying due to moving sources. Let divide the observed signal into T short intervals (blocks) based on time frames L . The length of each block is L_b , which satisfies $L_b = L/T$. Each block can be denoted by $t \in \{1, \dots, T\}$. In other words, for the observed signal of each block, we consider the source is approximately stationary and compute the corresponding mixing vector $\mathbf{a}_{k,t}$ and separation vector $\mathbf{w}_{k,t}$.

In the constant separating vector (CSV) mixing model, we separate the mixing vector $\mathbf{a}_{k,l}$ from the mixing matrix $\mathbf{A}_{k,l}$ and can derive a extremely effective parameterization involving only the mixing vector associated with the SOI. It considers that only the mixing vectors $\mathbf{a}_{k,l}$ are block-dependent ($\mathbf{a}_{k,t}$) while the separation filter $\mathbf{w}_{k,l}$ are constant over the blocks (\mathbf{w}_k). Specifically, for the t -th block, the mixing matrix $\mathbf{A}_{k,t}$ in (1) has the structure [10] (we omit frame index $l = 1, \dots, L$ from here):

$$\mathbf{A}_{k,t} = [\mathbf{a}_{k,t}, \mathbf{D}_{k,t}] = \begin{bmatrix} \gamma_{k,t} & \mathbf{h}_k^\text{H} \\ \mathbf{g}_{k,t} & \frac{1}{\gamma_{k,t}} (\mathbf{g}_{k,t} \mathbf{h}_k^\text{H} - \mathbf{I}_{M-1}) \end{bmatrix}. \quad (4)$$

Similarly, the structure of separation matrix $\mathbf{W}_{k,t}$ is

$$\mathbf{W}_{k,t} = \begin{bmatrix} \mathbf{w}_k^\text{H} \\ \mathbf{B}_{k,t} \end{bmatrix} = \begin{bmatrix} \beta_k^* & \mathbf{h}_k^\text{H} \\ \mathbf{g}_{k,t} & -\gamma_{k,t} \mathbf{I}_{M-1} \end{bmatrix}, \quad (5)$$

where $\gamma_{k,t}$ is the scalar, \mathbf{I}_M denotes the $M \times M$ identity matrix, the mixing vector $\mathbf{a}_{k,t}$ can be divided into $\mathbf{a}_{k,t} = [\gamma_{k,t}, \mathbf{g}_{k,t}^\top]^\top$, and the separation vector \mathbf{w}_k can be divided into $\mathbf{w}_k = [\beta_k, \mathbf{h}_k^\top]^\top$. The CSV mixing model assumes that the mixing vector $\mathbf{a}_{k,t}$ and the distribution of source can vary block by block, but the separation vector \mathbf{w}_k is constant between blocks. Combining the above equations and the relationship between the mixing matrix $\mathbf{A}_{k,t}$ and the separation matrix $\mathbf{W}_{k,t}$, it is easy to deduce that the relationship between the mixing vector $\mathbf{a}_{k,t}$ and the separation vector \mathbf{w}_k satisfies $\mathbf{w}_k^\text{H} \mathbf{a}_{k,t} = 1$ (called distortionless constraint), which means that

$$\beta_k^* \gamma_{k,t} + \mathbf{h}_k^\text{H} \mathbf{g}_{k,t} = 1. \quad (6)$$

In addition, the matrix $\mathbf{B}_{k,t}$ is called block matrix [9], which satisfies $\mathbf{B}_{k,t} \mathbf{a}_{k,t} = \mathbf{0}$. The remaining signal source can be represented by $\mathbf{z}_k = \mathbf{B}_k \mathbf{x}_k = \mathbf{B}_k \mathbf{y}_k$, where $\mathbf{y}_k = \mathbf{D}_k \mathbf{z}_k$. Therefore, for the t -th block, (1) can be rewritten as

$$\mathbf{x}_{k,t} = \begin{bmatrix} \gamma_{k,t} & \mathbf{h}_k^\text{H} \\ \mathbf{g}_{k,t} & \frac{1}{\gamma_{k,t}} (\mathbf{g}_{k,t} \mathbf{h}_k^\text{H} - \mathbf{I}_{M-1}) \end{bmatrix} \begin{bmatrix} s_{k,t} \\ \mathbf{z}_{k,t} \end{bmatrix}. \quad (7)$$

In the offline CSV-AuxIVE algorithm, the auxiliary function approach [5, 6] was used to update the separation filter \mathbf{w}_k efficiently. Here, the filter \mathbf{w}_k is updated with several iterations as follows:

$$\mathbf{w}_k = \left(\sum_{t=1}^T \frac{\mathbf{V}_{k,t}}{\hat{\sigma}_{k,t}^2} \right)^{-1} \sum_{t=1}^T \frac{\mathbf{w}_k^\text{H} \mathbf{V}_{k,t} \mathbf{w}_k}{\hat{\sigma}_{k,t}^2} \mathbf{a}_{k,t} = \mathbf{P}_k^{-1} \mathbf{Q}_k, \quad (8)$$

where $\hat{\sigma}_{k,t}^2$ denotes the variance of the source for the t -th block and $\mathbf{V}_{k,t}$ is the weighted spatial covariance matrix.

3. Proposed method: Online-CSVAuxIVE

In [9], although the authors divide the observed signal into several blocks in the frequency domain, the method does not process each block separately. It means that we must know all of the observed signals before extraction, which is called offline CSV-AuxIVE algorithm. However, for some devices that require real-time processing of signals, such as hearing aids, multi-person online conferencing devices, the offline CSV-AuxIVE algorithm may lead to high latency and is not applicable. Therefore, to further improve the applicability of this algorithm, this paper proposes an online version of the CSV-AuxIVE algorithm, with the specific idea of processing each piece of data sequentially and updating it to obtain a different mixing vector $\mathbf{a}_{k,t}$ for each of block. In this paper,

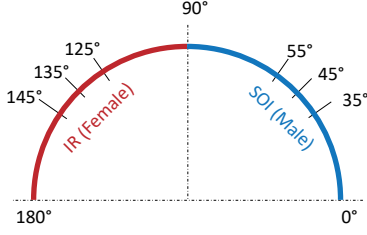


Figure 1: Layout of experimental environment.

Table 1: The conditions of sources (range & speed)

Index	Speed	Range (start angle → ending angle)		
			SOI (male)	IR (female)
1	Static	N/A	45°	135°
2	Slow	Small	35° → 55°	125° → 145°
3	Fast	Small	35° → 55°	125° → 145°
4	Slow	Large	0° → 90°	90° → 180°
5	Fast	Large	0° → 90°	90° → 180°
6	Very Fast	Large	0° → 90°	90° → 180°

we call the original algorithm the offline CSV-AuxIVE model and the proposed algorithm the online model.

In this paper, we introduce the recursive update in each block to calculate $P_{k,t}$ and $Q_{k,t}$ as:

$$P_{k,t} = \alpha P_{k,t-1} + (1 - \alpha) \frac{V_{k,t}}{\hat{\sigma}_{k,t}^2}, \quad (9)$$

$$Q_{k,t} = \alpha Q_{k,t-1} + (1 - \alpha) \frac{\mathbf{w}_{k,t}^H V_{k,t} \mathbf{w}_{k,t}}{\sigma_{k,t}^2} \mathbf{a}_{k,t}, \quad (10)$$

where α is the forgetting factor. $P_{k,t}$ and $Q_{k,t}$ are updated after several iterations. The pseudo code is summarized in Algorithm 1. For the proposed algorithm, there are three parameters that have a large impact on the results: the forgetting factor α , the length of the blocks L_b , and the initial value setting of the matrix Q_k .

4. Experiment

We verified the effectiveness of online CSV-AuxIVE for sources at different ranges and speeds, comparing with the conventional online AuxIVA. The parameters of the compared algorithms are:

- (1) online CSV-AuxIVE with $L_b = 150$ frames and $\alpha = 0.3$.
- (2) online AuxIVA with $L_b = 20$ frames and $\alpha = 0.95$.

4.1 Experiment conditions

The dataset is from Experiment 4.2 in [9]. The authors attached wireless speakers to a manually driven rotating arm, manually controlled the range and speed of movement of

Algorithm 1 Pseudo-code of online CSV-AuxIVE

Input: $\mathbf{x}_{k,t}$, L_b , α , Maxiter

Output: $\mathbf{y}_{k,t}$

```

1: (Initialization)
2:  $\mathbf{w}_{k,0} = [1, 0, \dots, 0]^T \in \mathbb{C}^{m \times n}$ 
3:  $\mathbf{P}_{k,0} = \mathbf{R}_{k,0} = 10^{-12} \times \text{Identity Matrix} \in \mathbb{C}^{m \times m}$ 
4:  $\mathbf{Q}_{k,0} = \text{Null Matrix} \in \mathbb{C}^{m \times n}$ 
5: for  $t = 1 \rightarrow T$  do
6:   (Enhancement)
7:    $\mathbf{y}_{k,t} = \mathbf{w}_{k,t-1}^H \mathbf{x}_{k,t}$ 
8:    $[\gamma_{k,t-1}, \mathbf{g}_{k,t-1}^T]^T = \mathbf{a}_{k,t-1}$ 
9:    $y_{k,t} = \gamma_{k,t-1} \mathbf{y}_{k,t}$ 
10:  (Updating)
11:  for  $k = 1 \rightarrow K$  do
12:     $\hat{\mathbf{C}}_{k,t-1} = E[\mathbf{x}_{k,t-1} \mathbf{x}_{k,t-1}^H]$ 
13:  end for
14:  for Iter = 1 → Maxiter do
15:    for  $k = 1 \rightarrow K$  do
16:       $\sigma_{k,t-1} = \sqrt{\mathbf{w}_{k,t-1}^H \hat{\mathbf{C}}_{k,t-1} \mathbf{w}_{k,t-1}}$ 
17:    end for
18:     $r_{t-1} = \sqrt{\sum_k |\mathbf{w}_{k,t-1}^H \mathbf{x}_{k,t-1}|^2}$ 
19:    for  $k = 1 \rightarrow K$  do
20:       $\mathbf{a}_{k,t-1} = \frac{\hat{\mathbf{C}}_{k,t-1} \mathbf{w}_{k,t-1}}{\mathbf{w}_{k,t-1}^H \hat{\mathbf{C}}_{k,t-1} \mathbf{w}_{k,t-1}}$ 
21:       $\mathbf{V}_{k,t-1} = E[\mathbf{x}_{k,t-1} \mathbf{x}_{k,t-1}^H / r_{t-1}]$ 
22:       $\mathbf{P}_{k,t-1} = \alpha \mathbf{P}_{k,t-2} + (1 - \alpha) \frac{\mathbf{V}_{k,t-1}}{\sigma_{k,t-1}^2}$ 
23:       $\mathbf{a}_{\text{tmp}} = \frac{\mathbf{w}_{k,t-1}^H \mathbf{V}_{k,t-1} \mathbf{w}_{k,t-1}}{\sigma_{k,t-1}^2} \mathbf{a}_{k,t-1}$ 
24:       $\mathbf{Q}_{k,t-1} = \alpha \mathbf{Q}_{k,t-2} + (1 - \alpha) \mathbf{a}_{\text{tmp}}$ 
25:       $\mathbf{R}_{k,t-1} = \alpha \mathbf{R}_{k,t-2} + (1 - \alpha) \mathbf{V}_{k,t-1}$ 
26:       $\mathbf{w}_{k,t-1} = \mathbf{P}_{k,t-1}^{-1} \mathbf{Q}_{k,t-1}$ 
27:       $\mathbf{w}_{k,t-1} = \frac{\mathbf{w}_{k,t-1}}{\sqrt{\mathbf{w}_{k,t-1}^H \mathbf{R}_{k,t-1} \mathbf{w}_{k,t-1}}}$ 
28:    end for
29:  end for
30: end for

```

the two sources, and recorded the mixed signals of the two source signals in different states using four linear microphones spaced 16 cm apart. The male voice represents the SOI and the female voice represents the interference (IR). The layout of the experimental environment is shown in Fig. 1, with two source signals moving on a semicircular trajectory with a radius of 1 m.

From Table 1, it can be seen that each source signal has 6 different conditions, and after combining with each other, there are 36 different mixed signals.

The signal-to-noise ratio (SDR) of segmentation followed by averaging is used as a criterion for judging the performance of the algorithm. SDR is calculated using BSS-eval [11]. The length of the segments is 1 s. The experimental conditions are shown in Table 2.

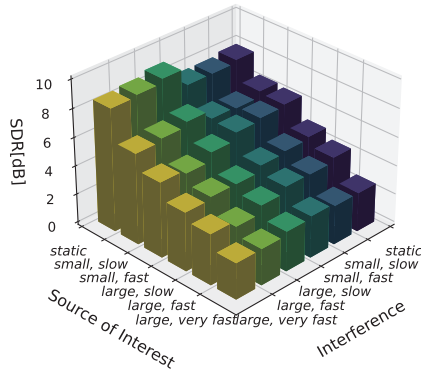


Figure 2: Average SDR of 36 mixtures by online AuxIVA. ($L_b = 20, \alpha = 0.95$)

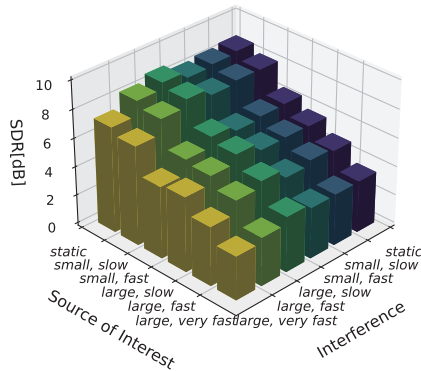


Figure 3: Average SDR of 36 mixtures by online CSV-AuxIVE. ($L_b = 150, \alpha = 0.3$)

4.2 Results

We compare the SDR between online AuxIVA and online CSV-AuxIVE using Figs. 2 and 3. As shown in Fig. 2, AuxIVA performs well only when the SOI is stationary. And when the SOI starts to move, its performance decreases. On the other hand in Fig. 3, online CSV-AuxIVE shows less sensitivity to the movement of the SOI. If the SOI moves slowly, online CSV-AuxIVE performs better than online AuxIVA regardless of the movement of the IR. Focusing only the condition of the SOI movement as {Slow, Small}, CSV-AuxIVE had about 1.5 dB higher SDR improvement than online AuxIVA regardless of the variation of the moving range and the moving speed of the interference.

5. Conclusion

In this paper, for the moving BSE, we found a novel and effective offline CSV mixing model. To make this algorithm applicable to devices with the requirement of real-time pro-

Table 2: Experimental conditions

Sampling rate	16 kHz
STFT window function	Hann
STFT window length	1024 samples
# of Iteration (Maxiter)	100
Initial value of Q	Zero matrix
Initial value of P	$10^{-12} \times$ Identity matrix
Room dimensions	$12 \times 8 \times 2.6$ m

cessing, we have proposed an online version of the CSV-AuxIVE algorithm. Unlike the offline algorithm, the proposed algorithm needs to know only a part of the mixed signal and then process the input observed signals sequentially. The results confirmed that the average SDR of the online CSV-AuxIVE was about 1.5 dB higher than that of the online AuxIVA under the same conditions of the target source, regardless of the variation of the moving range and the moving speed of the interference.

References

- [1] A. Hyvärinen and E. Oja, "Independent component analysis: Algorithms and applications," *Neural networks*, vol. 13, no. 4-5, pp. 411–430, 2000.
- [2] A. Tharwat, "Independent component analysis: An introduction," *Applied Computing and Informatics*, 2020.
- [3] T. Kim, T. Eltoft, and T.-W. Lee, "Independent vector analysis: An extension of ICA to multivariate components," in *Proc. ICA*, 2006, pp. 165–172.
- [4] T. Kim, I. Lee, and T.-W. Lee, "Independent vector analysis: Definition and algorithms," in *Fortieth Asilomar Conference on Signals, Systems and Computers*, 2006, pp. 1393–1396.
- [5] N. Ono, "Stable and fast update rules for independent vector analysis based on auxiliary function technique," in *Proc. WASPAA*, 2011, pp. 189–192.
- [6] N. Ono, "Auxiliary-function-based independent vector analysis with power of vector-norm type weighting functions," in *Proc. APSIPA ASC*, 2012, pp. 1–4.
- [7] T. Taniguchi, N. Ono, A. Kawamura, and S. Sagayama, "An auxiliary-function approach to online independent vector analysis for real-time blind source separation," in *Proc. HSCMA*, 2014, pp. 107–111.
- [8] Z. Koldovský, J. Málek, and J. Janský, "Extraction of independent vector component from underdetermined mixtures through block-wise determined modeling," in *Proc. ICASSP*, 2019, pp. 7903–7907.
- [9] J. Janský, Z. Koldovský, J. Málek, T. Kounovský, and J. Čmejla, "Auxiliary function-based algorithm for blind extraction of a moving speaker," *EURASIP JASMP*, vol. 2022, no. 1, pp. 1–16, 2022.
- [10] Z. Koldovský and P. Tichavský, "Gradient algorithms for complex non-gaussian independent component/vector extraction, question of convergence," *IEEE trans. SP*, vol. 67, no. 4, pp. 1050–1064, 2018.
- [11] E. Vincent, R. Gribonval, and C. Févotte, "Performance measurement in blind audio source separation," *IEEE trans. ASLP*, vol. 14, no. 4, pp. 1462–1469, 2006.