



Joint Optimization of Sampling Rate Offsets Based on Entire Signal Relationship Among Distributed Microphones

Yoshiki Masuyama¹, Kouei Yamaoka¹, Nobutaka Ono¹

¹Tokyo Metropolitan University, Tokyo, Japan

masuyama-yoshiki@ed.tmu.ac.jp

Abstract

In this paper, we propose to simultaneously estimate all the sampling rate offsets (SROs) of multiple devices. In a distributed microphone array, the SRO is inevitable, which deteriorates the performance of array signal processing. Most of the existing SRO estimation methods focused on synchronizing two microphones. When synchronizing more than two microphones, we select one reference microphone and estimate the SRO of each non-reference microphone independently. Hence, the relationship among signals observed by non-reference microphones is not considered. To address this problem, the proposed method jointly optimizes all SROs based on a probabilistic model of a multichannel signal. The SROs and model parameters are alternately updated to increase the log-likelihood based on an auxiliary function. The effectiveness of the proposed method is validated on mixtures of various numbers of speakers.

Index Terms: Wireless acoustic sensor network, distributed microphone array, sampling rate offset, auxiliary function.

1. Introduction

Microphone array signal processing, including blind source separation (BSS) [1, 2], is a fundamental technique with various applications such as automatic speech recognition [3] and sound event detection [4]. Although more microphones are desirable to improve the performance of BSS [5, 6], it is costly to prepare a large microphone array. This limits the number of microphones and the array size in many applications. To address this problem, distributed microphone array (DMA) processing has gained considerable attention [7, 8]. A DMA exploits a set of microphones on multiple devices, including tablets and smartphones, and does not require any specialized devices. By using a DMA, we can acquire a large number of observations and conduct array signal processing such as sound source localization [9], speech enhancement [10, 11], and BSS [12].

In a DMA, microphones are connected to device-dependent analog-to-digital converters, and their sampling rates are slightly different even when the nominal ones are the same. These sampling rate offsets (SROs) deteriorate the performance of array signal processing including BSS [13]. We should thus estimate and compensate for the SROs in advance [13–25]. Since array signal processing is often conducted in the time-frequency (T-F) domain, an SRO model in the T-F domain, called the linear phase drift (LPD) model, has been widely used [15]. The LPD model considers that the SRO changes the phase of the short-time Fourier transform (STFT) coefficients linearly with time and frequency. Various SRO estimation approaches have been developed based on this model [13–18].

The first approach, called the coherence-drift-based approach [15–17], computes the complex coherence between the signals observed by the reference and non-reference microphones. The SRO is estimated from the ratio of the complex

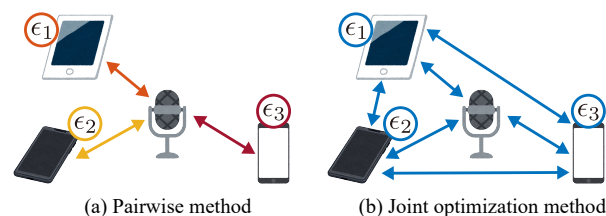


Figure 1: Illustration of (a) existing pairwise synchronization and (b) proposed joint-optimization-based synchronization. As an example, the center microphone is considered as the reference one. Each color corresponds to an optimization problem.

coherence between successive time frames. The second approach relies on a probabilistic model of the two-channel signals [13, 14]. This approach assumes that the STFT coefficients of the synchronized signals follow a multivariate complex Gaussian distribution. On the basis of this assumption, the SRO is estimated in a maximum likelihood manner. The third approach estimates the SRO to maximize the correlation between the STFT coefficients [18]. It is shown that the correlation takes the largest value when the SRO is accurately compensated for. The second and third approaches often require a higher computational cost to search for the optimal SRO, but they have achieved promising results.

The aforementioned approaches are based on the relationship between the signals observed by one reference microphone and one non-reference microphone [13–21]. When synchronizing more than two microphones, we should select the reference microphone and estimate the SRO of each non-reference microphone independently, as depicted in Fig. 1-(a). We call this type of synchronization the pairwise method. One of the drawbacks of this method is that its performance depends on the selection of the reference microphone. This is because the relationship among the non-reference microphones is not taken into account. Meanwhile, it is still a challenging task to select the optimal reference microphone in terms of synchronization. It is thus desirable to exploit full spatial information of the signals observed by a DMA regardless of the reference microphone.

To this end, we propose to consider the relationships between all microphone pairs and optimize all the SROs jointly as depicted in Fig. 1-(b). The proposed method estimates the SROs in a maximum likelihood manner through a probabilistic model of the entire multichannel signal. Although the golden section search has been used to optimize each SRO independently [13], it is not directly applicable to the joint optimization of all SROs. To address this problem, we present an iterative algorithm using an auxiliary function and guarantee the non-decrease property of the log-likelihood in the algorithm. Our experimental results showed that the proposed joint optimization method outperformed the existing pairwise methods [13, 18].

2. SRO Estimation in T-F Domain

2.1. SRO Model in T-F Domain

Let us assume that an M -channel signal is observed by a DMA, where the 0th microphone is selected as the reference one. The sampling rate of the m th microphone is given by

$$r_m = (1 + \epsilon_m)r_0, \quad (1)$$

where $m = 0, \dots, M - 1$ is the microphone index, ϵ_m is the SRO of the m th microphone, and $\epsilon_0 = 0$. Let $\tilde{\chi}_m$ be a continuous signal to be measured by the m th microphone. Then, the τ th entry of its discrete version is given by

$$\chi_m[\tau] = \tilde{\chi}_m \left(\frac{\tau}{(1 + \epsilon_m)r_0} + \Delta_m \right), \quad (2)$$

where Δ_m is the sampling time offset (STO) of the m th microphone. Although an accurate estimation of the STO is not easy, its small error is acceptable for BSS. We thus hereafter assume that the STO is already recovered by an existing method [13].

Since array signal processing is often conducted in the T-F domain owing to its efficiency, the LPD model has been widely used to estimate and compensate for the SRO. Let the STFT of the discrete signal χ_m with a window g of length L be

$$x_m[t, f] = \sum_{l=0}^{L-1} \chi_m[l + at]g[l]e^{-2\pi jfl/F}, \quad (3)$$

where j is the imaginary unit, a is the window shift, and $t = 0, \dots, T - 1$ and $f = 0, \dots, F - 1$ are the time frame and frequency bin indices, respectively. The LPD model represents the SRO by a phase modification of STFT coefficients and compensates for it as follows [15]:

$$\hat{x}_m[t, f] = x_m[t, f] \exp \left(\frac{2\pi jatf\epsilon_m}{F} \right). \quad (4)$$

That is, $\hat{x}_m[t, f]$ can be interpreted as the STFT coefficient of the synchronized signal when ϵ_m is accurately estimated.

2.2. Maximum Likelihood Estimation of SRO

Based on the LPD model, the probabilistic-model-based approach estimates the SRO in a maximum likelihood manner [13, 14]. In this approach, the compensated STFT coefficients $\hat{\mathbf{x}}[t, f] = [\hat{x}_0[t, f], \dots, \hat{x}_{M-1}[t, f]]^T$ are assumed to follow a multivariate complex Gaussian distribution:

$$\hat{\mathbf{x}}[t, f] \sim \mathcal{N}_c(\mathbf{0}, \mathbf{V}[f]), \quad (5)$$

where $\mathbf{V}[f]$ is the spatial covariance matrix (SCM) of the synchronized signals. The probabilistic model in (5) implies that the sound sources do not move and their powers are stationary.

The existing methods have addressed the case $M = 2$. In such a case, the log-likelihood for (5) can be reformulated to the following univariate objective function with respect to ϵ_1 [14]:

$$\mathcal{I}(\epsilon_1) = - \sum_f \log \left(\sum_t |x_0[t, f]|^2 \sum_t |\hat{x}_1[t, f]|^2 - \left| \sum_t x_0[t, f] \hat{x}_1[t, f] \right|^2 \right), \quad (6)$$

where $\hat{x}_1[t, f]$ depends on ϵ_1 as shown in (4). The SRO is estimated by maximizing this objective function. As the objective function is usually locally unimodal around the global optimum, the golden section search initialized by a coarse grid search can find the optimal ϵ_1 efficiently.

3. Proposed Iterative SRO Estimation

In this section, we propose a joint optimization method of all SROs of an arbitrary number of microphones. The proposed method alternately updates the SROs and SCMs to increase the log-likelihood for (5). In the update of SROs, we maximize an auxiliary function instead of the intractable log-likelihood.

3.1. Joint Optimization Problem of all SROs

Most of the existing SRO estimation methods focused on synchronizing two microphones [13–18]. A naïve extension of these methods to the case $M > 2$ is the pairwise method as illustrated in Fig. 1-(a). This method considers the signals observed by the 0th and $m \neq 0$ th microphones as two-channel signals. Then, the SRO of the m th microphone is estimated independently, where the relationship among signals observed by the non-reference microphones is not considered. To improve the performance of SRO estimation, the relationships among all observed signals should be exploited regardless of the reference microphone as depicted in Fig. 1-(b).

To this end, we propose to jointly optimize all SROs in a maximum likelihood manner. By denoting the SROs as $\epsilon = [\epsilon_0, \dots, \epsilon_{M-1}]^T$, the log-likelihood of the SROs and SCMs for the multivariate complex Gaussian model in (5) is given by

$$\begin{aligned} \mathcal{L}(\epsilon, \mathbf{V}[0], \dots, \mathbf{V}[F-1]) \\ = \sum_{f=0}^{F-1} \sum_{t=0}^{T-1} -\log \det(\mathbf{V}[f]) - \hat{\mathbf{x}}^H[t, f] \mathbf{V}^{-1}[f] \hat{\mathbf{x}}[t, f], \end{aligned} \quad (7)$$

where $(\cdot)^H$ is the Hermitian transpose, and a constant term is omitted. This log-likelihood is difficult to maximize with respect to the SROs in a closed form. Furthermore, the golden section search used in the existing method is not applicable to the maximization of such a multivariate function.

3.2. Alternative Updates of SROs and SCMs

To maximize (7) with respect to ϵ and $\mathbf{V}[f]$, we develop an iterative algorithm that alternately updates them. For fixed SROs, the SCMs that maximize the log-likelihood in (7) are easily obtained as

$$\mathbf{V}[f] \leftarrow \frac{1}{T} \sum_{t=0}^{T-1} \hat{\mathbf{x}}[t, f] \hat{\mathbf{x}}^H[t, f]. \quad (8)$$

Meanwhile, it is difficult to maximize the log-likelihood with respect to ϵ even with fixed $\mathbf{V}[f]$. Hence, we use the auxiliary function method [26, 27] that can handle multivariate non-convex optimization problems and has achieved promising results in array signal processing [21, 28, 29]. To be specific, we consider the following objective function by removing a term independent of ϵ from (7):

$$\mathcal{J}(\epsilon) = - \sum_{f=0}^{F-1} \sum_{t=0}^{T-1} \hat{\mathbf{x}}^H[t, f] \mathbf{V}^{-1}[f] \hat{\mathbf{x}}[t, f]. \quad (9)$$

Then, we introduce an auxiliary variable $\tilde{\epsilon}$ and derive an auxiliary function $\mathcal{Q}(\epsilon | \tilde{\epsilon})$ that satisfies the following properties:

- For all ϵ and $\tilde{\epsilon}$, $\mathcal{J}(\epsilon) \geq \mathcal{Q}(\epsilon | \tilde{\epsilon})$.
- For all ϵ , $\mathcal{Q}(\epsilon | \epsilon) = \mathcal{J}(\epsilon)$.

The auxiliary function method alternately updates ϵ and $\tilde{\epsilon}$ to maximize $\mathcal{Q}(\epsilon | \tilde{\epsilon})$. Thanks to the properties of the auxiliary function, $\arg\max_{\tilde{\epsilon}} \mathcal{Q}(\epsilon | \tilde{\epsilon})$ is obtained as $\tilde{\epsilon} \leftarrow \epsilon$. Meanwhile,

it depends on the auxiliary function whether $\arg\max_{\epsilon} \mathcal{Q}(\epsilon | \tilde{\epsilon})$ is obtained in a closed form or not. The detail of the proposed auxiliary function is explained in the next subsection.

3.3. Auxiliary Function for Jointly Updating SROs

To derive the auxiliary function $\mathcal{Q}(\epsilon | \tilde{\epsilon})$, we use the following tractable lower bound of the negative cosine function. Let $\alpha \in \mathbb{R}_+$, $\beta \in \mathbb{R}$, $\gamma \in \mathbb{R}$, $\theta \in \mathbb{R}$, and $\tilde{\theta} \in \mathbb{R}$. Then, the following inequality holds¹ [29]:

$$-\alpha \cos(\beta\theta + \gamma) \geq -\frac{\alpha}{2} \text{sinc}(\beta\tilde{\theta} - \phi)(\beta\theta - \phi)^2 + \eta, \quad (10)$$

where $\text{sinc}(\cdot) = \sin(x)/x$ if $x \neq 0$ and 1 otherwise, and

$$\phi = 2\pi \left\lfloor \frac{\beta\tilde{\theta} + \gamma}{2\pi} \right\rfloor + \pi - \gamma, \quad (11)$$

$$\eta = \frac{\alpha}{2} \text{sinc}(\beta\tilde{\theta} - \phi)(\beta\tilde{\theta} - \phi)^2 - \alpha \cos(\beta\tilde{\theta} + \gamma). \quad (12)$$

The equality in (10) holds when $\theta = \tilde{\theta}$.

To reformulate the objective function (9) by a sum of the negative cosine functions, we define

$$\Upsilon[t, f] = \text{diag}(\mathbf{x}[t, f])^H \mathbf{V}^{-1}[f] \text{diag}(\mathbf{x}[t, f]), \quad (13)$$

where $\text{diag}(\cdot)$ returns a diagonal matrix whose diagonal entries are its input. By leveraging the conjugate symmetry of $\Upsilon[t, f]$ and the Euler's formula [30], the objective function can be reformulated as:

$$\mathcal{J}(\epsilon) = \sum_{t=0}^{T-1} \sum_{f=0}^{F-1} \sum_{m=0}^{M-1} \sum_{n=0}^{M-1} \underline{\mathcal{J}}_{t,f,m,n}(\epsilon), \quad (14)$$

$$\underline{\mathcal{J}}_{t,f,m,n}(\epsilon) = -|\Upsilon_{m,n}[t, f]| \cos\left(\omega[t, f](\epsilon_n - \epsilon_m) + \angle\Upsilon_{m,n}[t, f]\right), \quad (15)$$

where $\omega[t, f] = 2\pi atf/F$, and $\angle\cdot$ denotes the principal value of the complex-argument.

Since the entry-wise objective function in (15) is a negative cosine function with respect to $\epsilon_n - \epsilon_m$, we obtain the following auxiliary function based on (10):

$$\underline{\mathcal{Q}}_{t,f,m,n}(\epsilon | \tilde{\epsilon}) = -\lambda_{m,n}[t, f] \left(\omega[t, f](\epsilon_n - \epsilon_m) - \mu_{m,n}[t, f] \right)^2 + \nu_{m,n}[t, f], \quad (16)$$

where $\nu_{m,n}[t, f]$ does not depend on ϵ , and

$$\xi_{m,n}[t, f] = \omega[t, f](\tilde{\epsilon}_n - \tilde{\epsilon}_m), \quad (17)$$

$$\lambda_{m,n}[t, f] = \frac{|\Upsilon_{m,n}[t, f]|}{2} \text{sinc}\left(\xi_{m,n}[t, f] - \mu_{m,n}[t, f]\right), \quad (18)$$

$$\mu_{m,n}[t, f] = 2\pi \left\lfloor \frac{\xi_{m,n}[t, f] + \angle\Upsilon_{m,n}[t, f]}{2\pi} \right\rfloor + \pi - \angle\Upsilon_{m,n}[t, f], \quad (19)$$

where $\xi_{m,n}[t, f] - \mu_{m,n}[t, f]$ is in $[-\pi, \pi)$, and thus $\lambda_{m,n}[t, f] \geq 0$. We stress that (15) and (16) correspond to the left and right side of (10), respectively. Finally, we obtain the auxiliary function $\mathcal{Q}(\epsilon | \tilde{\epsilon})$ by summing up $\underline{\mathcal{Q}}_{t,f,m,n}(\epsilon | \tilde{\epsilon})$ for all T-F bins and microphone pairs.

¹Although the original paper [29] derived an upper bound of a negative cosine function, the lower bound in (10) can be derived in a similar manner. We thus omit the detailed proof of the inequality.

Algorithm 1 Iterative Algorithm to Estimate SROs

Input: Initial estimate of SROs ϵ , \mathbf{D} , \mathbf{u} , $\omega[t, f]$

Output: Final estimate of SROs ϵ

for $k = 0, \dots, K - 1$ **do**

$$\hat{x}_m[t, f] = x_m[t, f] \exp\left(\frac{2\pi jatf\epsilon_m}{F}\right)$$

$$\mathbf{V}[f] \leftarrow (1/T) \sum_{t=0}^{T-1} \hat{\mathbf{x}}[t, f] \hat{\mathbf{x}}^H[t, f]$$

$$\Upsilon[t, f] = \text{diag}(\mathbf{x}[t, f])^H \mathbf{V}^{-1}[f] \text{diag}(\mathbf{x}[t, f])$$

for $k' = 0, \dots, K' - 1$ **do**

$$\tilde{\epsilon} \leftarrow \epsilon$$

$$\xi_{m,n}[t, f] = \omega[t, f](\tilde{\epsilon}_n - \tilde{\epsilon}_m)$$

$$\lambda_{m,n}[t, f] = \frac{|\Upsilon_{m,n}[t, f]|}{2} \text{sinc}\left(\xi_{m,n}[t, f] - \mu_{m,n}[t, f]\right)$$

$$\mu_{m,n}[t, f] = 2\pi \left\lfloor \frac{\xi_{m,n}[t, f] + \angle\Upsilon_{m,n}[t, f]}{2\pi} \right\rfloor + \pi - \angle\Upsilon_{m,n}[t, f]$$

$$\mathbf{A} = \sum_{f=0}^{F-1} \sum_{t=0}^{T-1} \omega^2[t, f] \mathbf{\Lambda}[t, f]$$

$$\mathbf{b} = \sum_{f=0}^{F-1} \sum_{t=0}^{T-1} \omega[t, f] \mathbf{\Lambda}[t, f] \boldsymbol{\mu}[t, f]$$

$$\begin{pmatrix} \epsilon \\ \rho \end{pmatrix} \leftarrow \begin{pmatrix} \mathbf{D}^T \mathbf{A} \mathbf{D} & \mathbf{u} \\ \mathbf{u}^T & 0 \end{pmatrix}^{-1} \begin{pmatrix} \mathbf{D}^T \mathbf{b} \\ 0 \end{pmatrix}$$

end for

end for

Since the auxiliary function $\mathcal{Q}(\epsilon | \tilde{\epsilon})$ is a sum of negative quadratic functions, we can easily maximize it under the constraint $\epsilon_0 = 0$. By considering the Karush–Kuhn–Tucker (KKT) condition, the optimal ϵ is obtained by solving the following linear equation [31]:

$$\begin{pmatrix} \mathbf{D}^T \mathbf{A} \mathbf{D} & \mathbf{u} \\ \mathbf{u}^T & 0 \end{pmatrix} \begin{pmatrix} \epsilon^* \\ \rho^* \end{pmatrix} = \begin{pmatrix} \mathbf{D}^T \mathbf{b} \\ 0 \end{pmatrix}, \quad (20)$$

where $\mathbf{u} = [1, 0, \dots, 0]^T \in \mathbb{R}^M$, $\rho^* \in \mathbb{R}$ is the KKT multiplier, $\mathbf{D} \in \mathbb{R}^{M^2 \times M}$ is a matrix that computes the difference in ϵ as

$$\epsilon_n - \epsilon_m = (\mathbf{D}\epsilon)_{mM+n}, \quad (21)$$

and $\mathbf{A} \in \mathbb{R}^{M^2 \times M^2}$ and $\mathbf{b} \in \mathbb{R}^{M^2}$ are given by

$$\mathbf{A} = \sum_{f=0}^{F-1} \sum_{t=0}^{T-1} \omega^2[t, f] \mathbf{\Lambda}[t, f], \quad (22)$$

$$\mathbf{b} = \sum_{f=0}^{F-1} \sum_{t=0}^{T-1} \omega[t, f] \mathbf{\Lambda}[t, f] \boldsymbol{\mu}[t, f]. \quad (23)$$

Here, $\mathbf{\Lambda}[t, f]$ is a diagonal matrix whose $(mM + n, mM + n)$ th entry is given by $\lambda_{m,n}[t, f]$, and the $(mM + n)$ th entry of $\boldsymbol{\mu}[t, f]$ is given by $\mu_{m,n}[t, f]$. In each inner iteration, ϵ is updated to maximize the auxiliary function $\mathcal{Q}(\epsilon | \tilde{\epsilon})$ as follows:

$$\begin{pmatrix} \epsilon \\ \rho \end{pmatrix} \leftarrow \begin{pmatrix} \mathbf{D}^T \mathbf{A} \mathbf{D} & \mathbf{u} \\ \mathbf{u}^T & 0 \end{pmatrix}^{-1} \begin{pmatrix} \mathbf{D}^T \mathbf{b} \\ 0 \end{pmatrix}. \quad (24)$$

The proposed algorithm is summarized in Algorithm 1, where $k = 0, \dots, K - 1$ and $k' = 0, \dots, K' - 1$ are the iteration counters. In each outer iteration, the auxiliary function method is used to update the SROs K' times. Owing to the property of the auxiliary function $\mathcal{Q}(\epsilon | \tilde{\epsilon})$, this algorithm ensures that the log-likelihood $\mathcal{L}(\epsilon, \mathbf{V}[0], \dots, \mathbf{V}[F - 1])$ does not decrease. When $M = 2$, the proposed method aims to maximize the same objective function considered in the existing methods [13, 14]. When $M > 2$, the proposed method can consider the consistency of the estimated SROs based on the entire relationship among signals observed by the non-reference microphones. On the other hand, the pairwise method cannot leverage the relationships due to their separate optimization.

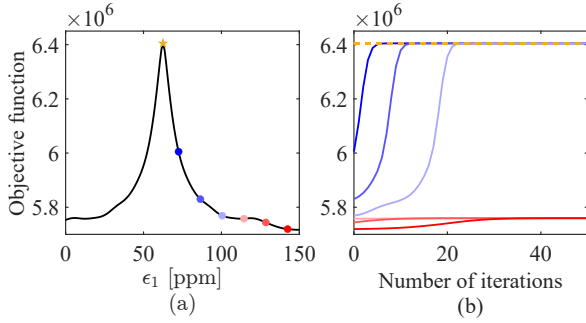


Figure 2: (a) Objective function with respect to ϵ_1 in parts per million (ppm). (b) Objective function with respect to number of iterations. Each solid line uses ϵ_1 at the circle of the same color in (a) as the initial value. The yellow star and dotted lines correspond to the oracle ϵ_1 .

4. Experimental Evaluations

In this section, we evaluated the proposed method on multichannel speech mixtures that imitate meeting recordings with some speakers. We first investigated the convergence of the proposed method. To visualize the log-likelihood, we used a two-channel speech signal where the log-likelihood can be reformulated as the univariate function with respect to ϵ_1 . Second, the effectiveness of the joint optimization of all SROs was validated on the four-channel speech mixtures.

4.1. Convergence of Proposed Method

In this experiment, we used a two-channel signal with a single speaker. The source signal of 10 s length was generated by concatenating utterances in the Voice Conversion Challenge (VCC) 2018 dataset [32]. The source signal was downsampled to 16000 Hz. To synthesize the reverberant signal, we performed room simulations using the `pyroomacoustics` toolbox [33]. The sound source and microphones were randomly located in a room of 6.0 m \times 8.0 m \times 4.0 m size. The reverberation time was also randomly sampled from [0.2, 0.4] s. The signal measured by the non-reference microphone was further resampled at 16001 Hz. For the STFT, the 2048-point-long Hann window was used with 1024-point shifts, where the number of DFT points was 4096. In Algorithm 1, K' was 1.

Fig. 2-(a) shows the log-likelihood in (7) as a function of ϵ_1 , where the SCMs were calculated by (8). The yellow star corresponds to the oracle ϵ_1 , and multiple initial values of the proposed method are depicted by circles. Fig. 2-(b) shows the convergence of the proposed method with different initial values. According to Fig 2-(a), the log-likelihood can be viewed as a unimodal function around the oracle ϵ_1 . As a result, the proposed method converged close to the oracle SRO when the initial value was appropriate. Even when the initial value was outside of the appropriate interval, the proposed method converged to a local maximum as theoretically guaranteed by the property of the auxiliary function method.

4.2. Synchronization of More Than Two Microphones

To confirm the effectiveness of the joint optimization of all SROs, we evaluated the performance of SRO estimation using four distributed microphones. In addition to the one-speaker signals, we synthesized two and three-speaker mixtures. In each case, 10 signals of 30 s length were synthesized. The signals

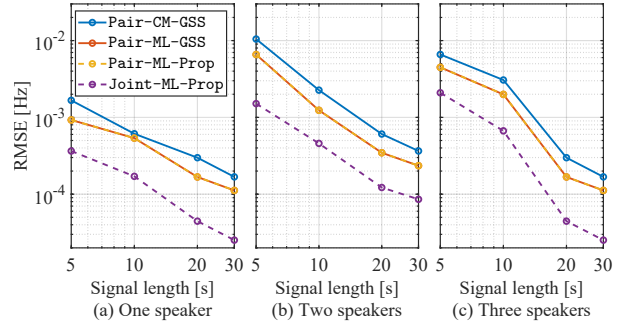


Figure 3: RMSEs of the estimated sampling rates for (a) one-speaker signals, (b) two-speaker mixtures, (c) three-speaker mixtures observed by four microphones.

measured by non-reference microphones were randomly resampled at $r_m \in [15999, 16001]$ Hz. Other conditions were the same as in the previous experiment.

The proposed method (`Joint-ML-Prop`) was compared with the pairwise maximum likelihood estimation method by the golden section search (`Pair-ML-GSS`) [13] and by Algorithm 1 (`Pair-ML-Prop`). We also investigated the performance of the pairwise correlation maximization method (`Pair-CM-GSS`) [18]. For all methods, we initialized SROs by a coarse grid search in a pairwise manner. The search range was from -100 ppm to 100 ppm with 100 grids, which is finer than the grids in Fig. 2. We expect that this initialization enabled us to avoid bad local optima. Then, the initial estimate was refined by the golden section search or Algorithm 1.

Fig. 3 shows the root mean square errors (RMSEs) of the estimated sampling rates for different signal lengths. `Pair-ML-GSS` and `Pair-ML-Prop` resulted in the same RMSE. That is, the difference of the optimization algorithms did not affect the performance in our experimental conditions. Meanwhile, `Joint-ML-Prop` outperformed all pairwise methods regardless of the number of speakers. This result confirmed the effectiveness of leveraging the relationships between all microphone pairs and optimizing all the SROs jointly. In all conditions, `Joint-ML-Prop` with 5-second-long signals was comparable to `Pair-ML-Prop` with 10-second-long signals. This result indicates that the proposed method can perform well with shorter signals, which is desirable to adopt to time-varying SROs [34] and unstationary environments [35].

5. Conclusion

In this paper, we propose a joint optimization method for all the SROs of multiple devices. The proposed method is based on the probabilistic model of the entire multichannel signal and estimates the SROs in a maximum likelihood manner. As the key idea, we maximize a tractable auxiliary function with respect to the SROs instead of the log-likelihood itself. Experimental results confirmed the effectiveness of the proposed joint optimization method compared with the pairwise methods. Future work includes an investigation of the robustness of the proposed method in real environments.

6. Acknowledgment

This work was supported by JSPS KAKENHI Grant Numbers JP20H00613 and JP21J21371, and JST CREST Grant Number JPMJCR19A3, Japan.

7. References

- [1] S. Makino, H. Sawada, and T. W. Lee, Eds., *Blind speech separation*. Springer, 2007.
- [2] H. Sawada, N. Ono, H. Kameoka, D. Kitamura, and H. Saruwatari, “A review of blind source separation methods: two converging routes to ILRMA originating from ICA and NMF,” *APSIPA Trans. Signal, Inf. Process.*, vol. 8, no. e12, pp. 1–14, May 2019.
- [3] C. Boeddeker, F. Rautenberg, and R. Haeb-Umbach, “A comparison and combination of unsupervised blind source separation techniques,” *arXiv*, Jan. 2021.
- [4] R. Scheibler, T. Komatsu, and M. Togami, “Multichannel separation and classification of sound events,” in *Proc. Eur. Signal Process. Conf. (EUSIPCO)*, Aug. 2021, pp. 1035–1039.
- [5] R. Scheibler and N. Ono, “Independent vector analysis with more microphones than sources,” in *Proc. IEEE Workshop Appl. Signal Process. Audio Acoust. (WASPAA)*, Oct. 2019, pp. 185–189.
- [6] R. Ikeshita, T. Nakatani, and S. Araki, “Overdetermined independent vector analysis,” in *Proc. IEEE Int. Conf. Acoust., Speech Signal Process. (ICASSP)*, May 2020, pp. 591–595.
- [7] A. Bertrand, “Applications and trends in wireless acoustic sensor networks: A signal processing perspective,” in *Proc. IEEE Symp. Commun., Veh. Technol. (SCVT)*, Nov. 2011, pp. 1–6.
- [8] A. Bertrand, S. Doclo, S. Gannot, N. Ono, and T. van Waterschoot, “Special issue on wireless acoustic sensor networks and ad hoc microphone arrays,” *Signal Process.*, vol. 107, no. C, pp. 1–3, Feb. 2015.
- [9] A. Brendel and W. Kellermann, “Distributed source localization in acoustic sensor networks using the coherent-to-diffuse power ratio,” *IEEE J. Sel. Top. Signal Process.*, vol. 13, no. 1, pp. 61–75, Mar. 2019.
- [10] S. Markovich-Golan, A. Bertrand, M. Moonen, and S. Gannot, “Optimal distributed minimum-variance beamforming approaches for speech enhancement in wireless acoustic sensor networks,” *Signal Process.*, vol. 107, pp. 4–20, Feb. 2015.
- [11] V. M. Tavakoli, J. R. Jensen, M. G. Christensen, and J. Benesty, “A framework for speech enhancement with ad hoc microphone arrays,” *IEEE/ACM Trans. Audio, Speech, Lang. Process.*, vol. 24, no. 6, pp. 1038–1051, Jun. 2016.
- [12] Y. Hioka and W. B. Kleijn, “Distributed blind source separation with an application to audio signals,” in *Proc. IEEE Int. Conf. Acoust., Speech Signal Process. (ICASSP)*, May 2011, pp. 233–236.
- [13] S. Miyabe, N. Ono, and S. Makino, “Blind compensation of interchannel sampling frequency mismatch with maximum likelihood estimation in STFT domain,” in *Proc. IEEE Int. Conf. Acoust., Speech Signal Process. (ICASSP)*, May 2013, pp. 674–678.
- [14] —, “Blind compensation of interchannel sampling frequency mismatch for ad hoc microphone array based on maximum likelihood estimation,” *Signal Process.*, vol. 107, pp. 185–196, Feb. 2015.
- [15] S. Markovich-Golan, S. Gannot, and I. Cohen, “Blind sampling rate offset estimation and compensation in wireless acoustic sensor networks with application to beamforming,” in *Proc. Int. Workshop Acoust. Signal Enhance. (IWAENC)*, Sep. 2012, pp. 1–4.
- [16] M. H. Bahari, A. Bertrand, and M. Moonen, “Blind sampling rate offset estimation for wireless acoustic sensor networks through weighted least-squares coherence drift estimation,” *IEEE/ACM Trans. Audio, Speech, Lang. Process.*, vol. 25, no. 3, pp. 674–686, Mar. 2017.
- [17] J. Schmalenstroer, J. Heymann, L. Drude, C. Boeddeker, and R. Haeb-Umbach, “Multi-stage coherence drift based sampling rate synchronization for acoustic beamforming,” in *Proc. IEEE Workshop Multimed. Signal Process. (MMSP)*, Oct. 2017, pp. 1–6.
- [18] L. Wang and S. Doclo, “Correlation maximization-based sampling rate offset estimation for distributed microphone arrays,” *IEEE/ACM Trans. Audio, Speech, Lang. Process.*, vol. 24, no. 3, pp. 571–582, Mar. 2016.
- [19] A. Chinaev, P. Thüne, and G.ENZNER, “A double-cross-correlation processor for blind sampling rate offset estimation in acoustic sensor networks,” in *Proc. IEEE Int. Conf. Acoust., Speech Signal Process. (ICASSP)*, May 2019, pp. 641–645.
- [20] —, “Double-cross-correlation processing for blind sampling-rate and time-offset estimation,” *IEEE/ACM Trans. Audio, Speech, Lang. Process.*, vol. 29, pp. 1881–1896, Apr. 2021.
- [21] K. Yamaoka, N. Ono, and Y. Wakabayashi, “Sampling frequency mismatch estimation by auxiliary-function-based iterative maximization of double-cross-correlation,” in *Proc. Eur. Signal Process. Conf. (EUSIPCO)*, Aug. 2021, pp. 1125–1129.
- [22] D. Cherkassky and S. Gannot, “Blind synchronization in wireless acoustic sensor networks,” *IEEE/ACM Trans. Audio, Speech, Lang. Process.*, vol. 25, no. 3, pp. 651–661, Mar. 2017.
- [23] J. Zhang and P. Wu, “Joint sampling synchronization and source localization for wireless acoustic sensor networks,” *IEEE Commun. Lett.*, vol. 24, no. 5, pp. 1020–1023, May 2020.
- [24] R. Wang, Z. Chen, and F. Yin, “Active sampling rate calibration method for acoustic sensor networks,” *IEEE/ACM Trans. Audio, Speech, Lang. Process.*, vol. 28, pp. 3095–3107, Nov. 2020.
- [25] J. Schmalenstroer and R. Haeb-Umbach, “Efficient sampling rate offset compensation—an Overlap-Save based approach,” in *Proc. Eur. Signal Process. Conf. (EUSIPCO)*, Sep. 2017, pp. 499–503.
- [26] D. R. Hunter and K. Lange, “A tutorial on MM algorithms,” *Am. Stat.*, vol. 58, no. 1, pp. 30–37, Jan. 2004.
- [27] Y. Sun, P. Babu, and D. P. Palomar, “Majorization-minimization algorithms in signal processing, communications, and machine learning,” *IEEE Trans. Signal Process.*, vol. 65, no. 3, pp. 794–816, Feb. 2017.
- [28] N. Ono, “Stable and fast update rules for independent vector analysis based on auxiliary function technique,” in *Proc. IEEE Workshop Appl. Signal Process. Audio Acoust. (WASPAA)*, Oct. 2011, pp. 189–192.
- [29] K. Yamaoka, R. Scheibler, N. Ono, and Y. Wakabayashi, “Sub-sample time delay estimation via auxiliary-function-based iterative updates,” in *Proc. IEEE Workshop Appl. Signal Process. Audio Acoust. (WASPAA)*, Oct. 2019, pp. 130–134.
- [30] K. Yamaoka, N. Ono, and Y. Wakabayashi, “Estimation of consistent time delays in subsample via auxiliary-function-based iterative updates,” *arXiv:2203.09723*, Mar. 2022.
- [31] S. Boyd and L. Vandenberghe, *Introduction to Applied Linear Algebra: Vectors, Matrices, and Least Squares*. Cambridge University Press., 2018.
- [32] J. Lorenzo-Trueba, J. Yamagishi, T. Toda, D. Saito, F. Villavicencio, T. Kinnunen, and Z. Ling, “The voice conversion challenge 2018: Promoting development of parallel and nonparallel methods,” in *Odyssey*, Jun. 2018, pp. 195–202.
- [33] R. Scheibler, E. Bezzam, and I. Dokmanić, “Pyroomacoustics: A Python package for audio room simulation and array processing algorithms,” in *IEEE Int. Conf. Acoust., Speech Signal Process. (ICASSP)*, Apr. 2018, pp. 351–355.
- [34] T. Gburrek, J. Schmalenstroer, and R. Haeb-Umbach, “On synchronization of wireless acoustic sensor networks in the presence of time-varying sampling rate offsets and speaker changes,” *arXiv:2110.12820*, Oct. 2021.
- [35] S. Araki, N. Ono, K. Kinoshita, and M. Delcroix, “Estimation of sampling frequency mismatch between distributed asynchronous microphones under existence of source movements with stationary time periods detection,” in *Proc. IEEE Int. Conf. Acoust., Speech Signal Process. (ICASSP)*, May 2019, pp. 785–789.