



Multichannel Audio Modeling with Elliptically Stable Tensor Decomposition

Mathieu Fontaine, Fabian-Robert Stöter, Antoine Liutkus, Umut Simsekli,
Romain Serizel, Roland Badeau

► To cite this version:

Mathieu Fontaine, Fabian-Robert Stöter, Antoine Liutkus, Umut Simsekli, Romain Serizel, et al.. Multichannel Audio Modeling with Elliptically Stable Tensor Decomposition. 14th International Conference on Latent Variable Analysis and Signal Separation (LVA/ICA 2018), Jul 2018, Surrey, United Kingdom. pp.13-23, 10.1007/978-3-319-93764-9_2 . lirmm-01766795

HAL Id: lirmm-01766795

<https://hal-lirmm.ccsd.cnrs.fr/lirmm-01766795>

Submitted on 14 Apr 2018

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

Multichannel Audio Modeling with Elliptically Stable Tensor Decomposition

Mathieu Fontaine¹, Fabian-Robert Stöter², Antoine Liutkus², Umut Şimşekli³,
Romain Serizel¹, Roland Badeau³

¹Université de Lorraine, CNRS, Inria, LORIA, F-54000 Nancy, France.

² Inria and LIRMM, Montpellier, France.

³ LTCI, Télécom ParisTech, Université Paris-Saclay, Paris, France.

Abstract. This paper introduces a new method for multichannel speech enhancement based on a versatile modeling of the residual noise spectrogram. Such a model has already been presented before in the single channel case where the noise component is assumed to follow an alpha-stable distribution for each time-frequency bin, whereas the speech spectrogram, supposed to be more regular, is modeled as Gaussian. In this paper, we describe a multichannel extension of this model, as well as a Monte Carlo Expectation - Maximisation algorithm for parameter estimation. In particular, a multichannel extension of the Itakura-Saito nonnegative matrix factorization is exploited to estimate the spectral parameters for speech, and a Metropolis-Hastings algorithm is proposed to estimate the noise contribution. We evaluate the proposed method in a challenging multichannel denoising application and compare it to other state-of-the-art algorithms.

1 Introduction

In many contexts, speech denoising is studied and applied in order to obtain, among other things, a comfortable listening or broadcast of a talk [2], by exploiting the observed noisy signal, obtained by several microphones. From an audio source separation perspective, this denoising is achieved through a probabilistic model, where the observed signal is divided into two latent sources: a noise component and a target source.

Both speech and noise components are usually considered in the *time frequency* (TF) domain where all TF-bins are supposed to be independent and follow a Gaussian law [5, 13]. A common approach to speech enhancement is the spectral subtraction method [6]. Its principle is to estimate an a priori *signal to noise ratio* (SNR) in order to infer a *short-time spectral amplitude* (STSA) estimator of the noise which will be subtracted from the STSA of the observations. Another popular trend is to decompose the *power spectral densities* (PSD) of sources into a product of two matrices. The *non-negative matrix factorization* (NMF) model assumes that the PSDs admit low-rank structures and it performs well in denoising .

To the best of our knowledge, NMF models for multichannel speech enhancement have been proposed only in a Gaussian probabilistic context, whereas a non-Gaussian approach could offer a more flexible model for noise and speech. Moreover, a good initialization in a Gaussian NMF model is crucial to avoid staying stuck in a local minimum [3]. Many studies in the single-channel case have shown a better robustness to initialization when the signal is modeled in the TF domain with as heavy tail distribution [22, 19].

Among this type of distributions, α -stable distributions preserve interesting properties satisfied by Gaussian laws, and they can model distributions ranging from light tails as in the *Gaussian case* to heavy tails as in the *Cauchy case*. Indeed, α -stable distributions are the only ones which admit a central limit theorem and stability by summation [16]. Various studies have been carried out on audio modeling using alpha-stable processes [19, 12]. Especially in the TF domain, a generalization of wide-sense stationary (WSS) processes [13] has been established in the α -stable case [12] and applied to noise reduction [8]. More precisely, in [20] it was considered that the target source is Gaussian and the residual noise is α -stable, in order to get a greater flexibility on noise modeling.

This paper introduces a generalization of [20] to the multichannel case. The goal is to develop a Gaussian NMF model for speech that assumes a low-rank structure for speech covariances [5], while the noise part is taken as an α -stable process. Parameters are estimated through a combination of the multichannel extension of Itakura Saito NMF (IS-NMF) [17] for speech and a Markov Chain Monte Carlo (MCMC) strategy for estimating the noise part. The proposed method is evaluated for multichannel denoising, and compared to other state-of-the-art algorithms.

2 Probabilistic and Filtering models

2.1 Mixture model

Let $\mathbf{x} \in \mathbb{C}^{F \times T \times K}$ be the observed data in the short-time Fourier transform (STFT) domain where F, T and K denote the number of frequency bands, time frames and microphones, respectively. The observation \mathbf{x} will be assumed to be the sum of two latent audio sources represented by two tensors: the first one is written $\mathbf{y} \in \mathbb{C}^{F \times T \times K}$ and accounts for the *speech signal*. The second one is written $\mathbf{r} \in \mathbb{C}^{F \times T \times K}$ and called the *residual component*. We have:

$$\mathbf{x}_{ft} = \mathbf{y}_{ft} + \mathbf{r}_{ft}, \quad (1)$$

where each term belongs to \mathbb{C}^K . The main goal in this paper is to estimate the tensors \mathbf{y} and \mathbf{r} knowing \mathbf{x} , by using a probabilistic model described below.

2.2 Source model

At short time scales, the speech signal may be assumed stationary and does not feature strong impulsiveness. This motivates modeling it as a locally stationary Gaussian process [13]. Furthermore, we also assume that the different channels

for \mathbf{y}_{ft} are correlated, accounting for the *spatial* characteristics of the signal. Consequently, we assume that each \mathbf{y}_{ft} is an isotropic complex Gaussian vector¹ of mean $\mathbf{0}$ and covariance matrix $\mathbf{C}_{ft}^{\mathbf{y}} \triangleq \mathbf{R}_f v_{f,t}$, where the *spatial covariance matrix* $\mathbf{R}_f \in \mathbb{C}^{K \times K}$ encodes the time-invariant correlations of speech in the different channels and v_{ft} is the PSD of the speech signal [5]. To exploit the redundancy of speech, we further decompose v_{ft} through NMF and obtain:

$$\forall f, t \quad \mathbf{y}_{ft} \sim \mathcal{N}_c \left(\mathbf{y}_{ft}; 0, \mathbf{R}_f v_{ft} \triangleq \mathbf{R}_f \sum_{l=1}^L w_{fl} h_{lt} \right). \quad (2)$$

where \triangleq means “equals by definition” and $\mathbf{W} \in \mathbb{R}_+^{F \times L}$, $\mathbf{H} \in \mathbb{R}_+^{L \times T}$ are matrices which respectively contain all non-negative scalars w_{fl} and h_{lt} . While \mathbf{W} is understood as L spectral bases, \mathbf{H} stands for their activations over time. To make notations simpler, let $\Theta \triangleq \{\mathbf{W}, \mathbf{H}, \mathbf{R}\}$ be the parameters to be estimated with $\mathbf{R} \triangleq \{\mathbf{R}_f\}_f$. Note that the decomposition of v_{ft} is not unique: it is defined up to multiplicative constant.

In contrast to the speech signal, the model of the residual component should allow for outliers and impulsiveness. To do so, the residual part is modeled by an heavy-tailed distribution in the time domain. Recent works proposed a stationary model called α -harmonizable process with $\alpha \in (0, 2]$ in the single-channel case. It is shown in [16, 12] that such a model is equivalent to assuming that the signal at every time-frequency bin f, t follows a complex isotropic symmetric α -stable distribution. With the aim of extending the previous model to a multichannel one, we take all \mathbf{r}_{ft} as distributed with respect to an *elliptically contoured multivariate stable distribution* (ECMS, denoted $\mathcal{E}\alpha S$) and independent of one another. These distributions, which are a particular case of α -stable distributions, have the particularity of requiring only two parameters [16, 11]:

1. A *characteristic exponent* $\alpha \in (0, 2]$: the smaller α , the heavier the tails of the distribution.
2. A positive definite Hermitian *scatter matrix* in $\mathbb{C}^{K \times K}$.

In this article, the scatter matrices for all \mathbf{r}_{ft} are taken equal to $\sigma_f \mathbf{I}_K$, where $\mathbf{I}_K \in \mathbb{R}^{K \times K}$ is the identity matrix and $\sigma_f > 0$ is a positive scalar that does not depend on time. We have:

$$\forall f, t \quad \mathbf{r}_{ft} \sim \mathcal{E}\alpha S^K (\sigma_f \mathbf{I}_K). \quad (3)$$

2.3 Filtering model

As mentioned in subsection 2.1, we aim to reconstruct the sources \mathbf{y} and \mathbf{r} from the observed data \mathbf{x} . From a signal processing point of view, when parameters $\sigma, \mathbf{W}, \mathbf{H}, \mathbf{R}$ are known, one would like to compute the Minimum Mean Squared

¹ The probability density function (PDF) of an isotropic complex Gaussian vector is $\mathcal{N}_C(\mathbf{x}|\mu, \mathbf{C}) = \frac{1}{\pi^K \det \mathbf{C}} \exp(-(\mathbf{x} - \mu)^* \mathbf{C}^{-1} (\mathbf{x} - \mu))$.

Error (MMSE) estimates of both sources. In our probabilistic context, these can be expressed as the posteriori expectations $\mathbb{E}(\mathbf{y}_{ft}|\mathbf{x}_{ft}, \boldsymbol{\Theta}, \boldsymbol{\sigma})$.

To compute such estimates, a property specific to ECMS distributions can be exploited to represent \mathbf{r} as a complex normal distribution \mathcal{N}_c of dimension K , whose variance is randomly multiplied by a positive random *impulse variable* ϕ_{ft} distributed as $\mathcal{P}_{\frac{\alpha}{2}}S\left(2\cos\left(\frac{\pi\alpha}{4}\right)^{2/\alpha}\right)$, where $\mathcal{P}_{\frac{\alpha}{2}}S$ is the *positive $\alpha/2$ -stable distribution* (see [19] for more details):

$$\forall f, t \quad \mathbf{r}_{ft}|\phi_{ft} \sim \mathcal{N}_c(\mathbf{r}_{ft}; 0, \phi_{ft}\sigma_f\mathbf{I}_K), \quad (4)$$

If we assume for now that $\boldsymbol{\Phi} \triangleq \{\phi_{ft}\}_{f,t}$ are known in (4), we get the distribution of the mixture as:

$$\forall f, t \quad \mathbf{x}_{ft}|\phi_{ft} \sim \mathcal{N}_c(\mathbf{x}_{ft}; 0, \mathbf{C}_{ft}^{\mathbf{x}|\phi}), \quad (5)$$

where $\mathbf{C}_{ft}^{\mathbf{x}|\phi} \triangleq \mathbf{R}_f \sum_{l=1}^L w_{fl}h_{lt} + \phi_{ft}\sigma_f\mathbf{I}_K$. This in turns allows to build a multichannel Wiener filter as [2]:

$$\mathbb{E}(\mathbf{y}_{ft}|\mathbf{x}_{ft}, \boldsymbol{\Phi}, \boldsymbol{\Theta}, \boldsymbol{\sigma}) = \mathbf{C}_{ft}^{\mathbf{y}} \left(\mathbf{C}_{ft}^{\mathbf{x}|\phi}\right)^{-1} \mathbf{x}_{ft}, \quad (6)$$

with \cdot^{-1} standing for matrix inversion.

Now, the strategy we adopt here is to marginalize this expression over $\boldsymbol{\Phi} | x$, to get:

$$\hat{\mathbf{y}}_{ft} = \mathbb{E}_{\boldsymbol{\Phi}|x} [\mathbb{E}[\mathbf{y}_{ft}|\mathbf{x}_{ft}, \boldsymbol{\Phi}, \boldsymbol{\Theta}, \boldsymbol{\sigma}]] = \mathbf{G}_{ft}\mathbf{x}_{ft},$$

where

$$\mathbf{G}_{ft} \triangleq \mathbf{C}_{ft}^{\mathbf{y}} \boldsymbol{\Xi}_{ft} \quad (7)$$

is the marginal Wiener filter, and $\boldsymbol{\Xi}_{ft} \triangleq \mathbb{E}_{\boldsymbol{\Phi}|x} \left[\left(\mathbf{C}_{ft}^{\mathbf{x}|\phi}\right)^{-1} \right]$ is the average inverse mixture covariance matrix. We will explain how to compute $\boldsymbol{\Xi}$ later in section 3.3.

3 Parameter Estimation

3.1 Expectation-Maximization (EM) algorithm

Assuming that the observations \mathbf{x} and the impulse variable ϕ are known, we first aim to estimate the parameters $\boldsymbol{\Theta}$. We choose a maximum likelihood estimator in order to get the most likely source NMF parameters \mathbf{W}, \mathbf{H} :

$$(\mathbf{W}^*, \mathbf{H}^*, \mathbf{R}^*) = \arg \max_{\mathbf{W}, \mathbf{H}, \mathbf{R}} \log \mathbb{P}(\mathbf{x}, \boldsymbol{\Phi} | \boldsymbol{\Theta}, \boldsymbol{\sigma}), \quad (8)$$

where $\boldsymbol{\Phi}$ is a latent variable and $\log \mathbb{P}(\mathbf{x}, \boldsymbol{\Phi} | \boldsymbol{\Theta}, \boldsymbol{\sigma})$ is the log-likelihood. As in [20], we propose an EM algorithm. This method aims to minimize an upper-bound of $\mathcal{L}_n(\mathbf{W}, \mathbf{H}, \mathbf{R}) = -\log \mathbb{P}(\mathbf{x}, \boldsymbol{\Phi} | \boldsymbol{\Theta}, \boldsymbol{\sigma})$. This approach is summarized in the following two steps:

$$\text{E-Step: } \quad \mathcal{Q}_n(\mathbf{W}, \mathbf{H}, \mathbf{R}) = -\mathbb{E}_{\Phi|\mathbf{x}, \mathbf{W}^{(n-1)}, \mathbf{H}^{(n-1)}} [\mathcal{L}_n(\mathbf{W}, \mathbf{H}, \mathbf{R})], \quad (9)$$

$$\text{M-Step: } \quad \left(\mathbf{W}^{(n)}, \mathbf{H}^{(n)}, \mathbf{R}^{(n)} \right) = \arg \max_{\mathbf{W}, \mathbf{H}, \mathbf{R}} \mathcal{Q}_n(\mathbf{W}, \mathbf{H}, \mathbf{R}). \quad (10)$$

E-Step: We first introduce a positive function that upper-bounds the negative log-likelihood $\mathcal{L}_n(\mathbf{W}, \mathbf{H}, \mathbf{R})$, which is equal to [17]:

$$\mathcal{L}_n(\mathbf{W}, \mathbf{H}, \mathbf{R}) = \sum_{f,t} \left[\text{tr} \left(\tilde{\mathbf{X}}_{ft} \left(\mathbf{C}_{ft}^{\mathbf{x}|\phi} \right)^{-1} \right) + \log \det \mathbf{C}_{ft}^{\mathbf{x}|\phi} \right] \quad (11)$$

where $\tilde{\mathbf{X}}_{ft} \triangleq \mathbf{x}_{ft} \mathbf{x}_{ft}^*$ and \cdot^* stands for the Hermitian transposition. A positive auxiliary function $\mathcal{L}_n^+(\mathbf{W}, \mathbf{H}, \mathbf{R}, \mathbf{U}, \mathbf{V}) = \sum_{f,t} \left[\sum_l \frac{\text{tr} \left(\tilde{\mathbf{X}}_{ft} \mathbf{U}_{lft} \left(\mathbf{C}_{lft}^{\mathbf{x}|\phi} \right)^{-1} \mathbf{U}_{lft} \right)}{w_{fl} h_{lt}} + \frac{\text{tr} \left(\tilde{\mathbf{X}}_{ft} \mathbf{U}_{lft}^2 \right)}{\sigma_f \phi_{ft}} + \log \det \mathbf{V}_{ft} + \frac{\det \mathbf{C}_{ft}^{\mathbf{x}|\phi} - \det \mathbf{V}_{ft}}{\det \mathbf{V}_{ft}} \right]$ which satisfies:

$$\mathcal{L}_n^+(\mathbf{W}, \mathbf{H}, \mathbf{R}, \mathbf{U}, \mathbf{V}) \geq \mathcal{L}_n(\mathbf{W}, \mathbf{H}, \mathbf{R}) \quad (12)$$

is introduced in [17]. Using (12) and the definition of \mathcal{Q}_n in (9), we obtain:

$$\mathbb{E}_{\Phi|\mathbf{x}} \mathcal{L}_n(\cdot) \leq \mathbb{E}_{\Phi|\mathbf{x}} \mathcal{L}_n^+(\cdot) \triangleq \mathcal{Q}_n^+(\cdot) \quad (13)$$

with:

$$\begin{aligned} \mathcal{Q}_n^+(\mathbf{W}, \mathbf{H}, \mathbf{R}, \mathbf{U}, \mathbf{V}) &= \sum_{f,t} \left[\sum_l \frac{\mathbb{E}_{\Phi|\mathbf{x}} \left(\text{tr} \left[\tilde{\mathbf{X}}_{ft} \mathbf{U}_{lft} \left(\mathbf{C}_{lft}^{\mathbf{x}|\phi} \right)^{-1} \mathbf{U}_{lft} \right] \right)}{w_{fl} h_{lt}} \right. \\ &\quad \left. + \mathbb{E}_{\Phi|\mathbf{x}} \left(\text{tr} \left[\tilde{\mathbf{X}}_{ft} \mathbf{U}_{lft}^2 \right] \right) \sigma_f^{-1} \phi_{ft}^{-1} + \mathbb{E}_{\Phi|\mathbf{x}} \left(\log \det \mathbf{V}_{ft} + \det \left(\mathbf{V}_{ft}^{-1} \mathbf{C}_{lft}^{\mathbf{x}|\phi} \right) - 1 \right) \right] \quad (14) \end{aligned}$$

The form in (14) admits partial derivatives that will be useful as part of a multiplicative update [7] in the M-Step.

M-Step: Solving the M-Step in (10) is equivalent to zeroing the partial derivatives $\frac{\partial \mathcal{Q}_n^+}{\partial w_{fl}}$ and $\frac{\partial \mathcal{Q}_n^+}{\partial h_{lt}}$ and to set \mathbf{U}, \mathbf{V} such that the equality in (13) is verified. A multiplicative update approach yields:

$$w_{fl} \leftarrow w_{fl} \sqrt{\frac{\sum_t h_{lt} \text{tr}(\mathbf{R}_f \mathbf{P}_{ft})}{\sum_t h_{lt} \text{tr}(\mathbf{R}_f \mathbf{\Xi}_{ft})}}; \quad h_{lt} \leftarrow h_{lt} \sqrt{\frac{\sum_f w_{fl} \text{tr}(\mathbf{R}_f \mathbf{P}_{ft})}{\sum_f w_{fl} \text{tr}(\mathbf{R}_f \mathbf{\Xi}_{ft})}} \quad (15)$$

where the quantity $\mathbf{\Xi}_{ft} = \mathbb{E}_{\Phi|\mathbf{x}} \left[\left(\mathbf{C}_{ft}^{\mathbf{x}|\phi_i} \right)^{-1} \right]$ has been used above in (7) and $\mathbf{P}_{ft} = \mathbb{E}_{\Phi|\mathbf{x}} \left[\left(\mathbf{C}_{ft}^{\mathbf{x}|\phi_i} \right)^{-1} \tilde{\mathbf{X}}_{ft} \left(\mathbf{C}_{ft}^{\mathbf{x}|\phi_i} \right)^{-1} \right]$. We will explain how to compute these expectations in subsection 3.3.

3.2 Estimation of spatial covariance matrices and noise gains σ

We update the spatial covariance matrix \mathbf{R} in the M-step as in [5], further using the trick proposed in [14] to use a weighted update, resulting in:

$$\mathbf{R}_f \leftarrow \left(\sum_t v_{ft} \right)^{-1} \times \sum_t \left(\mathbf{C}_{ft}^{\mathbf{y}\mathbf{y}^*|\mathbf{x}} \right), \quad (16)$$

where: $\mathbf{C}_{ft}^{\mathbf{y}\mathbf{y}^*|\mathbf{x}} \triangleq \mathbf{G}_{ft} \tilde{\mathbf{X}}_{ft} \mathbf{G}_{ft} + \mathbf{C}_{ft}^{\mathbf{y}} - \mathbf{G}_{ft} \mathbf{C}_{ft}^{\mathbf{y}}$ is the total posterior variance for the speech source.

Concerning the estimation of the noise gain σ in (3), we exploit a result in [4] that if $z \sim \mathcal{E}\alpha S(\sigma)$, then $\mathbb{E}[\|z\|^p] \propto \sigma$, for $p < \alpha$, with α standing for proportionality. The strategy we adopt is to apply this estimation only once at the beginning of the algorithm to the mixture itself, by taking a robust estimation like the median \mathbb{M} instead of the average, to account for the fact that not all TF bins pertain to the noise, but that a small portion also pertain to speech. We thus pick $p = \alpha/2$ and take:

$$\sigma_f \leftarrow \mathbb{M} \left(\left\| \sum_t \mathbf{x}(f, t) \right\|^{\alpha/2} \right)^2. \quad (17)$$

3.3 Expectation estimation using Metropolis-Hastings algorithm

We still have to calculate the expectations $\bar{\boldsymbol{\Xi}}_{ft}$ and \mathbf{P}_{ft} . Unfortunately, they cannot be calculated analytically. To address this issue, we set up a Markov Chain Monte Carlo (MCMC) algorithm in order to approximate the expectations for each iteration. We are focusing on the Metropolis-Hastings algorithm through an empirical estimation of $\bar{\boldsymbol{\Xi}}_{ft}$ and \mathbf{P}_{ft} as follows:

$$\bar{\boldsymbol{\Xi}}_{ft} \simeq \frac{1}{I} \sum_{i=1}^I \left(\mathbf{C}_{ft}^{\mathbf{x}|\varphi_i} \right)^{-1}; \quad \mathbf{P}_{ft} \simeq \frac{1}{I} \sum_{i=1}^I \left(\left(\mathbf{C}_{ft}^{\mathbf{x}|\varphi_i} \right)^{-1} \tilde{\mathbf{X}}_{ft} \left(\mathbf{C}_{ft}^{\mathbf{x}|\varphi_i} \right)^{-1} \right) \quad (18)$$

with $\left(\mathbf{C}_{ft}^{\mathbf{x}|\varphi_i} \right)^{-1} = [\sum_l (\mathbf{R}_{fl} w_{fl} h_{lt}) + \varphi_{ft, i} \sigma_f \mathbf{I}_k]^{-1}$ and $\varphi_{ft, i}$ are sampled as follows:

First Step (Sampling process): Generate a sampling via the prior distribution $\varphi'_{ft} \sim \mathcal{P}_{\frac{\alpha}{2}} S \left(2 \cos \left(\frac{\pi\alpha}{4} \right)^{2/\alpha} \right)$.

Second Step (Acceptance):

- Draw $u \sim \mathcal{U}([0, 1])$ where \mathcal{U} denotes the uniform distribution.
- Compute the following acceptance probability:

$$\text{acc}(\varphi_{ft} \rightarrow \varphi'_{ft}) = \min \left(1, \frac{\mathcal{N}_c \left(\mathbf{x}_{ft}; 0, \varphi'_{ft} \sigma_f \mathbf{I}_K + \mathbf{C}_{ft}^{\mathbf{y}} \right)}{\mathcal{N}_c \left(\mathbf{x}_{ft}; 0, \varphi_{ft} \sigma_f \mathbf{I}_K + \mathbf{C}_{ft}^{\mathbf{y}} \right)} \right)$$

- Test the acceptance:
 - if $u < \text{acc}(\varphi_{ft,i-1} \rightarrow \varphi'_{ft})$, then $\varphi_{ft,i} = \varphi'_{ft}$ (acceptance)
 - otherwise, $\varphi_{ft,i} = \varphi_{ft,i-1}$ (rejection)

4 Single-Channel Speech Signal Reconstruction

Let $\hat{\mathbf{y}}$ be the multichannel signal obtained after Wiener filtering (7). In the context of speech enhancement, the desired speech is rather a single-channel signal, that we write $\hat{\mathbf{s}} \in \mathbb{C}^{F \times T}$. In this study, we take $\hat{\mathbf{s}}$ as a linear combination of $\hat{\mathbf{y}}$ with a time-invariant *beamformer* $\mathbf{B}_f \in \mathbb{C}^K$ [21]:

$$\hat{\mathbf{s}}_{ft} \triangleq \mathbf{B}_f^* \hat{\mathbf{y}}_{ft},$$

Where $.*$ denotes the Hermitian transposition. There are many ways to devise the beamformer \mathbf{B}_f . In this study, we choose to maximize the energy of $\mathbf{B}_f^* \mathbf{y}_{ft} | \mathbf{x}$:

$$\begin{aligned} \frac{1}{T} \sum_t \mathbb{E} \left(|\mathbf{B}_f^* \mathbf{y}_{ft}|^2 | \mathbf{x}_{ft} \right) &= \mathbf{B}_f^* \mathbb{E} \left(\mathbf{y}_{ft} \mathbf{y}_{ft}^* | \mathbf{x} \right) \mathbf{B}_f. \\ &= \mathbf{B}_f^* \frac{1}{T} \sum_t \left(\mathbf{C}_{ft}^{\mathbf{y} \mathbf{y}^* | \mathbf{x}} \right) \mathbf{B}_f. \end{aligned}$$

This is solved by taking \mathbf{B}_f as the eigenvector associated to the largest eigenvalue of the Hermitian matrix $\frac{1}{T} \sum_t \left(\mathbf{C}_{ft}^{\mathbf{y} \mathbf{y}^* | \mathbf{x}} \right)$ [5].

5 Evaluation

We investigate both the quality of speech enhancement and the audio source separation performance. Our proposed approach will be compared to two baseline methods:

- ARC** The proposed Alpha Residual component. We take $N = 10$ iterations for the EM and pick $\alpha = 1.9$.
- MWF** The classic multi-channel Wiener filter [5] which assumes Gaussianity for both noise and speech.
- GEVD** The generalized eigenvalue decomposition [18] is based on a low-rank approximation of the autocorrelation matrix of the speech signal.

5.1 Experimental setup

The corpus for evaluation is made up of mono speech excerpts from Librispeech [15] and three different environmental noises taken from Aurora [10]: babble noise, restaurant and train. A groundtruth voice activity detection (VAD) is used on all three methods.

Mixtures were obtained for two 15 cm separated microphones, with the Roomsimove simulator with room dimensions of $5 \times 4 \times 3$ meters and $\text{RT60} = 0$ ms and 500 ms. The sources are taken 1 m from the microphones, with different SNR values of $-5, 0, 5, 10$ dB and an angular distance of 30° or 90° . This results in 48 experiments.

5.2 Performance measures

For the evaluation, two scores will be measured: the first one is a speech intelligibility weighted spectral distortion (SIW-SD) measure and the second one is a speech intelligibility-weighted SNR (SIW-SNR) [9].

The SIW-SD measure is defined as:

$$\text{SIW} - \text{SD} = \sum_i I_i \text{SD}_i \quad (19)$$

where I_i is the band importance function [1] and SD_i the average SD (in dB) in the i -th one third octave band,

$$\text{SD}_i = \frac{1}{(2^{1/6} - 2^{-1/6})f_i^c} \int_{2^{-1/6}f_i^c}^{2^{1/6}f_i^c} |10 \log_{10} G^y(f)| df \quad (20)$$

with center frequencies f_i^c and $G^y(f)$ is given by:

$$G^y(f) = \frac{P_{\mathbf{y}}(f)}{P_{\hat{\mathbf{y}}}(f)} \quad (21)$$

where $P_{\mathbf{y}}(f)$ and $P_{\hat{\mathbf{y}}}(f)$ are the power, for the frequency f , of the speech component of the input signal \mathbf{y} and the speech component output signal $\hat{\mathbf{y}}$, respectively.

The SIW-SNR [9] is used here to compute the *SIW-SNR improvement* which is defined as

$$\Delta \text{SNR}_{\text{intellig}} = \sum_i I_i (\text{SNR}_{i,\text{out}} - \text{SNR}_{i,\text{in}}) \quad (22)$$

where $\text{SNR}_{i,\text{out}}$ and $\text{SNR}_{i,\text{in}}$ represent the output SNR of the noise reduction filter and the SNR of the signal in the first microphone in the i^{th} band, respectively.

5.3 Results

Results are displayed on Figure 1 and present the SIW-SNR and SIW-SD scores averaged over noise types and spatial scenarios, against the input SNR.

We first investigate the impact of reverberation. While we see that ARC is outperformed by other methods in the anechoic case, we see it is much less sensitive to reverberation and becomes competitive compared to the other algorithms in terms of SIW-SD at low input SNR.

6 Conclusion

We proposed a new method ARC for denoising that is more robust to reverberation than competing approaches, although less effective in the anechoic case. It is based on modeling the speech signal as a Gaussian process and noise as an α -stable sub-Gaussian process. Interestingly, that approach can be combined with existing methods, which could be an interesting avenue for future work.

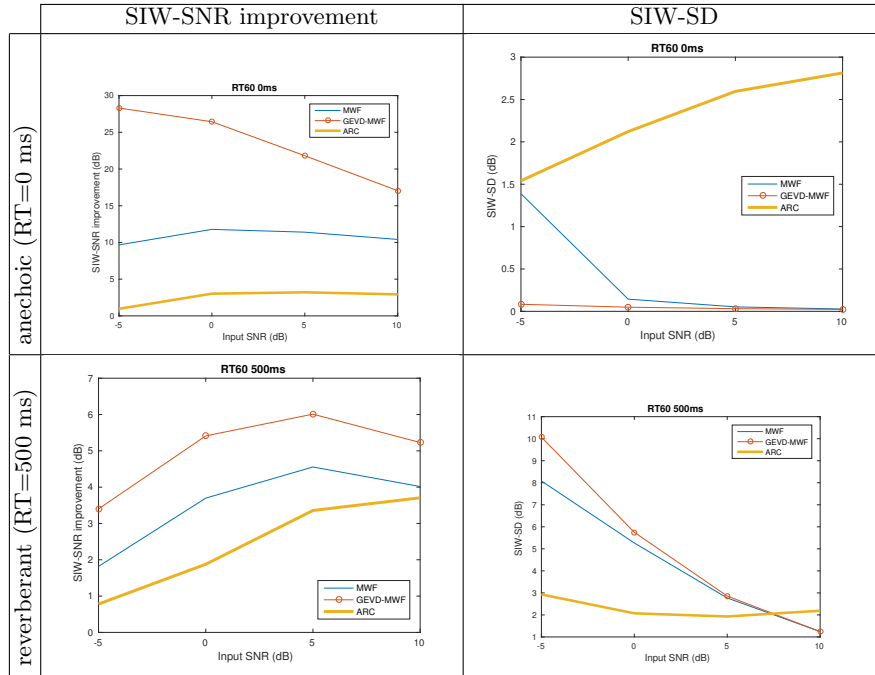


Fig. 1. SIW (left, higher is better) and SNR & SIW-SD (right, lower is better) for: (top) an anechoic scenario and (bottom) a reverberant room.

Acknowledgments. This work was partly supported by the research programme KAMoulox (ANR-15-CE38-0003-01), EDiSon3D (ANR-13-CORD-0008-01), FBIMATRIX (ANR-16-CE23-0014) funded by ANR, the French State agency for research.

References

1. ANSI: S3. 5-1997, methods for the calculation of the speech intelligibility index. New York: American National Standards Institute 19, 90–119 (1997)
2. Van den Bogaert, T., Doclo, S., Wouters, J., Moonen, M.: Speech enhancement with multichannel Wiener filter techniques in multimicrophone binaural hearing aids. *The Journal of the Acoustical Society of America* 125(1), 360–371 (2009)
3. Boutsidis, C., Gallopoulos, E.: SVD based initialization: A head start for nonnegative matrix factorization. *Pattern Recognition* 41(4), 1350–1362 (2008)
4. Cambanis, S., Keener, R., Simons, G.: On α -symmetric multivariate distributions. *Journal of Multivariate Analysis* 13(2), 213–233 (1983)
5. Duong, N., Vincent, E., Gribonval, R.: Under-determined reverberant audio source separation using a full-rank spatial covariance model. *IEEE/ACM Trans. Audio, Speech, Language Process.* 18(7), 1830–1840 (sept 2010)
6. Ephraim, Y., Malah, D.: Speech enhancement using a minimum-mean square error short-time spectral amplitude estimator. *IEEE Trans. Acoust., Speech, Signal Process.* 32(6), 1109–1121 (1984)

7. Févotte, C., Idier, J.: Algorithms for nonnegative matrix factorization with the β -divergence. *Neural computation* 23(9), 2421–2456 (2011)
8. Fontaine, M., Liutkus, A., Girin, L., Badeau, R.: Parameterized Wiener filtering for single-channel denoising. In: *Workshop on Applications of Signal Processing to Audio and Acoustics (WASPAA)* (2017)
9. Greenberg, J., Peterson, P., Zurek, P.: Intelligibility-weighted measures of speech-to-interference ratio and speech system performance. *The Journal of the Acoustical Society of America* 94(5), 3009–3010 (1993)
10. Hirsch, H., Pearce, D.: The aurora experimental framework for the performance evaluation of speech recognition systems under noisy conditions. In: *ASR2000-Automatic Speech Recognition: Challenges for the new Millenium ISCA Tutorial and Research Workshop (ITRW)* (2000)
11. Leglaive, S., Simsekli, U., Liutkus, A., Badeau, R., Richard, G.: Alpha-stable multichannel audio source separation. In: *42nd International Conference on Acoustics, Speech and Signal Processing (ICASSP)* (2017)
12. Liutkus, A., Badeau, R.: Generalized Wiener filtering with fractional power spectrograms. In: *40th International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. pp. 266–270. IEEE (2015)
13. Liutkus, A., Badeau, R., Richard, G.: Gaussian processes for underdetermined source separation. *IEEE Trans. Signal Process.* 59(7), 3155–3167 (2011)
14. Nugraha, A.A., Liutkus, A., Vincent, E.: Multichannel music separation with deep neural networks. In: *24th European Signal Processing Conference (EUSIPCO)*, 2016. pp. 1748–1752. IEEE (2016)
15. Panayotov, V., Chen, G., Povey, D., Khudanpur, S.: Librispeech: an ASR corpus based on public domain audio books. In: *2015 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. pp. 5206–5210. IEEE (2015)
16. Samoradnitsky, G., Taqqu, M.: *Stable non-Gaussian random processes: stochastic models with infinite variance*, vol. 1. CRC Press (1994)
17. Sawada, H., Kameoka, H., Araki, S., Ueda, N.: Efficient algorithms for multichannel extensions of Itakura-Saito nonnegative matrix factorization. In: *37th International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. pp. 261–264. IEEE (2012)
18. Serizel, R., Moonen, M., Van Dijk, B., Wouters, J.: Low-rank approximation based multichannel Wiener filter algorithms for noise reduction with application in cochlear implants. *IEEE/ACM Transactions on Audio, Speech, and Language Processing* 22(4), 785–799 (2014)
19. Şimşekli, U., Liutkus, A., Cemgil, A.: Alpha-stable matrix factorization. *IEEE Signal Process. Lett.* 22(12), 2289–2293 (2015)
20. Şimşekli, U., et al.: Alpha-stable low-rank plus residual decomposition for speech enhancement. In: *2018 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE (2018)
21. Van Veen, B.D., Buckley, K.M.: Beamforming: A versatile approach to spatial filtering. *IEEE assp magazine* 5(2), 4–24 (1988)
22. Yoshii, K., Itoyama, K., Goto, M.: Student’s t nonnegative matrix factorization and positive semidefinite tensor factorization for single-channel audio source separation. In: *International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. pp. 51–55. IEEE (2016)