

A MULTICHANNEL WIDELY LINEAR APPROACH TO BINAURAL NOISE REDUCTION USING AN ARRAY OF MICROPHONES

Jacob Benesty¹ and Jingdong Chen²

¹: INRS-EMT, University of Quebec
800 de la Gauchetiere Ouest, Suite 6900
Montreal, QC H5A 1K6, Canada

²: Northwestern Polytechnical University
127 Youyi West Road
Xi'an, Shaanxi 710072, China

ABSTRACT

This paper deals with the problem of binaural noise reduction using an array of microphones. This is a very important problem in applications such as teleconferencing and hearing aids where there is a need to mitigate the noise effect from the noisy signals picked up by multiple microphones and produce two “clean” outputs. The mitigation of the noise should be made in such a way that no audible distortion is added to the two outputs (this is the same as in the single-channel case) and meanwhile the spatial information of the desired sound source should be preserved so that, after noise reduction, the listener will still be able to localize the sound source thanks to his/her binaural hearing mechanism. In this paper, we present a novel approach to this problem where we first form a number of complex input signals from the multiple and real microphone observations. We also merge the two expected real outputs into a complex output signal. The widely linear estimation theory is then used to derive optimal noise reduction filters that can achieve noise reduction while preserving the desired signal (speech) and its spatial information. With this new formulation, the Wiener and minimum variance distortionless response (MVDR) filters are derived. Experiments are provided to justify the effectiveness of these filters.

Index Terms— Binaural noise reduction, microphone arrays, widely linear (WL) estimation, Wiener Filter, minimum variance distortionless response (MVDR) filter.

1. SIGNAL MODEL AND PROBLEM FORMULATION

Binaural noise reduction has a wide range of applications in areas such as teleconferencing, telecollaboration, social networks, and hearing aids. It has emerged as a very important research and engineering problem over the last two decades. Traditionally, this problem is tackled by extending the single-channel spectral modification based noise reduction technique to the binaural case by either posing some constraints on the suppression of each frequency band or using head-related transfer functions (HRTFs) to preserve the spatial information [1], [2]. Recently, we developed a new technique, which works for a stereo input and binaural output system [3]. The basic idea is to form a complex input signal from the stereo inputs and a complex output signal from the two expected real output channels. By doing so, the binaural noise reduction problem can be processed by a single-channel widely linear filter. In this paper, we attempt to extend the basic principle shown in [3] to a more general case where an array of microphones is used. Without loss of generality, we consider the signal model in which $2N$ microphones¹ capture a source signal convolved with acoustic impulse responses in some noise field. The signal received at the i th microphone is then expressed as

pressed as

$$\begin{aligned} y_{r,i}(t) &= g_i(t) * s(t) + v_{r,i}(t) \\ &= x_{r,i}(t) + v_{r,i}(t), \quad i = 1, 2, \dots, 2N, \end{aligned} \quad (1)$$

where $g_i(t)$ is the acoustic impulse response from the unknown speech source, $s(t)$, location to the i th microphone, $*$ stands for linear convolution, and $x_{r,i}(t)$ and $v_{r,i}(t)$ are, respectively, the convolved speech and additive noise received at microphone i . We assume that the impulse responses are time invariant. We also assume that the signals $x_{r,i}(t) = g_i(t) * s(t)$ and $v_{r,i}(t)$ are uncorrelated, zero mean, real, and broadband.

In binaural noise reduction, it is desired to simultaneously recover the convolved speech signals at two microphones. In this paper, we consider recovering the signals $x_{r,1}(t)$ and $x_{r,N+1}(t)$ given the observations $y_{r,i}(t)$, $i = 1, 2, \dots, 2N$. This means that the desired signals in our problem are the speech signals received at the first and $(N + 1)$ th microphones². It is clear then that we have two objectives. The first one is to attenuate the contribution of the noise terms $v_{r,1}(t)$ and $v_{r,N+1}(t)$ as much as possible. The second objective is to preserve $x_{r,1}(t)$ and $x_{r,N+1}(t)$ with their spatial information, so that with the enhanced signals, along with our binaural hearing process, we will still be able to localize the source $s(t)$. This is the well-known problem of binaural noise reduction.

We have $2N$ real input and two real output signals. It is convenient, however, to work in the complex domain in order that the original binaural problem is transformed to a single-output noise reduction problem with a microphone array. Indeed, from the $2N$ real microphone signals given in (1), we can form N complex microphone signals as

$$\begin{aligned} y_n(t) &= y_{r,n}(t) + jy_{r,N+n}(t) \\ &= x_n(t) + v_n(t), \quad n = 1, 2, \dots, N, \end{aligned} \quad (2)$$

where $j = \sqrt{-1}$,

$$x_n(t) = x_{r,n}(t) + jx_{r,N+n}(t), \quad n = 1, 2, \dots, N \quad (3)$$

is the complex convolved speech signal, and

$$v_n(t) = v_{r,n}(t) + jv_{r,N+n}(t), \quad n = 1, 2, \dots, N \quad (4)$$

is the complex additive noise. Now, our problem can be stated as follows: given the N complex microphone signals, $y_n(t)$, $n = 1, 2, \dots, N$, which are a mixture of the uncorrelated complex signals $x_n(t)$ and $v_n(t)$, our goal is to recover $x_1(t) = x_{r,1}(t) + jx_{r,N+1}(t)$ (i.e., our desired signal) the best way we can, including the phase, which is important for the localization of the source signal.

Effort of the second author is partially supported by the Anhui Science and Technology Project (11010202191).

¹The generalization to an odd number of microphones is straightforward.

²Note that the ideas and algorithms developed in this paper can be applied to any other pair of microphones.

The signal model given in (2) can be put into a vector form as

$$\mathbf{y}(t) = \mathbf{x}(t) + \mathbf{v}(t), \quad (5)$$

where

$$\mathbf{y}(t) = [\mathbf{y}_1^T(t) \quad \mathbf{y}_2^T(t) \quad \cdots \quad \mathbf{y}_N^T(t) \quad \cdots \quad \mathbf{y}_N^T(t)]^T,$$

the superscript T is the transpose operator,

$$\mathbf{y}_n(t) = [y_n(t) \quad y_n(t-1) \quad \cdots \quad y_n(t-L+1)]^T$$

is a vector of length L , and $\mathbf{x}(t)$ and $\mathbf{v}(t)$ are defined in a similar way to $\mathbf{y}(t)$. Since the clean and noise signals are uncorrelated by assumption, the correlation matrix (of size $NL \times NL$) of the noisy signal is

$$\Phi_{\mathbf{y}} = E[\mathbf{y}(t)\mathbf{y}^H(t)] = \Phi_{\mathbf{x}} + \Phi_{\mathbf{v}}, \quad (6)$$

where the superscript H is the transpose-conjugate operator, and $\Phi_{\mathbf{x}} = [\mathbf{x}(t)\mathbf{x}^H(t)]$ and $\Phi_{\mathbf{v}} = E[\mathbf{v}(t)\mathbf{v}^H(t)]$ are the correlation matrices of $\mathbf{x}(t)$ and $\mathbf{v}(t)$, respectively.

2. WIDELY LINEAR FILTERING

As seen from the model given in (2), we deal with complex random variables. A very important statistical characteristic of a complex random variable (CRV) is the so-called circularity property or lack of it (noncircularity) [4], [5]. A zero-mean CRV, z , is circular if and only if the only nonnull moments and cumulants are the moments and cumulants constructed with the same power in z and z^* [6], [7], where the superscript $*$ denotes complex conjugation. In particular, z is said to be a second-order circular CRV (CCRV) if its so-called pseudo-variance [4] is equal to zero, i.e., $E(z^2) = 0$, while its variance is nonnull, i.e., $E(|z|^2) \neq 0$. This means that the second-order behavior of a CCRV is well described by its variance. If the pseudo-variance $E(z^2)$ is not equal to 0, the CRV z is then noncircular. A good measure of the second-order circularity is the circularity quotient [4] defined as the ratio between the pseudo-variance and the variance, i.e.,

$$\gamma_z = \frac{E(z^2)}{E(|z|^2)}. \quad (7)$$

It is easy to show that $0 \leq |\gamma_z| \leq 1$. If $\gamma_z = 0$, z is a second-order CCRV; otherwise, z is noncircular.

Now, let us examine whether the complex desired signal, $x_1(t) = x_{r,1}(t) + jx_{r,N+1}(t)$, is second-order circular or not. We have

$$\begin{aligned} \gamma_{x_1} &= \frac{E[x_1^2(t)]}{E[|x_1(t)|^2]} \\ &= \frac{E[x_{r,1}^2(t)] - E[x_{r,N+1}^2(t)] + j2E[x_{r,1}(t)x_{r,N+1}(t)]}{\phi_{x_1}}, \end{aligned} \quad (8)$$

where $\phi_{x_1} = E[|x_1(t)|^2]$ is the variance of $x_1(t)$. One can check from (8) that the CRV $x_1(t)$ is second-order circular (i.e., $\gamma_{x_1} = 0$) if and only if

$$E[x_{r,1}^2(t)] = E[x_{r,N+1}^2(t)] \text{ and } E[x_{r,1}(t)x_{r,N+1}(t)] = 0. \quad (9)$$

Since the signals $x_{r,1}(t)$ and $x_{r,N+1}(t)$ come from the same source, they are in general correlated. As a result, the second condition in (9) should not be true. Therefore, we can safely state that the complex desired signal, $x_1(t)$, is noncircular, and so is the complex microphone signal, $y_1(t)$. If the noise terms at the two microphones are

assumed to be uncorrelated and have the same power, then $\gamma_{v_1} = 0$ [i.e., $v(t)$ is a second-order CCRV].

Since we deal with noncircular CRVs as demonstrated above, the classical linear estimation technique [8], which is developed for processing real signals or CCRVs, cannot be applied. Instead, an estimate of $x_1(t)$ should be obtained using the widely linear (WL) estimation theory as [5], [9]

$$\hat{x}_1(t) = \mathbf{h}^H \mathbf{y}(t) + \mathbf{h}'^H \mathbf{y}^*(t) = \tilde{\mathbf{h}}^H \tilde{\mathbf{y}}(t), \quad (10)$$

where \mathbf{h} and \mathbf{h}' are two complex FIR filters of length NL and

$$\tilde{\mathbf{h}} = \begin{bmatrix} \mathbf{h} \\ \mathbf{h}' \end{bmatrix} \quad (11)$$

$$\tilde{\mathbf{y}}(t) = \begin{bmatrix} \mathbf{y}(t) \\ \mathbf{y}^*(t) \end{bmatrix} \quad (12)$$

are the augmented WL filter and observation vector, respectively, both of length $2NL$. We can rewrite (10) as

$$\hat{x}_1(t) = \tilde{\mathbf{h}}^H [\tilde{\mathbf{x}}(t) + \tilde{\mathbf{v}}(t)] = x_f(t) + v_{rn}(t), \quad (13)$$

where $\tilde{\mathbf{x}}(t)$ and $\tilde{\mathbf{v}}(t)$ are defined in a similar way to $\tilde{\mathbf{y}}(t)$,

$$x_f(t) = \tilde{\mathbf{h}}^H \tilde{\mathbf{x}}(t) \quad (14)$$

is a filtered version of the desired signal vector and its conjugate, and

$$v_{rn}(t) = \tilde{\mathbf{h}}^H \tilde{\mathbf{v}}(t) \quad (15)$$

is the residual noise.

From (13), we see that $\hat{x}_1(t)$ depends on the vector $\tilde{\mathbf{x}}(t)$. However, our desired signal at time t is not the whole vector $\tilde{\mathbf{x}}(t)$ but only the sample $x_1(t)$; so we should decompose the vector $\tilde{\mathbf{x}}(t)$ into two orthogonal vectors: one correlated and the other uncorrelated with the desired signal sample, $x_1(t)$. This decomposition is given as

$$\tilde{\mathbf{x}}(t) = x_1(t)\boldsymbol{\rho}_{\tilde{\mathbf{x}}_{x_1}} + \tilde{\mathbf{x}}_i(t) = \tilde{\mathbf{x}}_d(t) + \tilde{\mathbf{x}}_i(t), \quad (16)$$

where

$$\tilde{\mathbf{x}}_d(t) = x_1(t)\boldsymbol{\rho}_{\tilde{\mathbf{x}}_{x_1}} \quad (17)$$

is the desired signal vector,

$$\tilde{\mathbf{x}}_i(t) = \tilde{\mathbf{x}}(t) - \tilde{\mathbf{x}}_d(t) \quad (18)$$

is the interference signal vector,

$$\boldsymbol{\rho}_{\tilde{\mathbf{x}}_{x_1}} = \frac{E[\tilde{\mathbf{x}}(t)x_1^*(t)]}{E[|x_1(t)|^2]} \quad (19)$$

is the normalized [with respect to $x_1(t)$] correlation vector between $\tilde{\mathbf{x}}(t)$ and $x_1(t)$, and

$$E[\tilde{\mathbf{x}}_i(t)x_1^*(t)] = \mathbf{0}_{2NL \times 1}. \quad (20)$$

Substituting (16) into (13), we obtain

$$\begin{aligned} \hat{x}_1(t) &= \tilde{\mathbf{h}}^H [\tilde{\mathbf{x}}_d(t) + \tilde{\mathbf{x}}_i(t) + \tilde{\mathbf{v}}(t)] \\ &= x_{fd}(t) + x_{ri}(t) + v_{rn}(t), \end{aligned} \quad (21)$$

where

$$x_{fd}(t) = \tilde{\mathbf{h}}^H \tilde{\mathbf{x}}_d(t) = x_1(t)\tilde{\mathbf{h}}^H \boldsymbol{\rho}_{\tilde{\mathbf{x}}_{x_1}} \quad (22)$$

is the filtered desired signal and

$$x_{ri}(t) = \tilde{\mathbf{h}}^H \tilde{\mathbf{x}}_i(t) \quad (23)$$

is the residual interference. We observe that the estimate of the desired signal at time t is the sum of three terms that are mutually uncorrelated. Therefore, the variance of $\hat{x}_1(t)$ is

$$\phi_{\hat{x}_1} = \phi_{x_{fd}} + \phi_{x_{ri}} + \phi_{v_{rn}}, \quad (24)$$

where

$$\phi_{x_{fd}} = \phi_{x_1} \left| \tilde{\mathbf{h}}^H \boldsymbol{\rho}_{\tilde{\mathbf{x}}x_1} \right|^2 = \tilde{\mathbf{h}}^H \boldsymbol{\Phi}_{\tilde{\mathbf{x}}d} \tilde{\mathbf{h}}, \quad (25)$$

$$\phi_{x_{ri}} = \tilde{\mathbf{h}}^H \boldsymbol{\Phi}_{\tilde{\mathbf{x}}i} \tilde{\mathbf{h}} = \tilde{\mathbf{h}}^H \boldsymbol{\Phi}_{\tilde{\mathbf{x}}} \tilde{\mathbf{h}} - \phi_{x_1} \left| \tilde{\mathbf{h}}^H \boldsymbol{\rho}_{\tilde{\mathbf{x}}x_1} \right|^2, \quad (26)$$

$$\phi_{v_{rn}} = \tilde{\mathbf{h}}^H \boldsymbol{\Phi}_{\tilde{\mathbf{v}}} \tilde{\mathbf{h}}, \quad (27)$$

$\boldsymbol{\Phi}_{\tilde{\mathbf{x}}d} = \phi_{x_1} \boldsymbol{\rho}_{\tilde{\mathbf{x}}x_1} \boldsymbol{\rho}_{\tilde{\mathbf{x}}x_1}^H$ is the correlation matrix (whose rank is equal to 1) of $\tilde{\mathbf{x}}_d(t)$, and $\boldsymbol{\Phi}_{\tilde{\mathbf{x}}i} = E[\tilde{\mathbf{x}}_i(t) \tilde{\mathbf{x}}_i^H(t)]$, $\boldsymbol{\Phi}_{\tilde{\mathbf{x}}} = E[\tilde{\mathbf{x}}(t) \tilde{\mathbf{x}}^H(t)]$, and $\boldsymbol{\Phi}_{\tilde{\mathbf{v}}} = E[\tilde{\mathbf{v}}(t) \tilde{\mathbf{v}}^H(t)]$ are the correlation matrices of $\tilde{\mathbf{x}}_i(t)$, $\tilde{\mathbf{x}}(t)$, and $\tilde{\mathbf{v}}(t)$, respectively.

It is clear from (21) that the objective of our noise reduction problem is to find optimal filters that can minimize the effect of $x_{ri}(t) + v_{rn}(t)$ while preserving the desired signal, $x_1(t)$.

3. OPTIMAL FILTERS

Before deriving the optimal filters, we first define the mean-square-error (MSE) criterion.

The error signal between the estimated and desired signals is defined as

$$e(t) = \hat{x}_1(t) - x_1(t) = x_{fd}(t) + x_{ri}(t) + v_{rn}(t) - x_1(t). \quad (28)$$

We can rewrite (28) as

$$e(t) = e_d(t) + e_r(t), \quad (29)$$

where

$$e_d(t) = x_{fd}(t) - x_1(t) = \left(\tilde{\mathbf{h}}^H \boldsymbol{\rho}_{\tilde{\mathbf{x}}x_1} - 1 \right) x_1(t) \quad (30)$$

is the speech distortion due to the WL filter, and

$$e_r(t) = x_{ri}(t) + v_{rn}(t) = \tilde{\mathbf{h}}^H \tilde{\mathbf{x}}_i(t) + \tilde{\mathbf{h}}^H \tilde{\mathbf{v}}(t) \quad (31)$$

represents the residual interference-plus-noise. The two error signals $e_d(t)$ and $e_r(t)$ are clearly uncorrelated.

The classical MSE criterion is then

$$\begin{aligned} J(\tilde{\mathbf{h}}) &= E[|e(t)|^2] = E[|e_d(t)|^2] + E[|e_r(t)|^2] \\ &= \phi_{x_1} \left| \tilde{\mathbf{h}}^H \boldsymbol{\rho}_{\tilde{\mathbf{x}}x_1} - 1 \right|^2 + \tilde{\mathbf{h}}^H \boldsymbol{\Phi}_{\tilde{\mathbf{x}}i} \tilde{\mathbf{h}} + \tilde{\mathbf{h}}^H \boldsymbol{\Phi}_{\tilde{\mathbf{v}}} \tilde{\mathbf{h}} \\ &= J_d(\tilde{\mathbf{h}}) + J_r(\tilde{\mathbf{h}}), \end{aligned} \quad (32)$$

where

$$J_d(\tilde{\mathbf{h}}) = E[|e_d(t)|^2] = \phi_{x_1} \left| \tilde{\mathbf{h}}^H \boldsymbol{\rho}_{\tilde{\mathbf{x}}x_1} - 1 \right|^2, \quad (33)$$

$$J_r(\tilde{\mathbf{h}}) = E[|e_r(t)|^2] = \tilde{\mathbf{h}}^H \boldsymbol{\Phi}_{\tilde{\mathbf{x}}i} \tilde{\mathbf{h}}, \quad (34)$$

and

$$\boldsymbol{\Phi}_{\tilde{\mathbf{in}}} = \boldsymbol{\Phi}_{\tilde{\mathbf{x}}i} + \boldsymbol{\Phi}_{\tilde{\mathbf{v}}} = \boldsymbol{\Phi}_{\tilde{\mathbf{y}}} - \phi_{x_1} \boldsymbol{\rho}_{\tilde{\mathbf{x}}x_1} \boldsymbol{\rho}_{\tilde{\mathbf{x}}x_1}^H \quad (35)$$

is the correlation matrix of the interference-plus-noise.

Given the above MSE criteria, the filters that can achieve noise reduction can be derived by either directly minimizing $J(\tilde{\mathbf{h}})$ or minimizing $J_d(\tilde{\mathbf{h}})$ or $J_r(\tilde{\mathbf{h}})$ with some constraints.

3.1. WL Wiener Filter

In our problem, the WL Wiener filter can be obtained by taking the gradient of (32) with respect to $\tilde{\mathbf{h}}$ and equating the result to zero. We easily get

$$\tilde{\mathbf{h}}_W = \boldsymbol{\Phi}_{\tilde{\mathbf{y}}}^{-1} \boldsymbol{\Phi}_{\tilde{\mathbf{x}}i} \tilde{\mathbf{i}} = [\mathbf{I}_{2NL} - \boldsymbol{\Phi}_{\tilde{\mathbf{y}}}^{-1} \boldsymbol{\Phi}_{\tilde{\mathbf{v}}}] \tilde{\mathbf{i}}, \quad (36)$$

where $\tilde{\mathbf{i}}$ is the first column of the identity matrix \mathbf{I}_{2NL} of size $2NL \times 2NL$. From (35), we can write the inverse of $\boldsymbol{\Phi}_{\tilde{\mathbf{y}}}$ according to the Woodbury's identity:

$$\boldsymbol{\Phi}_{\tilde{\mathbf{y}}}^{-1} = \boldsymbol{\Phi}_{\tilde{\mathbf{in}}}^{-1} - \frac{\boldsymbol{\Phi}_{\tilde{\mathbf{in}}}^{-1} \boldsymbol{\rho}_{\tilde{\mathbf{x}}x_1} \boldsymbol{\rho}_{\tilde{\mathbf{x}}x_1}^H \boldsymbol{\Phi}_{\tilde{\mathbf{in}}}^{-1}}{\phi_{x_1}^{-1} + \boldsymbol{\rho}_{\tilde{\mathbf{x}}x_1}^H \boldsymbol{\Phi}_{\tilde{\mathbf{in}}}^{-1} \boldsymbol{\rho}_{\tilde{\mathbf{x}}x_1}}. \quad (37)$$

Substituting (37) into (36), we get an alternative form of the WL Wiener filter, i.e.,

$$\tilde{\mathbf{h}}_W = \frac{\boldsymbol{\Phi}_{\tilde{\mathbf{in}}}^{-1} \boldsymbol{\rho}_{\tilde{\mathbf{x}}x_1}}{\phi_{x_1}^{-1} + \boldsymbol{\rho}_{\tilde{\mathbf{x}}x_1}^H \boldsymbol{\Phi}_{\tilde{\mathbf{in}}}^{-1} \boldsymbol{\rho}_{\tilde{\mathbf{x}}x_1}}, \quad (38)$$

which can be more useful for practical implementation as the matrix $\boldsymbol{\Phi}_{\tilde{\mathbf{in}}}$ is generally better conditioned than $\boldsymbol{\Phi}_{\tilde{\mathbf{y}}}$ in real-world applications.

3.2. WL MVDR Filter

The WL MVDR filter can be derived by minimizing either $E[|\hat{x}_1(t)|^2]$, $J(\tilde{\mathbf{h}})$, or $J_r(\tilde{\mathbf{h}})$ subject to the constraint that $J_d(\tilde{\mathbf{h}}) = 0$. Using $J(\tilde{\mathbf{h}})$, the problem can be mathematically written as

$$\min_{\tilde{\mathbf{h}}} J(\tilde{\mathbf{h}}) \quad \text{subject to} \quad \tilde{\mathbf{h}}^H \boldsymbol{\rho}_{\tilde{\mathbf{x}}x_1} = 1. \quad (39)$$

Using a Lagrange multiplier to adjoin the constraint to the cost function, taking the gradient with respect to $\tilde{\mathbf{h}}$, and then equating the result to zero, we obtain

$$\tilde{\mathbf{h}}_{\text{MVDR}} = \frac{\boldsymbol{\Phi}_{\tilde{\mathbf{y}}}^{-1} \boldsymbol{\rho}_{\tilde{\mathbf{x}}x_1}}{\boldsymbol{\rho}_{\tilde{\mathbf{x}}x_1}^H \boldsymbol{\Phi}_{\tilde{\mathbf{y}}}^{-1} \boldsymbol{\rho}_{\tilde{\mathbf{x}}x_1}}. \quad (40)$$

Now, using $J_r(\tilde{\mathbf{h}})$, the problem can be written into the following form:

$$\min_{\tilde{\mathbf{h}}} \tilde{\mathbf{h}}^H \boldsymbol{\Phi}_{\tilde{\mathbf{x}}i} \tilde{\mathbf{h}} \quad \text{subject to} \quad \tilde{\mathbf{h}}^H \boldsymbol{\rho}_{\tilde{\mathbf{x}}x_1} = 1, \quad (41)$$

for which the solution is

$$\tilde{\mathbf{h}}_{\text{MVDR}} = \frac{\boldsymbol{\Phi}_{\tilde{\mathbf{in}}}^{-1} \boldsymbol{\rho}_{\tilde{\mathbf{x}}x_1}}{\boldsymbol{\rho}_{\tilde{\mathbf{x}}x_1}^H \boldsymbol{\Phi}_{\tilde{\mathbf{in}}}^{-1} \boldsymbol{\rho}_{\tilde{\mathbf{x}}x_1}}. \quad (42)$$

It is easy to check that the two forms of the MVDR filter in (40) and (42) are identical.

4. EXPERIMENTAL EVALUATION

Experiments were conducted with the speech source recorded in a reverberant but quiet room. The room is 6 m long and 5 m wide. For ease of exposition, positions in the room are designated by (x, y) coordinates with reference to one corner of the room, $0 \leq x \leq 6$ and $0 \leq y \leq 5$. An equispaced linear array with six omnidirectional microphones is configured where the first and last microphones are, respectively, at (3.4, 0.5) and (3.9, 0.5), and spacing between two neighboring microphones is 0.1 m. A female talker reads a story

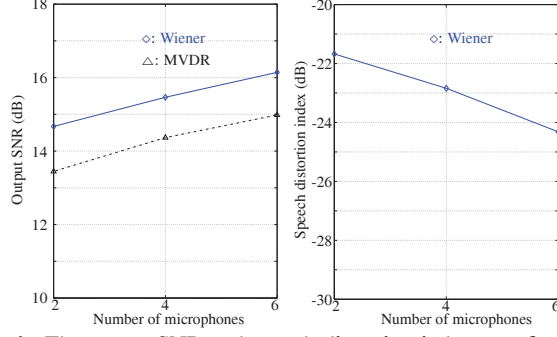


Fig. 1. The output SNR and speech distortion index as a function of the number of microphones for the Wiener and MVDR filters in white noise: $L = 40$ and input SNR of 10 dB. Note that the values of the speech distortion index for the MVDR filter are not displayed because they are very small.

while standing at (1.3, 3.0) and her voice is recorded with a sampling rate of 8 kHz. The recorded multichannel signals are treated as the convolved speech. White Gaussian noise is then added to the signals so that the input SNR is equal to 10 dB. The overall length of the signal is 30 s. We set L to 40 based on the results from a previous study [3].

We choose to implement the WL Wiener and MVDR filters according to (38) and (42), respectively, for which we need to know the parameters Φ_{in} , ϕ_{x_1} , and $\rho_{\tilde{x}x_1}$. In this paper, we first compute the two covariance matrices $\Phi_{\tilde{y}}$ and $\Phi_{\tilde{v}}$ directly from the noisy and noise signals using a short-time average. Specifically, at each time instant t , an estimate of $\Phi_{\tilde{y}}$, i.e., $\hat{\Phi}_{\tilde{y}}(t)$, is computed using the most recent 480 \tilde{y} vectors, and $\hat{\Phi}_{\tilde{v}}(t)$ is computed in the same manner but from the noise signal. An estimate of $\hat{\Phi}_{\tilde{x}}$ is calculated as $\hat{\Phi}_{\tilde{x}}(t) = \hat{\Phi}_{\tilde{y}}(t) - \hat{\Phi}_{\tilde{v}}(t)$. Then $\hat{\phi}_{x_1}(t)$ is taken as the first element of $\hat{\Phi}_{\tilde{x}}(t)$ and $\hat{\rho}_{\tilde{x}x_1}(t)$ is the first column of $\hat{\Phi}_{\tilde{x}}(t)$ normalized by $\hat{\phi}_{x_1}(t)$. The matrix $\hat{\Phi}_{\text{in}}(t)$ can then be computed according to (35).

Substituting $\hat{\Phi}_{\text{in}}(t)$, $\hat{\phi}_{x_1}(t)$, and $\hat{\rho}_{\tilde{x}x_1}(t)$ into (38) and (42), we implemented the WL Wiener and MVDR filters. To evaluate the noise reduction performance of these two filters, we used the output SNR and speech distortion index, which are defined as [3]:

$$\text{oSNR}(\tilde{\mathbf{h}}) = \frac{\phi_{x_{\text{fd}}}}{\phi_{x_{\text{ri}}} + \phi_{v_{\text{rn}}}}, \quad (43)$$

$$v_{\text{sd}}(\tilde{\mathbf{h}}) = \frac{E[|x_{\text{fd}}(t) - x_1(t)|^2]}{\phi_{x_1}}. \quad (44)$$

These two measures are computed using a long-time average. The results, as a function of the number of microphones, are plotted in Fig. 1. (Note that for a fair comparison, in the two-microphone case, we used the 1st and 4th microphones, and in the four-microphone case, we used the 1st, 2nd, 4th, and 5th microphones.) For both the WL Wiener and MVDR filters, we see that the output SNR increases while the speech distortion index decreases with the number of microphones. This clearly demonstrates the advantage of using multiple microphones for binaural noise reduction. It is also seen that the WL MVDR can improve SNR without adding speech distortion, though its SNR improvement is less than that of the WL Wiener filter.

To visualize the spatial sound effect, we computed the cross-correlation function between the two real output channels [i.e., the real and imaginary parts of $\hat{x}_1(t)$] every 64 ms using a short-time average with a frame size of 64 ms. The location of the maximum value

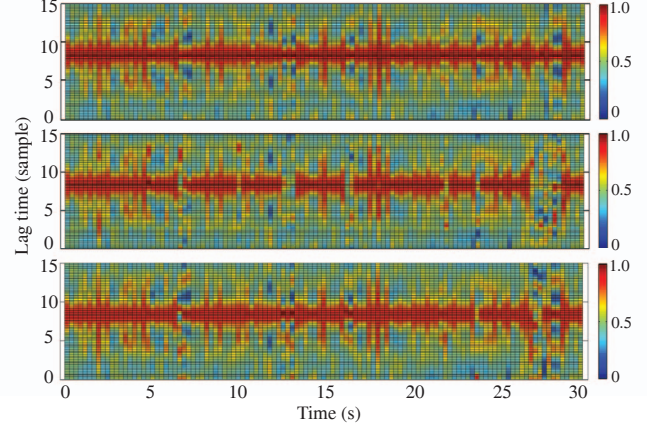


Fig. 2. The contours of the short-time cross-correlation coefficients between the left and right channels. Upper trace: clean speech; middle trace: noisy speech (with an input SNR of 10 dB); lower trace: enhanced speech using the MVDR filter with 6 microphones.

of this function indicates the position of the talker at time instant t . Figure 2 shows the contours of the time-varying cross-correlation function of the clean and noisy signals as well as that of the enhanced signal (due to space limitation, we only show the result with the MVDR filter using 6 microphones). One can notice that the noise has dramatically modified the sound spatial effect. It is clearly seen that the MVDR filter with 6 microphones has significantly recovered the spatial effect.

5. SUMMARY

In this paper, we presented a multichannel widely linear (WL) approach to the problem of binaural noise reduction using a microphone array. We first formed a number of complex input signals from the multiple and real microphone observations. We also merged the two expected real outputs into a complex output signal. By doing so, the binaural noise reduction problem is converted into one of multiple-input/single-output WL filtering. The WL estimation theory was then used to derive the optimal WL Wiener and MVDR filters that can achieve noise reduction while preserving the desired signal (speech) and its spatial information. Experimental results have justified the effectiveness of these filters.

6. REFERENCES

- [1] B. Kollmeier, J. Peissig, and V. Hohmann, "Binaural noise-reduction hearing aid scheme with real-time processing in the frequency domain," *Scand Audiol Suppl*, vol. 38, pp. 28–38, 1993.
- [2] J. Li, S. Sakamoto, S. Hongo, M. Akagi, and Y. Suzuki, "Two-stage binaural speech enhancement with Wiener filter based on equalization-cancellation model," in *Proc. IEEE WASPAA*, 2009, pp. 133–136.
- [3] J. Benesty, J. Chen, and Y. Huang, "Binaural noise reduction in the time domain with a stereo setup," *IEEE Trans. Audio, Speech, Lang. Process.*, vol. 19, pp. 2260–2272, Nov. 2011.
- [4] E. Ollila, "On the circularity of a complex random variable," *IEEE Signal Process. Lett.*, vol. 15, pp. 841–844, 2008.
- [5] D. P. Mandic and S. L. Goh, *Complex Valued Nonlinear Adaptive Filters: Noncircularity, Widely Linear and Neural Models*. Wiley, 2009.
- [6] P. O. Amblard, M. Gaeta, and J. L. Lacoume, "Statistics for complex variables and signals—Part I: variables," *Signal Process.*, vol. 53, pp. 1–13, 1996.
- [7] P. O. Amblard, M. Gaeta, and J. L. Lacoume, "Statistics for complex variables and signals—Part II: signals," *Signal Process.*, vol. 53, pp. 15–25, 1996.
- [8] J. Benesty, J. Chen, Y. Huang, and I. Cohen, *Noise Reduction in Speech Processing*. Berlin, Germany: Springer-Verlag, 2009.
- [9] B. Picinbono and P. Chevalier, "Widely linear estimation with complex data," *IEEE Trans. Signal Process.*, vol. 43, pp. 2030–2033, Aug. 1995.