

Evaluation of Robust Constrained MFMVDR Filtering for Single-Channel Speech Enhancement

Dörte Fischer, Simon Doclo¹.

Department of Medical Physics and Acoustics and Cluster of Excellence Hearing4all, University of Oldenburg, Germany.
Email:{doerte.fischer, simon.doclo}@uni-oldenburg.de

Abstract

By considering the multi-frame signal model, speech correlation between different time-frames can be exploited. Based on this signal model, the multi-frame minimum variance distortionless response (MFMVDR) filter for single-channel speech enhancement has been derived, which minimizes the total signal output power while avoiding speech distortion. It has been shown that the MFMVDR filter is very sensitive to estimation errors in the speech correlation vector resulting in correlated speech components being mistakenly suppressed. Inspired by robust beamforming approaches, in this paper we propose a robust constrained MFMVDR filter for single-channel speech enhancement by estimating the speech correlation vector that maximizes the total signal output power within a spherical uncertainty set. For the upper bound of the spherical uncertainty set, we propose to use a trained mapping function that depends on the a-priori signal-to-noise ratio (SNR). Experimental results for different noise types and SNRs show that the proposed robust approach yields a more accurate estimate of the speech correlation vector. A perceptual evaluation shows that the robust constrained MFMVDR filter leads to an improved speech quality but a lower noise reduction than the original non-robust MFMVDR filter, while still being preferred in overall quality.

1 Introduction

Using speech communication devices (e.g., hearing aids) in noisy environments may lead to a degraded speech quality and speech intelligibility. In these acoustic situations, speech enhancement algorithms are required to suppress the undesired noise while limiting speech distortion. Typically, single-channel speech enhancement algorithms are designed in the short-time Fourier transform (STFT) domain. When assuming that neighboring STFT coefficients are uncorrelated over time and frequency, the noisy STFT coefficients are processed by applying a (real-valued) gain to each time-frequency point [1]. When using the more realistic assumption that neighboring STFT coefficients are correlated over time, it has been proposed to process the noisy STFT coefficients by applying a (complex-valued) finite-impulse response (FIR) filter [2–5].

In [3, 4], a multi-frame signal model has been proposed, where a *speech correlation vector* contains the speech correlation between the current and previous time-frames. Conceptually, this multi-frame signal model is similar to a multi-channel signal model when interpreting time-frames as microphone inputs and the speech correlation vector as the steering vector. Similarly to the well-known minimum variance distortionless response (MVDR) beamformer [6, 7], the *multi-frame MVDR* (MFMVDR) filter for single-channel speech enhancement was derived in [3, 4], minimizing the total signal output power while not distorting correlated speech components.

In practice, obviously only the noisy speech signal is available such that the speech correlation vector needs to be estimated from the noisy STFT coefficients. In [8], a maximum-likelihood (ML) estimator for the speech correlation vector was proposed. It was shown that when using this speech correlation vector, the MFMVDR filter introduces low speech distortion and achieves a good noise reduction. In [9], we showed that accurately estimating the speech correlation vector is crucial in that even a small mismatch between the optimal and the estimated speech correlation vector may lead to a degraded performance of the MFMVDR filter.

¹This work was supported in part by the joint Lower Saxony-Israeli Project ATHENA and by the Cluster of Excellence EXC 1077 Hearing4all, funded by the German Research Foundation (DFG)

In the area of array processing several techniques have been proposed to increase the robustness of beamformers against estimation errors of the steering vector [10–13]. Inspired by the robust MVDR beamformer in [10], which estimates the steering vector that maximizes the total signal output power of the MVDR beamformer within a spherical uncertainty set, in this paper we propose a *robust constrained (RC) MFMVDR* filter by estimating the speech correlation vector that maximizes the total signal output power of the MFMVDR filter within a *spherical uncertainty set*. This spherical uncertainty imposes an upper bound on the norm of the mismatch vector, i.e., the difference between the speech correlation vector and the presumed speech correlation vector, e.g., calculated using the ML method in [8]. Since oracle simulations using several speech and noise signals at different signal-to-noise ratios (SNR) showed that the norm of the mismatch vector decreases with increasing SNR, we trained a mapping function to set the upper bound depending on the a-priori SNR for each time-frequency point. Evaluation results for different speech signals, noise types and SNRs show that the proposed RC speech correlation vector achieves a lower mean-square error than the ML estimate. Based on a pairwise preference listening experiment with 15 participants, we show that the proposed RC MFMVDR filter results in a significantly better speech quality than the ML MFMVDR filter although the amount of noise reduction is significantly lower. Nevertheless, the quality of the RC MFMVDR filter is preferred over the overall quality of the ML MFMVDR filter.

2 Multi-Frame Signal Model

We consider a single-channel setup, where a speech signal is degraded by additive noise. In the STFT domain, the complex-valued noisy speech coefficient $Y(k, m)$ at frequency-bin k and time-frame m is given by

$$Y(k, m) = X(k, m) + N(k, m), \quad (1)$$

with $X(k, m)$ and $N(k, m)$ denoting the speech and the noise coefficients, respectively. For conciseness, in the remainder of the paper the frequency-bin k will be omitted where possible.

In multi-frame single-channel speech enhancement approaches [2] the speech coefficient $X(m)$ is estimated by applying a complex-valued FIR filter $\mathbf{h}(m)$ to the noisy speech vector $\mathbf{y}(m)$, i.e.,

$$\hat{X}(m) = \mathbf{h}^H(m) \mathbf{y}(m), \quad (2)$$

where H denotes the Hermitian operator. The noisy speech vector $\mathbf{y}(m)$ contains L consecutive noisy speech coefficients and the filter $\mathbf{h}(m)$ contains L time-varying filter coefficients, i.e.,

$$\mathbf{y}(m) = [Y(m), Y(m-1), \dots, Y(m-L+1)]^T, \quad (3)$$

$$\mathbf{h}(m) = [H_0(m), H_1(m), \dots, H_{L-1}(m)]^T. \quad (4)$$

Using (1), the noisy speech vector $\mathbf{y}(m)$ is given by

$$\mathbf{y}(m) = \mathbf{x}(m) + \mathbf{n}(m), \quad (5)$$

where the speech vector $\mathbf{x}(m)$ and the noise vector $\mathbf{n}(m)$ are defined similarly as in (3). Assuming that the speech and noise signals are uncorrelated, the $L \times L$ -dimensional noisy speech correlation matrix $\Phi_{\mathbf{y}\mathbf{y}}(m) = \mathbb{E}[\mathbf{y}(m)\mathbf{y}^H(m)]$, with $\mathbb{E}[\cdot]$ the expectation operator, is given by

$$\Phi_{\mathbf{y}\mathbf{y}}(m) = \Phi_{\mathbf{x}\mathbf{x}}(m) + \Phi_{\mathbf{n}\mathbf{n}}(m), \quad (6)$$

with $\Phi_{xx}(m)$ the speech correlation matrix and $\Phi_{nn}(m)$ the noise correlation matrix.

To exploit the speech correlation across time-frames, it has been proposed in [3, 4] to decompose the speech vector $\mathbf{x}(m)$ into a temporally correlated speech component $\mathbf{s}(m)$ and a temporally uncorrelated speech component $\mathbf{x}'(m)$ with respect to the speech coefficient $X(m)$, i.e.,

$$\mathbf{x}(m) = \mathbf{s}(m) + \mathbf{x}'(m) \quad (7)$$

with

$$\mathbf{s}(m) = \gamma_x(m)X(m). \quad (8)$$

The (highly time-varying) speech correlation vector $\gamma_x(m)$ is defined as

$$\gamma_x(m) = \frac{\mathbb{E}[\mathbf{x}(m)X^*(m)]}{\mathbb{E}[|X(m)|^2]} = \frac{\Phi_{xx}(m)\mathbf{e}}{\mathbf{e}^T\Phi_{xx}(m)\mathbf{e}}, \quad (9)$$

where $*$ denotes the complex-conjugate operator and $\mathbf{e} = [1, 0, \dots, 0]^T$ is an L -dimensional selection vector. Through the normalization with $\mathbf{e}^T\Phi_{xx}(m)\mathbf{e}$, which corresponds to the speech power spectral density (PSD) $\phi_X(m)$, the first element of the speech correlation vector is equal to 1, i.e.,

$$\mathbf{e}^T\gamma_x(m) = 1. \quad (10)$$

Substituting (7) and (8) into (5) and considering the uncorrelated speech component $\mathbf{x}'(m)$ as an interference, the *multi-frame signal model* is given by

$$\boxed{\mathbf{y}(m) = \gamma_x(m)X(m) + \mathbf{u}(m)} \quad (11)$$

with $\mathbf{u}(m) = \mathbf{x}'(m) + \mathbf{n}(m)$ the undesired signal vector.

Similarly to (9), the noisy speech correlation vector $\gamma_y(m)$ and the noise correlation vector $\gamma_n(m)$ can be defined as

$$\gamma_y(m) = \frac{\Phi_{yy}(m)\mathbf{e}}{\mathbf{e}^T\Phi_{yy}(m)\mathbf{e}}, \quad \gamma_n(m) = \frac{\Phi_{nn}(m)\mathbf{e}}{\mathbf{e}^T\Phi_{nn}(m)\mathbf{e}}, \quad (12)$$

with $\mathbf{e}^T\Phi_{yy}(m)\mathbf{e}$ and $\mathbf{e}^T\Phi_{nn}(m)\mathbf{e}$ corresponding to the noisy speech PSD $\phi_Y(m)$ and the noise PSD $\phi_N(m)$, respectively. Using (6), it can be easily shown that

$$\phi_Y(m)\gamma_y(m) = \phi_X(m)\gamma_x(m) + \phi_N(m)\gamma_n(m), \quad (13)$$

such that

$$\gamma_x(m) = \frac{\xi(m)+1}{\xi(m)}\gamma_y(m) - \frac{1}{\xi(m)}\gamma_n(m), \quad (14)$$

where $\xi(m) = \frac{\phi_X(m)}{\phi_N(m)}$ denotes the a-priori SNR.

3 MFMVDR Filter

In this section, we briefly review the single-channel MFMVDR filter proposed in [3, 4] and the estimation of the noisy speech correlation matrix and the speech correlation vector proposed in [8].

3.1 Optimization Problem

The MFMVDR filter aims at minimizing the total signal output power while not distorting the correlated speech component, i.e.,

$$\min_{\mathbf{h}(m)} \mathbf{h}^H(m)\Phi_{yy}(m)\mathbf{h}(m), \text{ s.t. } \mathbf{h}^H(m)\gamma_x(m) = 1. \quad (15)$$

Solving this optimization problem yields the MFMVDR filter [3, 4]

$$\mathbf{h}_{\text{MFMVDR}}(m) = \frac{\Phi_{yy}^{-1}(m)\gamma_x(m)}{\gamma_x^H(m)\Phi_{yy}^{-1}(m)\gamma_x(m)} \quad (16)$$

with the signal output power $\phi_Y^{\text{out}}(m) = \mathbb{E}[|\mathbf{h}_{\text{MFMVDR}}^H(m)\mathbf{y}(m)|^2]$, i.e.,

$$\phi_Y^{\text{out}}(m) = \frac{1}{\gamma_x^H(m)\Phi_{yy}^{-1}(m)\gamma_x(m)}. \quad (17)$$

As can be seen from (16), the MFMVDR filter requires an estimate of the speech correlation vector $\gamma_x(m)$ and the noisy speech correlation matrix $\Phi_{yy}(m)$. Typically, both quantities are highly time-varying, making it difficult to accurately estimate especially the speech correlation vector when only having access to the noisy speech signal.

3.2 Estimation of the Noisy Speech Correlation Matrix

Estimating the noisy speech correlation matrix $\Phi_{yy}(m)$ can be performed by applying first-order recursive smoothing with smoothing parameter α_y , i.e.,

$$\hat{\Phi}_{yy}(m) = \alpha_y\hat{\Phi}_{yy}(m-1) + (1-\alpha_y)\mathbf{y}(m)\mathbf{y}^H(m). \quad (18)$$

To improve the robustness of the MFMVDR filter, we apply diagonal loading with a regularization parameter of 0.04 before computing the inverse matrix, as suggested in [8].

3.3 Estimation of the Speech Correlation Vector

Based on (9) and (14), the optimal estimate of the speech correlation vector $\hat{\gamma}_x^{\text{Opt}}(m)$ is given by

$$\hat{\gamma}_x^{\text{Opt}}(m) = \frac{\hat{\Phi}_{xx}(m)\mathbf{e}}{\mathbf{e}^T\hat{\Phi}_{xx}(m)\mathbf{e}}, \quad (19a)$$

$$= \frac{\hat{\xi}^{\text{Opt}}(m)+1}{\hat{\xi}^{\text{Opt}}(m)}\hat{\gamma}_y(m) - \frac{1}{\hat{\xi}^{\text{Opt}}(m)}\hat{\gamma}_n(m), \quad (19b)$$

with $\hat{\Phi}_{xx}(m) = \hat{\Phi}_{yy}(m) - \hat{\Phi}_{nn}(m)$ the estimated speech correlation matrix, where the noise correlation matrix $\hat{\Phi}_{nn}(m)$ is computed similarly as in (18). The estimates of the noisy speech correlation vector $\hat{\gamma}_y(m)$ and the noise correlation vector $\hat{\gamma}_n(m)$ can be computed similarly as in (19a) and the optimal estimate of the a-priori SNR is given by $\hat{\xi}^{\text{Opt}}(m) = (\mathbf{e}^T\hat{\Phi}_{xx}(m)\mathbf{e})/(\mathbf{e}^T\hat{\Phi}_{nn}(m)\mathbf{e})$.

Similarly to (19), the ML estimate of the speech correlation vector $\gamma_x(m)$ proposed in [8] is given by

$$\hat{\gamma}_x^{\text{ML}}(m) = \frac{\hat{\xi}(m)+1}{\hat{\xi}(m)}\hat{\gamma}_y(m) - \frac{1}{\hat{\xi}(m)}\boldsymbol{\mu}_{\gamma_n}, \quad (20)$$

with $\hat{\xi}(m)$ an estimate of the a-priori SNR. In comparison to (19), the estimated noise correlation vector $\hat{\gamma}_n(m)$ is assumed to be constant for all time-frequency points, such that it can be replaced by its mean value $\boldsymbol{\mu}_{\gamma_n}$, which is determined by the frame overlap and the STFT analysis window [8]. To estimate the a-priori SNR $\hat{\xi}(m)$ we use the noise PSD estimator proposed in [14], i.e.,

$$\hat{\phi}_N(m) = \min[\hat{\phi}_Y(m), \hat{\phi}_N(m-1)](1+\nu), \quad (21)$$

where the parameter ν is set to 5 dB/s as in [8]. For the speech PSD we use the ML estimator proposed in [15], i.e.,

$$\hat{\phi}_X(m) = \max[\hat{\phi}_Y(m) - \hat{\phi}_N(m), 0]. \quad (22)$$

It should be noted that the ML estimate of the speech correlation vector strongly depends on the a-priori SNR estimate $\hat{\xi}(m)$. Especially for low a-priori SNRs, the ML estimate may become very large, such that the estimation error between $\gamma_x(m)$ and $\hat{\gamma}_x^{\text{ML}}(m)$ may become very large. This may cause the ML-based MFMVDR filter to result in unpleasant artifacts in the background noise or introduce speech distortion [8, 16].

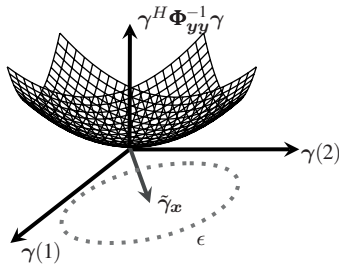


Figure 1: Quadratic cost function in (25) with exemplary presumed speech correlation vector $\tilde{\gamma}_x$ and bound ϵ .

4 Robust Constrained MFMVDR Filter

In [9] we showed that the performance of the MFMVDR filter is very sensitive to estimation errors of the speech correlation vector, such that correlated speech components may be mistakenly interpreted as uncorrelated and be suppressed rather than preserved. Inspired by the robust MVDR beamformer in [10], in this section we propose to estimate the speech correlation vector as the vector maximizing the total signal output power of the MFMVDR filter within a spherical uncertainty set. For conciseness, also the index m is omitted in this section. It should be noted that the calculations are performed for each frequency-bin k and time-frame m .

Given a presumed speech correlation vector $\tilde{\gamma}_x$, e.g., the ML estimate $\hat{\gamma}_x^{\text{ML}}$ in (20), the mismatch vector to the (unknown) speech correlation vector γ_x is defined as $\delta_x = \gamma_x - \tilde{\gamma}_x$, with $\epsilon_x = \|\gamma_x - \tilde{\gamma}_x\|_2^2$. We now define the spherical uncertainty set comprising all vectors whose squared distance to the presumed speech correlation vector $\tilde{\gamma}_x$ is smaller than or equal to a bound ϵ , i.e.,

$$\Gamma = \left\{ \gamma = \tilde{\gamma}_x + \delta \mid \|\delta\|_2^2 \leq \epsilon \right\}. \quad (23)$$

Similarly to the method proposed in [10] to robustly estimate the steering vector for the MVDR beamformer, the RC speech correlation vector is computed as the vector maximizing the total signal output power of the MFMVDR filter in (17) within the spherical uncertainty set in (23), i.e.,

$$\hat{\gamma}_x^{\text{RC}} = \underset{\gamma}{\operatorname{argmax}} \min_h h^H \Phi_{yy} h, \text{ s.t. } h^H \gamma = 1, \quad (24)$$

$$\|\gamma - \tilde{\gamma}_x\|_2^2 \leq \epsilon.$$

Using (17), this optimization problem can be reformulated as

$$\hat{\gamma}_x^{\text{RC}} = \underset{\gamma}{\operatorname{argmin}} \gamma^H \Phi_{yy}^{-1} \gamma, \text{ s.t. } \|\gamma - \tilde{\gamma}_x\|_2^2 \leq \epsilon \quad (25)$$

For $L = 2$, the quadratic cost function $\gamma^H \Phi_{yy}^{-1} \gamma$ in (25) is visualized in Figure 1 for an exemplary noisy speech correlation matrix Φ_{yy} , together with an exemplary presumed speech correlation vector $\tilde{\gamma}_x$ and bound ϵ . Obviously, the bound ϵ in (25) plays an important role and should be chosen in accordance with the accuracy of the presumed speech correlation vector $\tilde{\gamma}_x$, i.e., if $\|\gamma_x - \tilde{\gamma}_x\|_2^2$ is small, then ϵ should be small, whereas if $\|\gamma_x - \tilde{\gamma}_x\|_2^2$ is large, then ϵ should be large.

In order to avoid the (undesired) solution $\hat{\gamma}_x^{\text{RC}} = \mathbf{0}$, the bound ϵ should be chosen such that

$$\epsilon < \|\tilde{\gamma}_x\|_2^2. \quad (26)$$

Under this condition and considering the convex nature of the quadratic cost function in (25), the inequality constraint in (25) can be replaced by an equality constraint, i.e.,

$$\hat{\gamma}_x^{\text{RC}} = \underset{\gamma}{\operatorname{argmin}} \gamma^H \Phi_{yy}^{-1} \gamma, \text{ s.t. } \|\gamma - \tilde{\gamma}_x\|_2^2 = \epsilon. \quad (27)$$

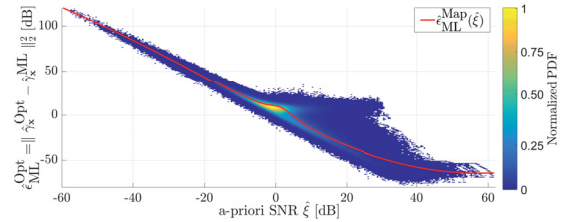


Figure 2: Normalized joint PDF of the optimal uncertainty parameter $\hat{\epsilon}_{\text{ML}}^{\text{Opt}}$ and the estimated a-priori SNR $\hat{\xi}$.

This constrained optimization problem can be solved using the method of Lagrange multipliers. The Lagrangian function is given by

$$f(\gamma, \lambda) = \gamma^H \Phi_{yy}^{-1} \gamma + \lambda \left(\|\gamma - \tilde{\gamma}_x\|_2^2 - \epsilon \right), \quad (28)$$

with λ the Lagrange multiplier. Setting the gradient of $f(\gamma, \lambda)$ with respect to γ equal to zero and applying the matrix inversion lemma, we obtain the RC speech correlation vector $\hat{\gamma}_x^{\text{RC}}(\lambda)$ as

$$\hat{\gamma}_x^{\text{RC}}(\lambda) = \tilde{\gamma}_x - (\lambda \Phi_{yy} + \mathbf{I})^{-1} \tilde{\gamma}_x \quad (29)$$

with \mathbf{I} denoting the $L \times L$ -dimensional identity matrix. Setting the derivative of $f(\gamma, \lambda)$ with respect to λ equal to zero and substituting (29) results in

$$\frac{\partial f(\gamma, \lambda)}{\partial \lambda} = \|\lambda \Phi_{yy} + \mathbf{I}\|^{-1} \tilde{\gamma}_x\|_2^2 - \epsilon = 0. \quad (30)$$

Let the eigenvalue decomposition of the noisy speech correlation matrix be given by

$$\Phi_{yy} = \mathbf{Q} \mathbf{U} \mathbf{Q}^H, \quad (31)$$

where the columns of \mathbf{Q} contain the orthogonal eigenvectors of Φ_{yy} and the diagonal elements of \mathbf{U} are the corresponding eigenvalues, with $u_0 \geq u_1 \geq \dots \geq u_{L-1}$. In addition, let

$$\mathbf{z}_x = \mathbf{Q}^H \tilde{\gamma}_x, \quad (32)$$

where $\mathbf{z}_x(l)$ denotes the l th element of \mathbf{z}_x . Using (31) and (32) in (30), we obtain

$$g(\lambda) = \sum_{l=0}^{L-1} \frac{|\mathbf{z}_x(l)|^2}{(1 + \lambda u_l)^2} = \epsilon \quad (33)$$

This non-linear equation in the Lagrange multiplier λ can be solved, e.g., using Newton's method. The solution is then used in (29), yielding the RC speech correlation vector $\hat{\gamma}_x^{\text{RC}}$. Since the condition in (9) may not be satisfied, possibly resulting in a scaling inaccuracy, a normalization is performed by dividing $\hat{\gamma}_x^{\text{RC}}$ with its first element. Using the normalized RC speech correlation vector $\hat{\gamma}_x^{\text{RC}}$ in (16) results in the RC MFMVDR filter.

As already mentioned, the bound ϵ should be chosen in accordance with the accuracy of the presumed speech correlation vector, e.g., assuming that the optimal bound is equal to $\hat{\epsilon}^{\text{Opt}} = \|\hat{\gamma}_x^{\text{Opt}} - \tilde{\gamma}_x\|_2^2$ with the optimal estimate of the speech correlation vector defined in (19). When using the ML estimate $\hat{\gamma}_x^{\text{ML}}$ as the presumed speech correlation vector $\tilde{\gamma}_x$, simulations have shown that the accuracy of the ML estimate (not unexpectedly) depends on the a-priori SNR estimate $\hat{\xi}$. Hence, we propose to train a mapping function $\hat{\epsilon}_{\text{ML}}^{\text{Map}}(\hat{\xi})$ between the a-priori SNR estimate $\hat{\xi}$, computed using (21) and (22), and the optimal bound $\hat{\epsilon}_{\text{ML}}^{\text{Opt}} = \|\hat{\gamma}_x^{\text{Opt}} - \hat{\gamma}_x^{\text{ML}}\|_2^2$. Figure 2 shows the normalized joint probability density function (PDF) of the optimal bound $\hat{\epsilon}_{\text{ML}}^{\text{Opt}}$ and the a-priori SNR estimate $\hat{\xi}$ for a wide range of speech and noise signals (30 TIMIT sentences [17], speech-shaped noise, two traffic and babble noise signals), for a broadband SNR range of 0 to 15 dB. It can be observed that with increasing a-priori SNR the optimal bound decreases. The mapping function $\hat{\epsilon}_{\text{ML}}^{\text{Map}}(\hat{\xi})$ (shown in red in Figure 2) is based on the maximum value of the normalized PDF for each a-priori SNR estimate $\hat{\xi}$.

5 Evaluation

In this section, we compare the estimation accuracy error of the proposed RC speech correlation vector in (29) with the ML speech correlation vector in (20). Furthermore, we evaluate the perceptual quality of both MFMVDR filters, i.e., the proposed RC MFMVDR filter and the ML MFMVDR filter.

5.1 STFT Framework

Similarly as in [8], to increase the exploitable speech correlation across time-frames, we use a highly temporally resolved STFT framework with a frame length of 4 ms and an overlap of 75 %, resulting in a frame shift of 1 ms. As the STFT analysis and synthesis window we use a square-root Hann window. The sampling frequency is 16 kHz. The number of consecutive time-frames is set to $L=18$, resulting in 21 ms of data used in each filtering operation. The smoothing parameter in (18) is set to $\alpha_y = 0.9$, resulting in a smoothing window of 10 ms.

5.2 Evaluation of the Speech Correlation Vector Estimation

The performance of both estimated speech correlation vectors, i.e., the RC estimator in (29) and the ML estimator in (20), is evaluated in terms of the mean-square error (MSE) between the optimal speech correlation vector in (19) and the estimated speech correlation vector, i.e.,

$$\text{MSE} = \frac{1}{\mathbb{F}} \sum_{k,m \in \mathbb{F}} \|\hat{\gamma}_{\mathbf{x}}^{\text{Opt}}(k,m) - \hat{\gamma}_{\mathbf{x}}(k,m)\|_2^2, \quad (34)$$

where \mathbb{F} is the set of time-frequency points that contain either noise-only or speech-and-noise (i.e., with a-priori SNR $\hat{\xi}^{\text{Opt}}(m)$ larger than -5 dB). Furthermore, we classify time-frequency points whose squared error is larger than 200 as outliers and exclude them from the MSE calculation. As clean speech signals, we use 60 sentences from the TIMIT database [17], spoken by different speakers (10 male, 10 female). As noise signals, we use speech-shaped noise, traffic and babble noise signals. The considered SNR range is -5 dB to 20 dB. We make sure that the evaluation data differs from the training set.

For different SNRs, Figure 3 depicts the performance (averaged over all combinations of speech and noise signals) in terms of the MSE and the percentage of outliers in speech-and-noise and noise-only points. It can be observed that compared to the ML estimate the RC estimate yields a considerably lower MSE. Moreover, the outliers are completely removed, indicating that the RC estimate is more accurate and much stabler than the ML estimate.

5.3 Perceptual Evaluation

The perceptual evaluation is performed using a pairwise preference test, where the quality of both MFMVDR filters, i.e., the proposed RC MFMVDR filter and the ML MFMVDR filter, and the unprocessed noisy speech signal is compared to each other. As speech signals, we used two TIMIT sentences, spoken by a female and a male speaker. As noise signals, we used speech-shaped noise, traffic and babble noise at an SNR of 0 dB. For each of the three noise types 15 self-reported normal hearing listeners aged between 21 and 38 were asked to judge which of the two presented signals is preferred in terms of speech quality (SQ), noise reduction (NR) and overall quality (OQ). The order of the signal pairs as well as the order of the noise types was randomized. As a reference, the clean speech signal was also available. The signals were presented diotically over Sennheiser HDA 200 headphones using an RME Babyface sound card in a quiet office. Similarly as in [18], the ratings are expressed in the percentage of times each algorithm was preferred, i.e., a value of 100% indicates that the algorithm won every pairwise comparison it was involved in. To determine statistical significance, we used the non-parametric Friedman test followed by a Wilcoxon signed-rank test with Bonferroni correction. We assume statistical significance for $p < 0.05$.

Since the ratings between the different noise types were not statistically significant, Figure 4 shows the preference results

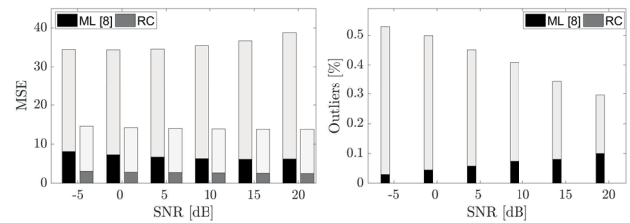


Figure 3: Average MSE and percentage of outliers for the ML and RC speech correlation vectors for different SNRs. The lower and upper parts of the bars represent the performance in speech-and-noise and noise-only, respectively.

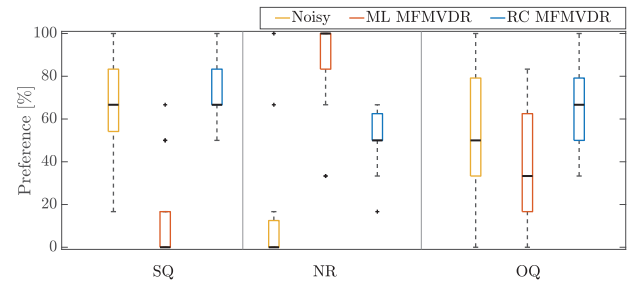


Figure 4: Boxplots of the average preference ratings of the unprocessed noisy speech signal, the ML MFMVDR filter and the proposed RC MFMVDR filter for the criteria speech quality, noise reduction and overall quality.

averaged over all three noise types for the three evaluation criteria. In terms of the speech quality, the proposed RC MFMVDR filter and the unprocessed noisy speech signal are preferred over the ML MFMVDR filter, while there is no statistically significant difference between the RC MFMVDR filter and the noisy speech signal. In terms of the noise reduction, the ML MFMVDR filter is preferred over the RC MFMVDR filter but the RC MFMVDR filter is still preferred over the noisy speech signal. In terms of the overall quality, the RC MFMVDR filter is preferred over the ML MFMVDR filter and the noisy signal, while there is no statistically significant difference between the ML MFMVDR filter and the noisy speech signal.

The results indicate that the RC MFMVDR filter leads to less speech distortion and a more natural speech quality than the ML MFMVDR filter. Although the RC MFMVDR filter is more conservative in suppressing the background noise than the ML MFMVDR filter, the proposed RC approach leads to less musical noise such that the overall signal quality is significantly preferred over the ML MFMVDR filter.

6 Conclusions

In this paper, we proposed a robust constrained (RC) multi-frame minimum variance distortionless response (MFMVDR) filter for single-channel speech enhancement. Inspired by robust beamforming approaches, we proposed to estimate the speech correlation vector as the vector maximizing the total signal output power within a spherical uncertainty set. The spherical uncertainty set imposes an upper bound on the norm of the mismatch vector between the speech correlation vector and the presumed speech correlation vector. We proposed to set this bound by training a mapping function depending on the a-priori SNR estimate. Simulation results show that the proposed RC approach leads to a more accurate and stable estimate of the speech correlation vector than the ML approach. Furthermore, a pairwise preference test indicates that although the proposed RC MFMVDR filter leads to a more conservative noise reduction than the ML MFMVDR filter, it produces less speech and noise distortions, such that the speech sounds more natural and less musical noise is present.

References

- [1] R. C. Hendriks, T. Gerkmann, and J. Jensen, *DFT-Domain Based Single-Microphone Noise Reduction for Speech Enhancement: A Survey of the State of the Art*. Morgan & Claypool, 2013.
- [2] J. Benesty, J. Chen, and E. A. Habets, *Speech enhancement in the STFT domain*. Springer Science & Business Media, 2011.
- [3] J. Benesty and Y. Huang, “A single-channel noise reduction MVDR filter,” in *Proc. IEEE Int. Conf. Acoust., Speech, Signal Process. (ICASSP)*, (Prague, Czech Republic), pp. 273–276, May 2011.
- [4] Y. Huang and J. Benesty, “A multi-frame approach to the frequency-domain single-channel noise reduction problem,” *IEEE Trans. Audio, Speech, Language Process.*, vol. 20, pp. 1256–1269, May 2012.
- [5] K. T. Andersen and M. Moonen, “Robust speech-distortion weighted interframe Wiener filters for single-channel noise reduction,” *IEEE Trans. Audio, Speech, Language Process.*, vol. 26, pp. 97–107, Jan. 2018.
- [6] B. D. Van Veen and K. M. Buckley, “Beamforming: A versatile approach to spatial filtering,” *IEEE ASSP Mag.*, vol. 5, pp. 4–24, Apr. 1988.
- [7] S. Doclo, W. Kellermann, S. Makino, and S. E. Nordholm, “Multichannel signal enhancement algorithms for assisted listening devices: Exploiting spatial diversity using multiple microphones,” *IEEE Signal Processing Magazine*, vol. 32, pp. 18–30, Mar. 2015.
- [8] A. Schasse and R. Martin, “Estimation of subband speech correlations for noise reduction via MVDR processing,” *IEEE Trans. Audio, Speech, Language Process.*, vol. 22, pp. 1355–1365, Sept. 2014.
- [9] D. Fischer and S. Doclo, “Sensitivity analysis of the multi-frame MVDR filter for single-microphone speech enhancement,” in *Proc. of Europ. Signal Process. Conf. (EUSIPCO)*, (Kos, Greece), pp. 603–607, Aug. 2017.
- [10] J. Li, P. Stoica, and Z. Wang, “On robust Capon beamforming and diagonal loading,” *IEEE Trans. Signal Process.*, vol. 51, pp. 1702–1715, July 2003.
- [11] A. Khabbazibasmenj, S. A. Vorobyov, and A. Hassanien, “Robust adaptive beamforming based on steering vector estimation with as little as possible prior information,” *IEEE Trans. Signal Process.*, vol. 60, pp. 2974–2987, Feb. 2012.
- [12] S. Vorobyov, “Principles of minimum variance robust adaptive beamforming design,” *IEEE Trans. Signal Process.*, vol. 93, pp. 3264–3277, Dec. 2013.
- [13] Y. Zhao, J. R. Jensen, M. G. Christensen, S. Doclo, and J. Chen, “Experimental study of robust beamforming techniques for acoustic applications,” in *Proc. IEEE Workshop Appl. Signal Process. Audio, Acoust. (WASPAA)*, (New Paltz, NY, USA), pp. 86–90, 2017.
- [14] E. Hänsler and G. Schmidt, *Acoustic echo and noise control: a practical approach*, vol. 40. John Wiley & Sons, 2005.
- [15] R. McAulay and M. Malpass, “Speech enhancement using a soft-decision noise suppression filter,” *IEEE Trans. Acoust., Speech, Signal Process.*, vol. 28, pp. 137–145, Apr. 1980.
- [16] D. Fischer, S. Doclo, E. A. P. Habets, and T. Gerkmann, “Combined single-microphone Wiener and MVDR filtering based on speech interframe correlations and speech presence probability,” in *Proc. ITG Conference on Speech Commun.*, (Paderborn, Germany), pp. 292–296, Oct. 2016.
- [17] J. S. Garofolo, “DARPA TIMIT acoustic-phonetic speech database,” in *National Institute of Standards and Technology (NIST)*, 1988.
- [18] M. Krawczyk-Becker and T. Gerkmann, “An evaluation of the perceptual quality of phase-aware single-channel speech enhancement,” *J. Acoust. Soc. Amer.*, vol. 140, pp. EL364–EL369, Oct. 2016.