

DELFT UNIVERSITY OF TECHNOLOGY

EE4715 ARRAY PROCESSING

Speech Enhancement

Authors:

Weijia Yi (5511127) & Xuan Gao (5232481)

April 28, 2023



1 Objective

This project aims to design and build a multi-microphone speech enhancement/beamforming system for far-end noise reduction. It contains following requirements:

- Generate signals according to the signal model discussed in class using the given audio files and impulse responses.
- Build a multi-microphone speech enhancement system for noise reduction.
- Perform an evaluation of the speech enhancement system.

There are five speech signals, containing two clean signals and three different types of noise signals (babble noise, artificial non-stationary noise, and stationary speech shaped noise). The impulse responses have five sets and each set contains the impulse responses from a source location to the four microphone locations. All datasets are provided from array processing course: [Project 2: Speech enhancement](#). The references of the implementation algorithm are course slides [1]. The evaluation method is referred to [2].

2 Signal model

The received signals by a single microphone can be modeled as

$$x[n] = \sum_{i=1}^d (h_i * s)[n] + n[n].$$

The received speech signals with noise are considered as the sum of additive noise signals n and the convolution of the input clean signal s and the impulse responses h of different channels. Take one channel as the target signal (in practical the early reflection), while signals from other channels are all considered as interference, then the model can be rewritten as

$$x[n] = \underbrace{\sum_{i=1}^d (h_{target} * s_i)[n]}_{\text{target}} + \underbrace{\sum_{i=1}^d (h_{inter} * s_i)[n]}_{\text{interference}} + n[n]. \quad (1)$$

The additive noise utilized in this project includes artificial non-stationary noise, babble noise and speech-shaped noise. Since the artificial noise signal is not perceivable, we amplify it to make it detectable. The noisy speech signals are showed in [Figure 1](#).

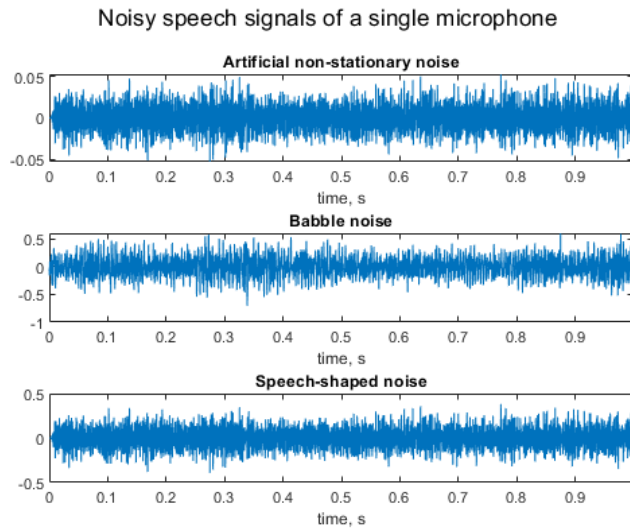


Figure 1: Noisy speech signals of a single microphone with 3 types of additive noise

3 Multi-microphone speech enhancement system implement

3.1 STEP 1: short-time stationary Fourier transform (STFT)

As speech signals change over time, the spectral and temporal characteristics of the speech waveform are time-varying. Therefore, only short segments of speech waveform can be assumed to have similar acoustic properties. This is what it means by "short-term stationary". Speech signals are typically assumed stationary over 20-30ms.

The processing is done in STFT domain, which we do fast Fourier transform (FFT) on short windowed time frames. All time frames obey short time wide-sense stationary (WSS) assumption.

For M microphones, speech signals with noise and noise signals are stacked vector notation as Equation 2. In our project, we use hamming window of the size 320. The number of frames is 6718 and the size of each frame is 320.

$$\begin{aligned}\mathbf{x}(k, l) &= \underbrace{\mathbf{a}_1(k, l)s_1(k, l)}_{\text{target}} + \underbrace{\sum_{i=2}^d \mathbf{a}_i(k, l)s_i(k, l) + \mathbf{n}'(k, l)}_{\text{interference+noise}} \\ &= \mathbf{a}(k, l)s(k, l) + \mathbf{n}(k, l)\end{aligned}\quad (2)$$

3.2 STEP 2: Estimation of \mathbf{R}_s using GEVD

Firstly, we estimate \mathbf{R}_x and \mathbf{R}_n , then we use generalized eigenvalue decomposition (GEVD) to estimate \mathbf{R}_s . The autocorrelation sequence r_x is calculated as $r_x(k, l) = E\{x(k)x^*(l)\}$. Therefore,

$$r_x = [r_x(-3), r_x(-2), r_x(-1), r_x(0), r_x(1), r_x(2), r_x(3)]$$

The autocorrelation matrix \mathbf{R}_x is a Hermitian and Toeplitz matrix and shown in Equation 3.

$$\mathbf{R}_x = E\{\mathbf{x}\mathbf{x}^H\} = \begin{bmatrix} r_x(0) & r_x(1) & r_x(2) & r_x(3) \\ r_x(1) & r_x(0) & r_x(1) & r_x(2) \\ r_x(2) & r_x(1) & r_x(0) & r_x(1) \\ r_x(3) & r_x(2) & r_x(1) & r_x(0) \end{bmatrix}\quad (3)$$

The autocorrelation sequence is $r_n = [r_n(-3), r_n(-2), r_n(-1), r_n(0), r_n(1), r_n(2), r_n(3)]$. The autocorrelation matrix \mathbf{R}_n is a Hermitian and Toeplitz matrix and shown in Equation 4.

$$\mathbf{R}_n = E\{\mathbf{n}\mathbf{n}^H\} = \begin{bmatrix} r_n(0) & r_n(1) & r_n(2) & r_n(3) \\ r_n(1) & r_n(0) & r_n(1) & r_n(2) \\ r_n(2) & r_n(1) & r_n(0) & r_n(1) \\ r_n(3) & r_n(2) & r_n(1) & r_n(0) \end{bmatrix}\quad (4)$$

Next, we apply GEVD to \mathbf{R}_x and \mathbf{R}_n . As $\mathbf{R}_n = \mathbf{U}^{-H}\mathbf{\Lambda}_n\mathbf{U}^{-1}$ and $\mathbf{R}_x = \mathbf{U}^{-H}\mathbf{\Lambda}_x\mathbf{U}^{-1}$, therefore,

$$\mathbf{U}^H\mathbf{R}_x\mathbf{U} = \text{diag}(a_1, \dots, a_n) \quad \text{and} \quad \mathbf{U}^H\mathbf{R}_n\mathbf{U} = \text{diag}(b_1, \dots, b_n)\quad (5)$$

Hence, we have $\mathbf{R}_n\mathbf{U} = \mathbf{U}^{-H}\mathbf{\Lambda}_{R_n}$ so that

$$\mathbf{R}_x\mathbf{U} = \mathbf{U}^{-H}\mathbf{\Lambda}_{R_x} = \mathbf{U}^{-H}\mathbf{\Lambda}_{R_n}\mathbf{\Lambda}_{R_n}^{-1}\mathbf{\Lambda}_{R_x} = \mathbf{R}_n\mathbf{U}\mathbf{\Lambda}_{n_x}\quad (6)$$

$$\mathbf{R}_n^{-1}\mathbf{R}_x = \mathbf{U}\mathbf{\Lambda}_{n_x}\mathbf{U}^{-1}\quad (7)$$

Hence, the generalised eigenvalues and eigenvectors of \mathbf{R}_x and \mathbf{R}_n are the eigenvalues and eigenvectors of

the matrix $\mathbf{R}_n^{-1}\mathbf{R}_x$.

Since Equation 8, we subtract 1 from the diagonal values of $\mathbf{\Lambda}_{nx}$ to get $\mathbf{\Lambda}_{ns}$.

$$\mathbf{R}_n^{-1}\mathbf{R}_x = \mathbf{R}_n^{-1}(\mathbf{R}_n + \mathbf{R}_s) = \mathbf{I} + \mathbf{R}_n^{-1}\mathbf{R}_s = \mathbf{U}\mathbf{\Lambda}_{nx}\mathbf{U}^{-1} = \mathbf{U}(\mathbf{\Lambda}_{ns} + \mathbf{I})\mathbf{U}^{-1} \quad (8)$$

After we get $\mathbf{\Lambda}_{ns}$, \mathbf{R}_s is calculated by Equation 9:

$$\begin{aligned} \mathbf{R}_n^{-1}\mathbf{R}_s &= \mathbf{U}\mathbf{\Lambda}_{ns}\mathbf{U}^{-1} \\ \mathbf{R}_s &= \mathbf{R}_n\mathbf{U}\mathbf{\Lambda}_{ns}\mathbf{U}^{-1} \end{aligned} \quad (9)$$

3.3 STEP 3: Beamforming

To solve out multi-channel wiener beamformer and MVDR beamformer, we first derive the expression of signal distortion and residual noise variance by the mean square error of the beamformer output $\mathbf{w}^H\mathbf{x}$ signals and the desired target signal at the reference microphone s_1 :

$$\begin{aligned} \mathbb{E}|\mathbf{w}^H\mathbf{x} - s_1|^2 &= \mathbb{E}|\mathbf{w}^H\mathbf{s} + \mathbf{w}^H\mathbf{n} - s_1|^2 \\ &= \mathbb{E}|\mathbf{w}^H\mathbf{s} - s_1|^2 + \mathbb{E}|\mathbf{w}^H\mathbf{n}|^2 \end{aligned}$$

The first term is the signal distortion and the second is the residual noise variance. First begin with the constraint optimization problem to minimize the distortion constrained by reducing noise power

$$\begin{aligned} &\text{minimise} \quad \mathbb{E}|\mathbf{w}^H\mathbf{s} - s_1|^2 \\ &\text{subject to} \quad \mathbb{E}|\mathbf{w}^H\mathbf{n}|^2 \leq c, \end{aligned}$$

where $0 \leq c \leq \sigma_{n_1}^2$, $\sigma_{n_1}^2$ is the noise variance at the reference microphone. The solution of this optimization problem is

$$\mathbf{w}^* = (\mathbf{R}_s + \mu\mathbf{R}_n)^{-1} \mathbf{R}_s\mathbf{e}_1, \quad (10)$$

where μ is the Lagrange multiplier.

- **Multi-Channel Wiener beamformer**

When $\mu = 1$, the classical multi-channel Wiener filter is given as

$$\mathbf{w}_{MWF} = \mathbf{R}_x^{-1}\mathbf{R}_s\mathbf{e}_1 \quad (11)$$

where $\mathbf{e}_1 = [1, 0, \dots, 0]^T$.

- **MVDR beamformer**

When $\mu = 0$, it leads to

$$\mathbf{w}_{MVDR} = \frac{\mathbf{R}_n^{-1}\mathbf{a}}{\mathbf{a}^H\mathbf{R}_n^{-1}\mathbf{a}} \quad (12)$$

\mathbf{a} is one column of \mathbf{Q} ($\mathbf{Q} = \mathbf{U}^{-H}$) corresponding to the largest eigenvalue.

Therefore, source signal is reconstructed as Equation 13:

$$\mathbf{s} = \mathbf{w}^H\mathbf{x} \quad (13)$$

3.4 STEP 4: Inverse STFT (ISTFT)

The source signal \mathbf{s} is still in Fourier domain. To convert it back to time domain, ISTFT is applied to \mathbf{s} .

4 Evaluation

4.1 Speech intelligibility in bits

Speech intelligibility in bits (SIIB) is an estimate of the amount of information shared between a talker and a listener in bits per second [2]. The communication process consisting of a message $\{\mathbf{M}_t\}$, speech signal $\{\mathbf{X}_t\}$, and degraded speech signal $\{\mathbf{Y}_t\}$ is described by a Markov chain as 14. $\{\mathbf{M}_t\} \rightarrow \{\mathbf{X}_t\}$ is the speech production channel and $\{\mathbf{X}_t\} \rightarrow \{\mathbf{Y}_t\}$ is the environmental channel.

$$\{\mathbf{M}_t\} \rightarrow \{\mathbf{X}_t\} \rightarrow \{\mathbf{Y}_t\} \quad (14)$$

The mutual information rate of the environmental channel is:

$$\begin{aligned} I(\{\mathbf{X}_t\}; \{\mathbf{Y}_t\}) &= \lim_{K \rightarrow \infty} \frac{1}{K} I(\mathbf{X}^K; \mathbf{Y}^K) \\ &= \lim_{K \rightarrow \infty} \frac{1}{K} I(\tilde{\mathbf{X}}^K; \tilde{\mathbf{Y}}^K) \\ &= \lim_{K \rightarrow \infty} \frac{1}{K} \sum_{j=1}^{KJ} I(\tilde{\mathbf{X}}_j^K; \tilde{\mathbf{Y}}_j^K) \end{aligned} \quad (15)$$

The information rate between $\{\mathbf{M}_t\}$ and $\{\mathbf{X}_t\}$ is:

$$\begin{aligned} I(\{\mathbf{M}_t\}; \{\mathbf{X}_t\}) &= \lim_{K \rightarrow \infty} \frac{1}{K} \sum_{j=1}^{KJ} I(\tilde{\mathbf{M}}_j^K; \tilde{\mathbf{X}}_j^K) \\ &= \lim_{K \rightarrow \infty} -\frac{1}{K} \sum_{j=1}^{KJ} \frac{1}{2} \log_2(1 - r_j^2) \end{aligned} \quad (16)$$

The mutual information rate of the speech production channel is:

$$I(\{\mathbf{M}_t\}; \{\mathbf{Y}_t\}) \leq \min(I(\{\mathbf{M}_t\}; \{\mathbf{X}_t\}), I(\{\mathbf{X}_t\}; \{\mathbf{Y}_t\})) \quad (17)$$

The proposed intelligibility metric combines 15, 16, and 17 to estimate the amount of information shared between $\{\mathbf{M}_t\}$ and $\{\mathbf{Y}_t\}$ in bits per second as shown in Equation 18.

$$\text{SIIB} = \frac{F}{K} \sum_{j=1}^{KJ} \min\left(-\frac{1}{2} \log_2(1 - r_j^2), I(\tilde{\mathbf{X}}_j^K; \tilde{\mathbf{Y}}_j^K)\right) \quad (18)$$

SIIB has high performance. Given a clean acoustic speech signal and a distorted acoustic speech signal, SIIB estimates the amount of information shared between two signals. In our project, we input the clean acoustic signal and the reconstructed acoustic signal to output the SIIB value.

4.2 Results and analysis

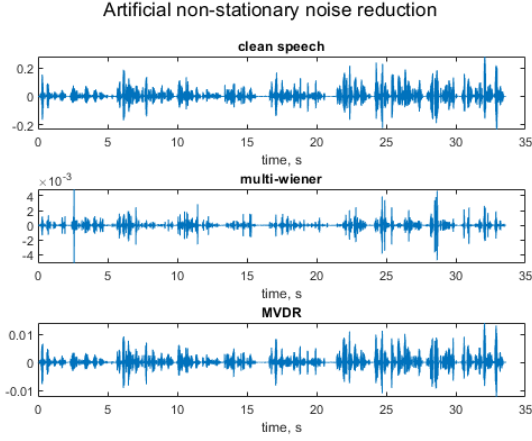


Figure 2: The comparison of artificial non-stationary noise reduction performance

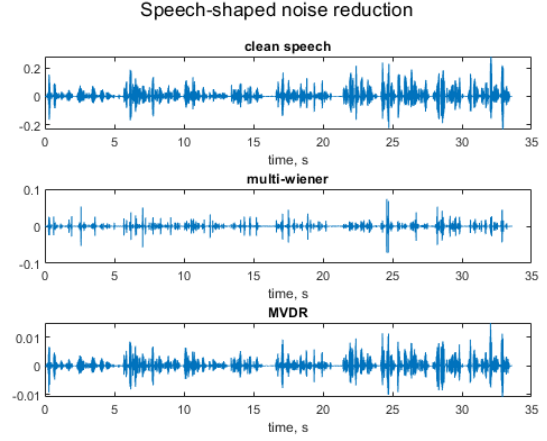


Figure 4: The comparison of speech-shaped noise reduction performance

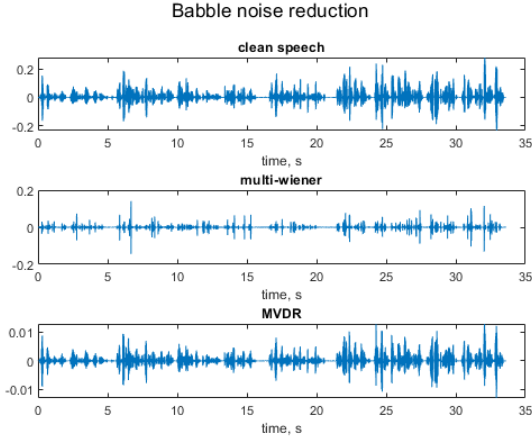


Figure 3: The comparison of babble noise reduction performance

Noise types	artificial	babble	speech-shaped
MVDR	100.5258	101.3844	101.2689
Multi-Wiener	19.0688	23.3159	29.2890

Table 1: Comparison of noise reduction effects of two types of beamformers by SIIB

There are two beamformers Multi-Wiener and MVDR, and three types of noise. Figure 2 to Figure 4 shows that Multi-Wiener beamformer causes more distortion on signal while MVDR beamformer may contain more noise but less distortion. Table 1 shows the comparison of noise reduction effects by SIIB. The bigger the value is, the more amount of information shared will be. In general, MVDR beamformer performs much better than Multi-Wiener beamformer in all types of noise. The value from MVDR beamformer is 4.37 times in average greater than the value gained from Multi-Wiener beamformer. For MVDR beamformer, it shows similar performance regardless of the noise type, as the value difference is less than 1 among all values. While for Multi-Wiener beamformer, it performs better in speech-shaped noise reduction and shows the least performance in artificial-nonstationary noise reduction. The SIIB value difference between the highest and smallest value of Multi-Wiener beamformer is 10.22. Therefore, reducing noise by MVDR beamformer leads to more understandable speech signals.

A Appendix - Create dataset

```
1 clear all
2
3 amp_arti = 8.8967;
4 filefolder = 'data\';
5 filenames = {'clean_speech.wav'; 'clean_speech_2.wav'; 'babble_noise.wav';...
6             'artificial_nonstat_noise.wav'; 'speech_shaped_noise.wav'};
7 datanames = {'clean_1'; 'clean_2'; 'babble';...
8             'artificial_nonstat'; 'speech_shaped'};
9 responses = {'h_target'; 'h_inter1'; 'h_inter2'; 'h_inter3'; 'h_inter4'};
10 for k = 1:length(filenames)
11     filepath = strcat(filefolder, filenames{k});
12     importfile(filepath)
13     assignin('base', datanames{k}, data);
14     clear data filepath
15 end
16
17 load('data\impulse_responses.mat')
18
19 h_tar=eval(char(responses(1)));
20 for M=1:4
21     for source=2:5
22         h_inter=eval(char(responses(source)));
23         s_2(M,:)=conv(h_target(M,:), clean_2);
24         s2(M,:)=s_2(M, size(s_2,2)-length(clean_2)+1:size(s_2,2));
25         inter_2(M,:,source)=conv(h_inter(M,:), clean_2);
26         inter2(M,:,source)=inter_2(M, size(inter_2,2)-length(clean_2)+1:size(inter_2,2), source);
27         clear h_inter
28     end
29     clear source
30     inter2=sum(inter2,3);
31     noisy_arti_2(M,:)=(artificial_nonstat(1:length(s2(M,:))) .* amp_arti)'+inter2(M,:);
32     noisy_babble_2(M,:)=(babble(1:length(s2(M,:)))'+inter2(M,:);
33     noisy_spee_2(M,:)=(speech_shaped(1:length(s2(M,:)))'+inter2(M,:);
34     x_arti_2(M,:)=noisy_arti_2(M,:)+s2(M,:);
35     x_babble_2(M,:)=noisy_babble_2(M,:)+s2(M,:);
36     x_spee_2(M,:)=noisy_spee_2(M,:)+s2(M,:);
37 end
38 clear M
39
40 save('data.mat', 'clean_2', 'x_babble_2',...
41     'x_arti_2', 'x_spee_2', 'fs',...
42     'noisy_arti_2', 'noisy_babble_2', 'noisy_spee_2')
43
44 function importfile(fileToRead1)
45 %IMPORTFILE(FILETOREAD1)
46 % Imports data from the specified file
47 % FILETOREAD1: file to read
48
49 newData1 = importdata(fileToRead1);
50
51 % Create new variables in the base workspace from those fields.
52 vars = fieldnames(newData1);
53 for i = 1:length(vars)
54     assignin('base', vars{i}, newData1.(vars{i}));
55 end
56 end
```

B Appendix - Speech enhancement system implement

```
1 clear all
2 load('data.mat')
3
```

```

4  frame_num=6718;
5  frame_size=320;
6
7  %% STFT
8  win = hamming(frame_size,'periodic');%16kHz,0.02s
9
10 % noise n, speech signals with noise x
11
12 for M=1:4
13     noisy_1(M, :, :)=stft(noisy_arti_2(M, :), fs, 'Window', win);
14     noisy_2(M, :, :)=stft(noisy_babble_2(M, :), fs, 'Window', win);
15     noisy_3(M, :, :)=stft(noisy_spee_2(M, :), fs, 'Window', win);
16     x_1(M, :, :)=stft(x_arti_2(M, :), fs, 'Window', win);
17     x_2(M, :, :)=stft(x_babble_2(M, :), fs, 'Window', win);
18     x_3(M, :, :)=stft(x_spee_2(M, :), fs, 'Window', win);
19 end
20
21 %% s estimation
22 [s_arti_wie, s_arti_mvdr]=beamformer(x_1, noisy_1, frame_num, frame_size);
23 [s_babble_wie, s_babble_mvdr]=beamformer(x_2, noisy_2, frame_num, frame_size);
24 [s_spee_wie, s_spee_mvdr]=beamformer(x_3, noisy_3, frame_num, frame_size);
25
26 %% ISTFT
27
28 esti_arti_1=istft(s_arti_wie, fs, 'Window', win);
29 esti_arti_2=istft(s_arti_mvdr, fs, 'Window', win);
30 esti_babble_1=istft(s_babble_wie, fs, 'Window', win);
31 esti_babble_2=istft(s_babble_mvdr, fs, 'Window', win);
32 esti_spee_1=istft(s_spee_wie, fs, 'Window', win);
33 esti_spee_2=istft(s_spee_mvdr, fs, 'Window', win);
34
35 save('output.mat', 'esti_arti_1', 'esti_arti_2', ...
36     'esti_babble_1', 'esti_babble_2', 'esti_spee_1', 'esti_spee_2')
37
38 %% beamformer func
39 function [s_wie, s_mvdr]=beamformer(x, n, frame_num, frame_size)
40 for k=1:frame_num
41     for l=1:frame_size
42 % Calculate Rn, Rx
43         rx_1=xcorr(x(:, l, k));
44         rx_1=rx_1(4:7);
45         Rx_1=toeplitz(rx_1);
46         clear rx_1
47
48         rn_1=xcorr(n(:, l, k));
49         rn_1=rn_1(4:7);
50         Rn_1=toeplitz(rn_1);
51         clear rn_1
52
53 % GEVD Rs
54         % A=Rx_1; B=Rn_1;
55         [U, Diag, Q]=eig(Rx_1, Rn_1);
56         Diag=Diag-eye(size(Diag, 1));
57         Rs_1=Rn_1*U*Diag/U;
58
59 % Calculate Wiener
60         e1=zeros(size(Rs_1, 2), 1);
61         e1(1)=1;
62         w1_wie=Rx_1\Rs_1*e1;
63 % Calculate MVDR
64         [~, index_row]=max(abs(diag(Diag)));
65         a=Q(:, index_row(1));
66         w1_mvdr=Rn_1\ a / (a' / Rn_1 * a);
67
68 % solve s
69         s_1_wie=w1_wie'*x(:, l, k);
70         s_1_mvdr=w1_mvdr'*x(:, l, k);
71         s_wie(l, k)=s_1_wie;
72         s_mvdr(l, k)=s_1_mvdr;
73     end

```



```
74 end
75 end
```

C Appendix - Evaluation

```
1 load output.mat
2 load data.mat
3
4 %% visualization
5 len=length(esti_spee_1);
6 clean=clean_2(1:len);
7 figure(1)
8 subplot(3,1,1); plot((1:length(clean))/fs,clean); title('clean speech')
9 xlabel('time, s')
10 subplot(3,1,2); plot((1:len)/fs,10*esti_arti_1); title('multi-wiener')
11 xlabel('time, s')
12 subplot(3,1,3); plot((1:len)/fs,10*esti_arti_2); title('MVDR')
13 xlabel('time, s')
14
15 figure(2)
16 subplot(3,1,1); plot((1:length(clean))/fs,clean); title('clean speech')
17 xlabel('time, s')
18 subplot(3,1,2); plot((1:len)/fs,10*esti_babble_1); title('multi-wiener')
19 xlabel('time, s')
20 subplot(3,1,3); plot((1:len)/fs,10*esti_babble_2); title('MVDR')
21 xlabel('time, s')
22
23 figure(3)
24 subplot(3,1,1); plot((1:length(clean))/fs,clean); title('clean speech')
25 xlabel('time, s')
26 subplot(3,1,2); plot((1:len)/fs,esti_spee_1); title('multi-wiener')
27 xlabel('time, s')
28 subplot(3,1,3); plot((1:len)/fs,esti_spee_2); title('MVDR')
29 xlabel('time, s')
30
31 %% SIIB_Gauss
32 siib_a_w = SIIB_Gauss(clean,esti_arti_1,fs);
33 siib_a_m = SIIB_Gauss(clean,esti_arti_2,fs);
34 siib_b_w = SIIB_Gauss(clean,esti_babble_1,fs);
35 siib_b_m = SIIB_Gauss(clean,esti_babble_2,fs);
36 siib_s_w = SIIB_Gauss(clean,esti_spee_1,fs);
37 siib_s_m = SIIB_Gauss(clean,esti_spee_2,fs);
```

References

- [1] R. C. Hendriks, *Microphone Array Processing*. TU Delft, June 2022.
- [2] S. Van Kuyk, W. B. Kleijn, and R. C. Hendriks, “An instrumental intelligibility metric based on information theory,” *IEEE Signal Processing Letters*, vol. 25, no. 1, pp. 115–119, 2017.