

DELFT UNIVERSITY OF TECHNOLOGY

ESTIMATION AND DETECTION
ET4386

Multi-Microphone Speech Enhancement

Authors:

Yanqi Hong (5884683) y.hong-8@student.tudelft.nl
Kaishen Lin (5971950) K.Lin-11@student.tudelft.nl

January 7, 2024



1 Introduction

Multi-Microphone Speech Enhancement is a crucial technology in audio processing, particularly for applications in noisy environments such as mobile phones, hearing aids, and speech recognition systems. The technique employs multiple microphones at different locations to enhance speech signals. This spatial arrangement allows for better differentiation between the desired speech and background noise, a challenge for single-microphone systems. By mathematically modeling the signals from various microphones, these systems effectively isolate and enhance speech, improving clarity and intelligibility. This technology is increasingly vital in our noise-filled world, enhancing the user experience in everyday applications and specialized devices.

Multi-channel speech enhancement is a complex problem that requires many factors to be considered, such as delay and damping. However, in this problem, we assume that a key factor acoustic transfer function(ATF) of the target source always equals 1. It's an idealized scenario where the target sound source is located in the far-field of the microphone array where the sound waves can be considered parallel (with minimal divergence). Meanwhile, the target source is directly facing the array, resulting in the same sound delay and damping at all microphones. This greatly simplifies the process of processing speech signals. In our project, we first divide the signal into frames using a 50% overlap, 20 ms duration and apply a Hann window to each frame, then perform a Fast Fourier Transform (FFT) on each frame. Then we use the time frames during the first second, as these are noise only, and apply Voice Activity Detection (VAD) to better extract noise characteristics. Finally, We use MLE when assuming the speech signal is deterministic and MMSE when assuming the speech signal is stochastic.

2 Signal Model and assessment metrics

2.1 Signal Model

According to the description of our project, A noisy signal can be modeled as

$$Y_m(n) = S_m(n) + W_m(n) \quad (1)$$

where $S_m(n)$ is the target speech signal at sample-time index n and microphone $m \in \{1, 2, \dots, M\}$ and $W_m(n)$ is the noise signal. After the framing and windowing, we do frequency transformation on each signal, then the signal model in the frequency domain is given by

$$Y_m(l, k) = A_m(l, k)S(l, k) + W_m(l, k) \quad (2)$$

where $S(l, k)$ is now the target source at the source location, and k indicates the frequency-bin index. we stack the DFT coefficients for the different microphones in a vector and assume that $A_m(l, k) = 1 \forall k, l, m$, we obtain

$$\mathbf{Y}(l, k) = \mathbf{S}(l, k) + \mathbf{W}(l, k) \quad (3)$$

The goal is to estimate the target speech signal in Python by inputting the .csv file converted from the given Matlab data. In the preprocessing of the data, we first divide the speech signal into frames using a 50% overlap, 20 ms duration and apply a Hann window to each frame, then perform an FFT on each frame. After FFT the spectrum of clean speech signal is shown in Figure 1, and the spectrum of channel 1 is shown in Figure 2. We obtain the spectrum of signal in a dB scale to make noise signals more prominent and visible.

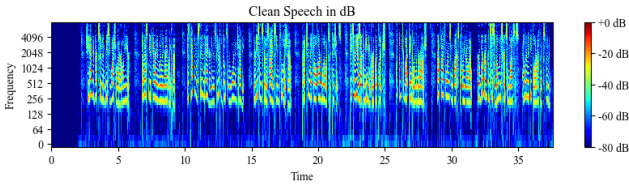


Figure 1: Clean speech spectrum

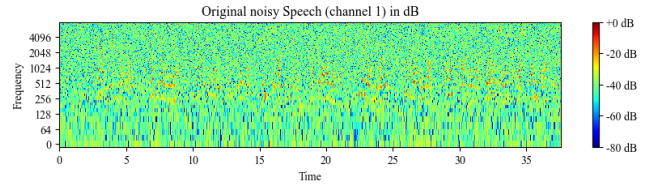


Figure 2: Noisy speech spectrum (channel 1)

2.2 Assessment metrics

In this project, we use three assessment metrics to evaluate the performance of our estimator.

2.2.1 Signal-to-Noise Ratio: SNR

SNR is a straightforward way to quantify the amount of noise reduction achieved by a speech enhancement model. Also, there is often a correlation between SNR and the perceived quality of speech. Higher SNR generally suggests clearer speech with less intrusive noise, which typically aligns with a better subjective listening experience. Equation 4 calculates SNR (γ_k) using $\sigma_{Y_k}^2$ and $\sigma_{W_k}^2$ for noisy speech and noise variances, respectively. As variance in a zero-mean signal signifies energy, γ_k represents the ratio of speech energy to noise energy. The noise energy is determined from the first-second data in the noisy speech, where no speech signal is present in the clean speech.

$$\gamma_k = 10\log_{10} \left(\frac{\sigma_{Y_k}^2}{\sigma_{W_k}^2} \right) = 10\log_{10} \left(\frac{E_Y}{E_W} \right) \quad (4)$$

2.2.2 Cramér-Rao Lower Bound: CRLB

An estimator is considered unbiased if its expected value is equal to the true value of the parameter it is estimating. The Cramér-Rao Lower Bound provides a theoretical lower bound on the variance of an unbiased estimator. In speech enhancement, it can be used to determine the theoretical limit of the performance of an estimator in terms of its variance.

$$\text{CRLB}(\theta) = \frac{1}{I(\theta)} \quad (5)$$

$$I(\theta) = -E \left[\frac{\partial^2 \ln f(X; \theta)}{\partial \theta^2} \right] \quad (6)$$

$I(\theta)$ is the Fisher Information for the parameter θ . $f(X; \theta)$ is the probability density function of the observed data X parameterized by θ . In this project, it's very hard to calculate the fisher information, so we use a general expression of CRLB. When it is a linear Gaussian model, like $x = h\theta + w$, $w \sim \mathcal{N}(0, C_w)$, where θ is an unknown parameter, x is the observed signal, and w is Gaussian noise with a mean of 0 and a covariance matrix C_w . The CRLB can be simplified as $\frac{1}{h^T C_w^{-1} h}$. Additionally, since it is a zero-mean signal, we employ the Welch method to compute the average power, which in this case is equivalent to the variance.

2.2.3 Variance of the estimator

To determine the variance of the estimator, we use Equation 7 to calculate, where K and L denote the frequency-bin index and time frame, respectively. A lower variance signifies higher precision in the estimation process. Therefore, the closer the variance approaches the CRLB, the better the performance of the estimator.

$$\text{var}_{\text{emp}} = \frac{1}{KL} \sum_{k=1}^K \sum_{l=1}^L |\hat{S}(l, k) - S(l, k)|^2 \quad (7)$$

3 Deterministic speech signal

In this chapter, we assume that clean speech in noisy speech is the deterministic signal. The signal type of the clean speech determines how to estimate the clean method from the noisy speech.

3.1 Model selection

The common estimators for deterministic signals include Maximum Likelihood Estimation (MLE), Best Linear Unbiased Estimator (BLUE), and Least Squares (LS). Notably, MLE stands out for its robustness and versatility, requiring fewer distribution assumptions compared to BLUE and being less susceptible to outliers than LS. These features make MLE the preferred option for precise and reliable parameter estimation, particularly in diverse or intricate datasets like those found in speech signal data. In conclusion, we choose the MLE as the estimator to denoise the signal.

3.2 Implementation of MLE

The definition of MLE is to find the value \mathbf{S} that maximizes $p(\mathbf{Y}; \mathbf{S})$ for a fixed \mathbf{Y} . In the MLE method proposed in [1], the amplitude and phase spectra of the clean speech signal at frequency bin k , denoted as S_k and $\theta_s(k)$, are assumed to be uncertain but deterministic. Because the noise signal $W_m(l, k) \sim \mathcal{CN}(0, \sigma_{W_m(l, k)}^2)$, the distribution of $\mathbf{Y}(\omega_k)$ is also Gaussian distribution with variance $\sigma_{W_m(l, k)}^2$ and mean $S_k e^{j\theta_s(k)}$. In the upcoming equations in this report, for the sake of simplicity in notation, we will focus on the estimation process in one time frame and one channel. Thus, the representation of the variance in noise can be denoted as $\sigma_{W_k}^2$, signifying the variance of noise at frequency bin k . So, the probability density function of $\mathbf{Y}(\omega_k)$ can be formulated in Equation 8 [1].

$$p(Y(\omega_k); S_k, \theta_s(k)) = \frac{1}{\pi \sigma_{W_k}^2} \exp \left[-\frac{|Y(\omega_k) - S_k e^{j\theta_s(k)}|^2}{\sigma_{W_k}^2} \right] \quad (8)$$

To achieve the maximum estimation of \mathbf{S} , we need to simplify and differentiate the log-likelihood function $\log(p(Y(\omega_k); S_k, \theta_s(k)))$ with respect to S_k and set the derivative to zero. Then, we can get the estimation function of S_k in Equation 9 [2]. In Equation 9, γ_k denotes the posteriori SNR which can be calculated by Equation 4. Because Equation 9 is solely designed to estimate the magnitude of \hat{S} , it is necessary to add the phase information of $\mathbf{Y}(\omega_k)$ after the magnitude estimation.

$$\hat{S}(\omega_k) = \left[\frac{1}{2} + \frac{1}{2} \sqrt{\frac{\gamma_k - 1}{\gamma_k}} \right] Y(\omega_k) \quad (9)$$

In equations Equation 9 and 4, the critical aspect lies in estimating the magnitude of \hat{S} to compute the noise energy E_W . To address the variability in noise variance throughout the speech signal, voice activity detection (VAD) is employed to improve the estimation of the noise energy. VAD captures short-time energy from the input signal, comparing it to a threshold value [2], typically determined during non-speech periods. The step of VAD is outlined in algorithm 1. In algorithm 1, the parameter M is configured for the frames of the first second, considering the absence of speech during this initial period. The smoothing parameter μ is set to 0.9, falling within the recommended range of 0.3 to 0.96 [2]. Additionally, the SNR threshold of 5 dB is established to ensure accurate identification of speech frames [3].

To apply MLE across multiple channels, we choose to combine data from these channels by calculating the average value. Because the noise signal W_m has a zero-mean property, averaging across more channels makes the mean of Y approach zero. This can weaken the impact of noise on the speech signal to a certain extent, to achieve a certain degree of noise reduction purposes and convert the multi-channel data into one channel. Besides, the Python code to implement MLE is shown in Appendix B.

Algorithm 1: Voice Activity Detection

Data: Number of speech frames N , the first M frames

- 1 **Set parameter value** SNR threshold: θ , smoothing parameter: μ ;
 - 2 **Calculate noise energy** Calculate E_W as the average energy over the first M frames ($E_W = \sum_k \sigma_{W_k}^2$);
 - 3 **for** $t = 1$ **to** N **do**
 - 4 **Calculate energy of each frame** E_Y ;
 - 5 **Calculate SNR** $SNR = \frac{E_Y}{E_W}$;
 - 6 **if** $SNR < \theta$ **then**
 - 7 $E_W = \mu E_W + (1 - \mu) E_Y$
-

3.3 Results of MLE

In this section, we will initially present the results of MLE using only a single-channel data input. Then, we will showcase the MLE results in the context of a multi-channel case.

3.3.1 Case 1: MLE result in a single-channel

To test MLE in a single channel, we select the first channel from the noisy speech data. Following the MLE algorithm and VAD, we obtain the spectrum of enhanced speech in a dB scale, as illustrated in Figure 3. In ??, it is clear that the speech signal is clearer compared to the original noisy speech (shown in Figure 2). Then, we reverse the FFT process to obtain the signal in the time domain to calculate the SNR based on Chapter 2.2.1. The enhanced speech SNR is found to be 3.325. This represents a 9.88% improvement compared to the original noisy data in the first channel (SNR=3.026).

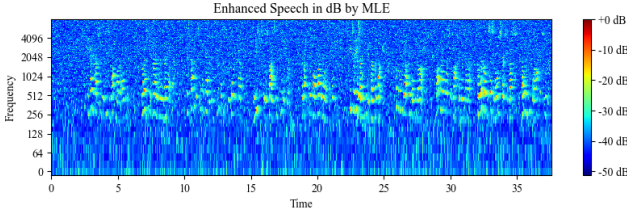


Figure 3: Enhanced speech spectrum of MLE in single-channel

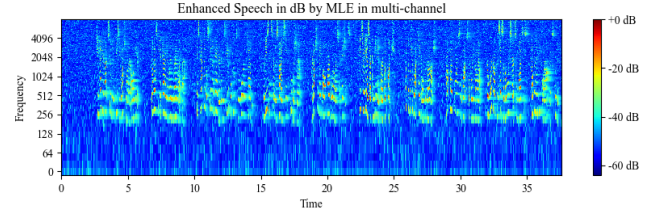


Figure 4: Enhanced speech spectrum of MLE in 16-channels

3.3.2 Case 2: MLE result in multi-channel

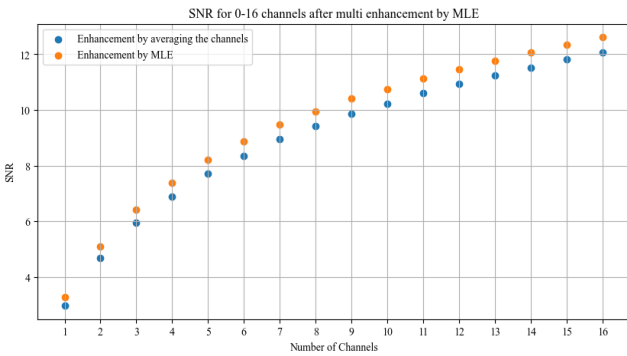


Figure 5: SNR for 1-16 channels after multi-channel enhancement by MLE

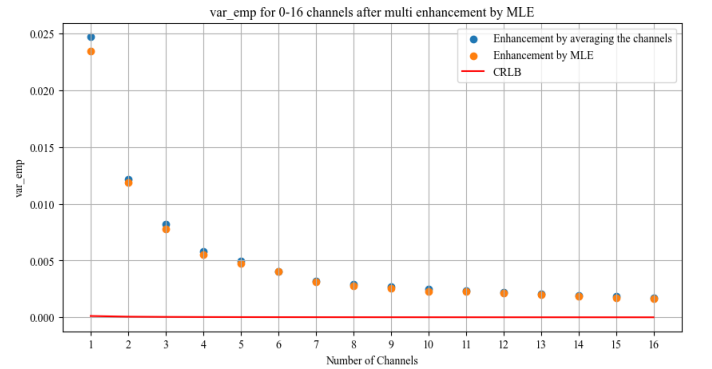


Figure 6: var_{emp} for 1-16 channels after multi-channel enhancement by MLE

After successfully applying MLE to denoise single-channel noisy speech, the next step involves testing MLE in a multi-channel environment. As outlined in Chapter 3.2, the multi-channel data is converted into a single-channel format by averaging. Following this conversion, the single-channel MLE process is employed to obtain the estimated clean speech. In Figure 4, the spectrum of the enhanced speech, utilizing data from all 16 channels, exhibits a clearer speech signal compared to both the original signal and the enhanced speech from a single channel. The effectiveness of MLE is evident in improving the quality of the speech signal. In terms of SNR, as shown in Figure 5, the orange dots representing MLE's SNR across 0 to 16 channels consistently outperform the blue dots, which represent SNR obtained solely by averaging multi-channel data. This highlights MLE's superiority in enhancing SNR

compared to the simple averaging method. Notably, the SNR of MLE increases with the growing number of channels, reaching its highest value of 12.61 when utilizing all 16 channels.

The variance of the estimator across 1 to 16 channels is depicted in Figure 6 based on the calculation of Equation 7. The blue dots illustrate the enhancement achieved by averaging the data, while the orange dots represent the MLE in multi-channel estimation. The red line signifies the CRLB, representing the minimum variance from the estimator. In Figure 6, it's clear that the variance of the MLE estimator decreases and tends to approach the CRLB as the number of microphones increases. Nonetheless, there remains a discernible deviation from the CRLB. Moreover, the variance of the MLE is slightly lower than that achieved by simply averaging the data. This observation suggests that the MLE has the potential to bring about a slight improvement in the variance compared to the simple averaging of the data.

4 Stochastic speech signal

In this section, we assume that the speech signal is not deterministic but rather a statistical signal. Therefore, we cannot continue to use MLE to reduce noise which is suitable to the deterministic signal. To estimate the random speech signal, we need an estimator that considers the statistical properties of both the signal and the noise.

4.1 Model selection

The common estimator of the random signal is the Maximum A Posteriori Probability (MAP) estimator and Minimum Mean Square Error (MMSE) estimator. The MAP estimator, emphasizes the use of prior information. However, MMSE often outperforms MAP in practice. This is because MMSE is more flexible in modeling noise and demonstrates better adaptability to the variations in actual noise, thereby offering more stable and accurate results in the process of speech signal restoration. As the clean speech is unknown in the speech noise reduction progress. It is hard to use the prior information to estimate the clean signal by MAP estimator. In conclusion, we decide to use the MMSE to denoise the noisy speech.

4.2 Implementation of MMSE

The goal of the MMSE algorithm is to minimize the expected value of the Bayesian mean square error (Bmse). After setting the derivative of Bmse to zero we obtain, we can get the estimated value $\hat{S}_k = E[S_k|Y_k]$. This equation is to calculate the mean of the a posteriori probability density function of S which is not easy because the speech signal is neither a stationary nor an ergodic process [2]. After computation, the gain function for estimating the magnitude of S_k in Equation 10 is derived [2]. This process is based on two assumptions: 1. Fourier transform coefficients exhibit a Gaussian probability distribution with a zero mean. 2. The Fourier transform coefficients are statistically independent.

$$\hat{S}_k = \frac{\sqrt{\pi}}{2} \frac{\sqrt{\nu_k}}{\gamma_k} \exp\left(-\frac{\nu_k}{2}\right) \left[(1 + \nu_k) I_0\left(\frac{\nu_k}{2}\right) + \nu_k I_1\left(\frac{\nu_k}{2}\right) \right] Y_k \quad (10)$$

In Equation 10, I_0 and I_1 denote the modified Bessel functions of zero and first order, respectively. The variable ν_k is composed of the priori SNR ξ_k and the posteriori SNR γ_k , as defined in Equation 11. The equation for computing ξ_k is presented in Equation 12 ($\lambda_S(k)$ denotes the variance of clean speech signal), while the equation for calculating γ_k is given in Equation 4. In Equation 10, ξ_k plays a primary role in noise suppression, and the MMSE method is particularly sensitive to the fluctuations of ξ_k [2]. Therefore, an accurate and relatively smooth estimation of ξ_k is essential.

$$\nu_k = \frac{\xi_k}{1 + \xi_k} \gamma_k \quad (11)$$

$$\xi_k = \frac{\lambda_S(k)}{\sigma_{W_k}^2} \quad (12)$$

To estimate ξ_k , we can use its relationship with γ_k called decision-directed method [4]. The equation to calculate ξ_k by the decision-directed method is shown in Equation 13. In this equation, the priori SNR ξ_k in time frame l is updated using the value from the previous time frame $(l - 1)$ and the posteriori SNR γ_k in current time frame l . The parameter α is called the smoothing constant and is typically set to 0.98, as recommended in [2]. The value ξ_{\min} is important for reducing low-level musical noise and is advised to be set to $-15dB$ according to [5]. When determining the posteriori SNR γ_k , the VAD method can be employed for estimation, utilizing Equation 4 and following algorithm 1.

$$\hat{\xi}_k(l) = \max \left[\alpha \frac{\hat{S}_k^2(l-1)}{\sigma_{W_k}^2(l-1)} + (1 - \alpha) \max[\gamma_k(l) - 1, 0], \xi_{\min} \right] \quad (13)$$

Similar to the approach employed in the MLE, we can calculate the average value across all channels. This process helps transform the multi-channel signal into a single-channel signal, and the detailed procedure is outlined in Chapter 3.2. Besides, the Python code to implement MLE is shown in Appendix B.

4.3 Results of MMSE

4.3.1 Case 1: MMSE result in a single-channel

Like the test for the first channel in Chapter 3.3.1, after we estimate the first channel noisy speech signal based on MMSE, VAD and the decision-directed method, we can get the spectrum of the enhanced speech. The spectrum is shown in Figure 7. We can find that the speech signal can be recovered from the noisy data successfully. After reversing the FFT process, we calculated the

SNR following the steps outlined in Chapter 3.3.1. The SNR for MMSE in the single-channel is 9.57, which is almost three times higher than the SNR for MLE in the single-channel test (SNR = 3.325). Comparing the results of MLE and MMSE, we find that MMSE performs better in reducing noise in the same test environment.

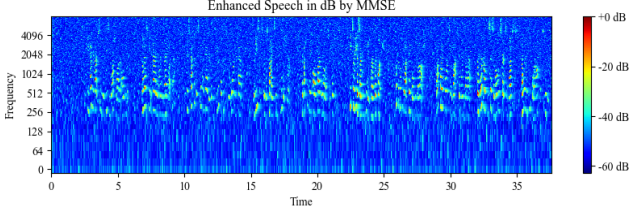


Figure 7: Enhanced speech spectrum of MMSE in single-channel

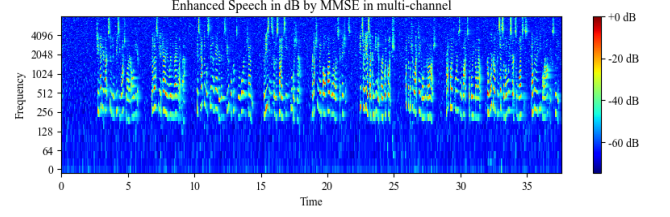


Figure 8: Enhanced speech spectrum of MMSE in 16-channels

4.3.2 Case 2: MMSE result in multi-channel

In the multi-channel test, similar to the single-channel MMSE test and the MLE multi-channel test, we averaged signals across varying channel numbers (1 to 16) and applied the MMSE algorithm for clean speech estimation. The resulting spectrum in Figure 8 demonstrates that the speech signal is noticeably clearer compared to the 16-channel MLE spectrum (Figure 4). In terms of SNR, as shown in Figure 9, the 16-channels MMSE significantly outperforms both 16-channels MLE and the average of 16 channels. Notably, MMSE achieves this high performance with just two microphones, rivaling the effectiveness of MLE with all 16 channels. Moreover, SNR increases with additional channels, peaking at 16 channels with a value of 23.008, nearly double that of MLE with 16 channels. The variance analysis in Figure 10 reveals that MMSE consistently maintains significantly lower variance compared to MLE and the average data. With the increasing number of channels, MMSE's variance decreases and nearly approaches the CRLB. Besides, the variance in one-channel MMSE even surpasses that of 16-channels MLE. The graph in Appendix A further illustrates that the difference between 16-channels MMSE and CRLB is numerically less than 0.0002.

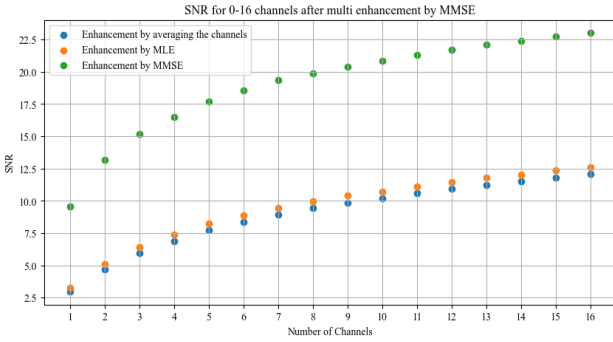


Figure 9: SNR for 1-16 channels after multi-channel enhancement by MMSE

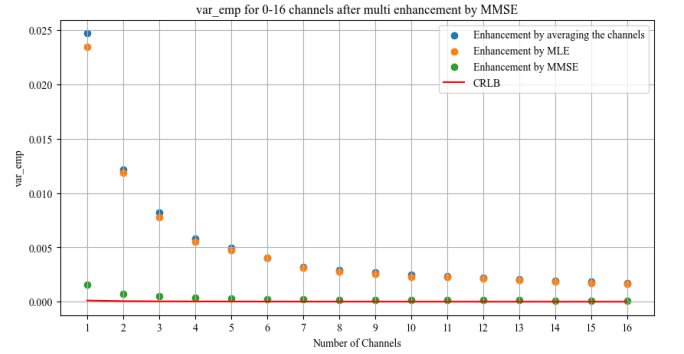


Figure 10: var_{emp} for 1-16 channels after multi-channel enhancement by MMSE

5 Conclusion

In this report, we first introduce the background and theoretically analyze the speech system. Next in the following two scenarios—Deterministic Speech Signal and Stochastic Speech Signal—we conduct a comparative analysis of various models, elucidating our rationale for selecting the MLE and the MMSE methods. When the speech signal is deterministic, we choose MLE as our estimator because MLE is more robust and flexible than BLUE and LS, needing fewer assumptions and being less affected by outliers. When the speech signal is Stochastic, we choose MMSE because of its flexibility in noise modeling and better adaptability to noise variations. What's more, the MAP estimator's reliance on prior information is less effective in the context of unknown clean speech during noise reduction.

In the solution of this problem, We use the variance of the estimator to assess our model and test if it reaches the Cramer-Rao bound. What's more, in order to measure the clarity of a signal in the presence of noise, We use SNR to assess our model. After preprocessing the data, we use VAD to extract the power of the noise and apply our estimator to the speech signal. Then We detail the implementation process of speech enhancement.

As a result, we can conclude that simply combining the data from one to sixteen channels has an improvement in SNR and variance. In the scenario of a deterministic signal, our MLE estimator performs better than just combining the data but doesn't reach the CRLB. However, when we use MMSE in the scenario of stochastic signal, the estimator performs much better than MLE, and it is very close to the CRLB. We assume the reason is that speech signals are generally considered as random signals in reality. MMSE's superiority in speech enhancement applications primarily stems from its suitability for handling stochastic signals like speech, as well as its ability to adapt to variations in speech signals in noisy environments. MLE, on the other hand, might be more effective in other applications where the signals are more deterministic and noise is less of a factor.

References

- [1] McAulay, R. J. and Malpass, M. L. (1980), Speech enhancement using a soft-decision noise suppression filter, *IEEE Trans. Acoust. Speech Signal Process.*, 28, 137–145
- [2] Philipos C. Loizou. 2013. *Speech Enhancement: Theory and Practice* (2nd. ed.). CRC Press, Inc., USA.
- [3] T. Bäckström, “Voice activity detection (VAD),” 8.1. Voice Activity Detection (VAD) - Introduction to Speech Processing, https://speechprocessingbook.aalto.fi/Recognition/Voice_activity_detection.html (accessed Jan. 4, 2024).
- [4] Ephraim, Y. and Malah, D. (1984), Speech enhancement using a minimum mean square error short-time spectral amplitude estimator, *IEEE Trans. Acoust. Speech Signal Process.*, 32(6), 1109–1121
- [5] Cappe, O. (1994), Elimination of the musical noise phenomenon with the Ephraim and Malah noise suppressor, *IEEE Trans. Speech Audio Process.*, 2(2), 346–349.

Appendix A Zoomed-in graph of the variance of 16-channel MMSE and CRLB

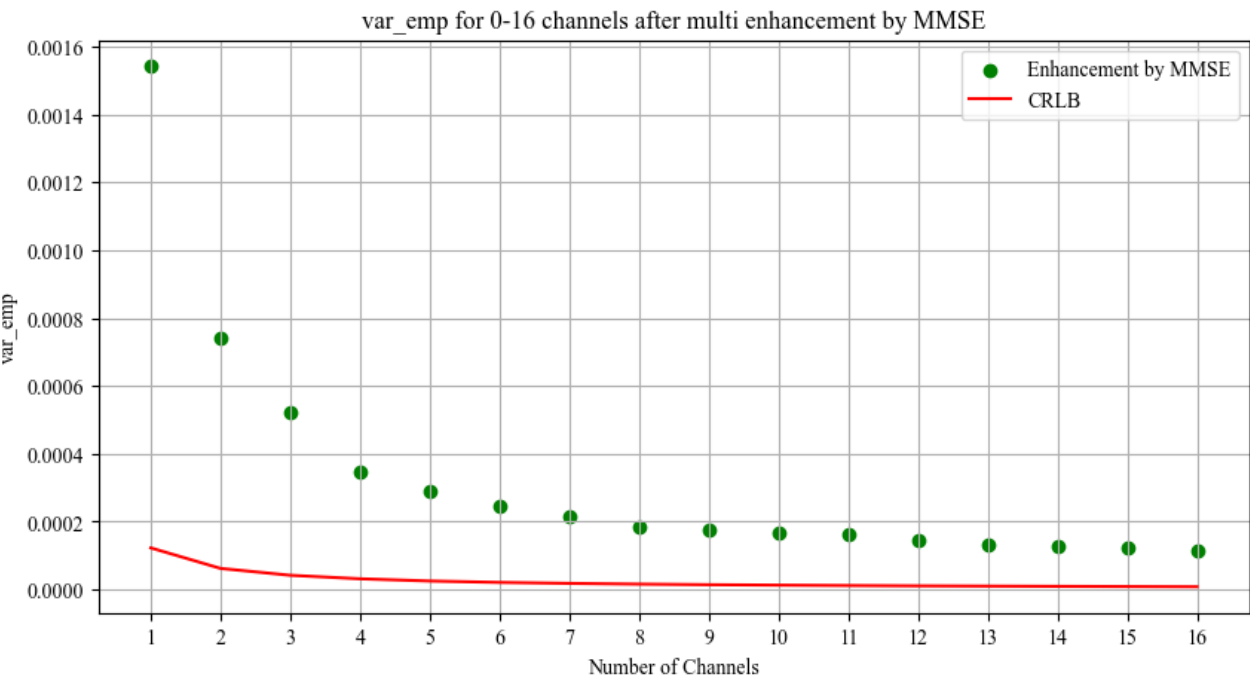


Figure 11: Zoomed-in graph of variance of 16-channel MMSE and CRLB

Appendix B Program code in python

```
1 #####
2
3 # # Multi-Microphone Speech Enhancement
4
5 #####
6 # ***ET4386 Estimation and Detection Project***
7 #
8 # *Author: Yanqi Hong | Kaishen Lin*
9 #####
10
11 #import library
12 import matplotlib.pyplot as plt
13 import numpy as np
14 import pandas as pd
15 import numpy as np
16 from scipy.signal import hamming
17 from scipy.signal import welch
18 from scipy.fft import fft
19 from scipy.fftpack import fft
20 import librosa
21 import scipy.special as sp
22 import soundfile as sf
23 import statistics
24
25 # Set the default font for plots
26 plt.figure(dpi=200)
27 plt.rcParams['font.family'] = 'Times New Roman'
28
29
30
31 #####
32 # ## 1. Data analysis
33 #####
34
35 #-----#
36 # ### 1.1 Data input and general analysis
37 #-----#
38
39 #####
40 #####Code to conver .mat file to .csv file in matlab#####
41 # writematrix(Clean,'Clean.csv');
42 # writematrix(Data,'Data.csv');
43 #####
44
45 #import csv file
46 Data = pd.read_csv('Data.csv', header=None)      #noisy speech with 16 channels
47 Clean=pd.read_csv('Clean.csv', header=None)      #clean speech
48 display(Data)
49 display(Clean)
50 # plot all 16 channels in 16 subplots
51 plt.figure(figsize=(20,10))
52 for i in range(1,17):
53     t=np.linspace(0, len(Data[i-1])/fs, len(Data[i-1]))
54     plt.subplot(4,4,i)
55     plt.plot(Data[i-1], linewidth=0.5) # Set the linewidth to 0.5
56     plt.title('Channel '+str(i))
57     plt.xlabel('Time (s)')
58     plt.ylabel('Amplitude')
59
60 plt.subplots_adjust(hspace=0.5, wspace=0.5)
61 plt.show()
62
63 # plot clean speech
64 t=np.linspace(0, len(Clean)/fs, len(Clean))
65 plt.figure(figsize=(10,5))
66 plt.plot(t,Clean, linewidth=0.5) # Set the linewidth to 0.5
67 plt.title('Clean Speech')
68 plt.xlabel('Time (s)')
69 plt.ylabel('Amplitude')
70 plt.show()
71
72 Clean = np.transpose(Clean.to_numpy())[0]
73 Data = Data.to_numpy()
74
75
76 #-----#
77 # ### 1.2 SNR evaluation
```

```

78 #-----#
79
80 # Evaluation by SNR
81 def SNR(NoisyData):
82     """
83     Calculate the Signal-to-Noise Ratio (SNR) of the given noisy data.
84
85     Parameters:
86     NoisyData (array-like): The input noisy data.
87
88     Returns:
89     float: The calculated SNR value.
90     """
91     noise_segment = NoisyData[:15000]
92     snr = 10*np.log10(np.var(NoisyData)/np.var((noise_segment)))
93     return snr
94
95
96
97
98
99 # Calculate SNR for each channel
100 SNR_each_channel = np.zeros(16)
101 for i in range(16):
102     SNR_each_channel[i] = SNR(Data[:,i])
103
104 # Plot SNR for each channel as points
105 plt.figure(figsize=(10,5))
106 plt.scatter(range(1, 17), SNR_each_channel)
107 plt.title('original SNR for each channel')
108 plt.xticks(range(1, 17))
109 plt.xlabel('Channel')
110 plt.ylabel('SNR')
111 plt.grid()
112 plt.show()
113
114 print('average SNR : ', statistics.mean(SNR_each_channel))
115
116
117 #-----#
118 # ### 1.3 Spectrogram analysis
119 #-----#
120
121 #specrum in clean speech
122 Fs = 16000
123 # speech_noisy_freq = np.zeros((320, 3760, 1))
124
125
126 sample_rate = 16000
127 frame_size = 2*int(sample_rate * 0.02)
128 hop_size = frame_size // 4
129
130 clean_freq = librosa.stft(Clean, n_fft=frame_size, hop_length=hop_size, window='hann')
131
132
133 # subplots of the clean speech in linear and dB scale spectrogram
134 plt.figure(figsize=(10,5))
135 plt.subplot(2,1,1)
136 librosa.display.specshow(abs(clean_freq), sr=sample_rate, hop_length=hop_size, x_axis='time', y_axis='log', cmap='jet')
137 plt.title('Clean Speech in linear scale')
138 plt.xlabel('Time')
139 plt.ylabel('Frequency')
140 plt.colorbar()
141 plt.subplots_adjust(hspace=0.5, wspace=0.5)
142
143
144 plt.subplot(2,1,2)
145 librosa.display.specshow(librosa.amplitude_to_db(abs(clean_freq), ref=np.max), sr=sample_rate, hop_length=hop_size, x_axis='time', y_axis='log', cmap='jet')
146 plt.title('Clean Speech in dB')
147 plt.xlabel('Time')
148 plt.ylabel('Frequency')
149 plt.colorbar(format='%+2.0f dB')
150 plt.show()
151
152
153 #specrum in noisy speech
154 Fs = 16000
155 speech_noisy_freq = np.zeros((321, 3761, 16))
156

```

```

157 #calculate the fft for each channel
158 for i in range(16):
159     speech_noisy_freq[:, :, i] = librosa.stft(Data[:, i], n_fft=frame_size, hop_length=hop_size, window='hann')
160
161
162 # # subplots of the noisy speech (channel 1) in linear and dB scale spectrogram
163 plt.figure(figsize=(10, 5))
164 plt.subplot(2, 1, 2)
165 librosa.display.specshow(librosa.amplitude_to_db(abs(speech_noisy_freq[:, :, 0]), ref=np.max), sr=sample_rate, hop_length=
hop_size, x_axis='time', y_axis='log', cmap='jet')
166 plt.title('Original noisy Speech (channel 1) in dB')
167 plt.xlabel('Time')
168 plt.ylabel('Frequency')
169 plt.colorbar(format='%+2.0f dB')
170
171 plt.subplot(2, 1, 1)
172 librosa.display.specshow(abs(speech_noisy_freq[:, :, 0]), sr=sample_rate, hop_length=hop_size, x_axis='time', y_axis='log',
cmap='jet')
173 plt.title('Original noisy Speech (channel 1) in linear scale')
174 plt.xlabel('Time')
175 plt.ylabel('Frequency')
176 plt.colorbar()
177 plt.subplots_adjust(hspace=0.5, wspace=0.5)
178 plt.show()
179
180
181 # subplots for all 16 channels in dB scale spectrogram
182 plt.figure(figsize=(45, 25))
183 for i in range(16):
184     plt.subplot(4, 4, i+1)
185     librosa.display.specshow(librosa.amplitude_to_db(abs(speech_noisy_freq[:, :, i]), ref=np.max), sr=sample_rate,
hop_length=hop_size, x_axis='time', y_axis='log', cmap='jet')
186     plt.title('Channel ' + str(i+1))
187     plt.xlabel('Time')
188     plt.ylabel('Frequency')
189     plt.colorbar(format='%+2.0f dB')
190
191 plt.subplots_adjust(hspace=0.5, wspace=0.25) # Adjust hspace to reduce the amount of whitespace below the main title
192 plt.suptitle('Noisy Speech in dB for all 16 channels')
193 plt.show()
194
195 plt.figure(figsize=(45, 25))
196 for i in range(16):
197     plt.subplot(4, 4, i+1)
198     librosa.display.specshow(abs(abs(speech_noisy_freq[:, :, i])), sr=sample_rate, hop_length=hop_size, x_axis='time',
y_axis='log', cmap='jet')
199     plt.title('Channel ' + str(i+1))
200     plt.xlabel('Time')
201     plt.ylabel('Frequency')
202     plt.colorbar()
203
204
205 plt.subplots_adjust(hspace=0.5, wspace=0.25) # Adjust hspace to reduce the amount of whitespace below the main title
206 plt.suptitle('Noisy Speech in linear scale for all 16 channels')
207 plt.show()
208
209
210
211 #-----#
212 # ### 1.4 CRLB analysis
213 #-----#
214
215 #Enhancement for all channels only by calculating the mean of all channels
216 snr_all_mean = np.zeros(16)
217 speech_enhanced_all_mean = np.zeros((601600, 16))
218 speech_enhanced_freq_all_mean = np.zeros((321, 3761, 16))
219
220 for i in range(16):
221     speech_enhanced_freq_all_mean[:, :, i] = np.mean(speech_noisy_freq[:, :, 0:i+1], axis=2)
222     speech_enhanced_all_mean[:, i] = librosa.istft(speech_enhanced_freq_all_mean[:, :, i], n_fft=frame_size, hop_length=
hop_size, window='hann')
223     snr_all_mean[i] = SNR(speech_enhanced_all_mean[:, i])
224
225
226 # #CRLB
227 CRLB = np.zeros(16)
228 fs = 16000
229 for i in range(16):
230     channel_speech = np.mean(Data[:, 0:i+1], axis=1)
231     noise_segment = channel_speech[:fs]
232     frequencies, power_spectral_density = welch(noise_segment, fs, nperseg=320) #Calculate power spectral density

```

```

233     CRLB[i] = np.sum(power_spectral_density) / len(power_spectral_density) #Calculate average power
234
235
236 # Plot CRLB for each channel as points
237 plt.figure(figsize=(10,5))
238 plt.scatter(range(1, 17), CRLB)
239 plt.title('CRLB for channel 0-16')
240 plt.xticks(range(1, 17))
241 plt.xlabel('Channel')
242 plt.ylabel('CRLB')
243 plt.grid()
244 plt.show()
245
246
247
248 #####
249 # ## 2. Denoise
250 #####
251
252 #-----#
253 # ### 2.1. Denoise with MLE
254 #-----#
255
256
257 #.....#
258 # #### 2.1.1 MLE for single channel
259 #.....#
260
261 # #### Estimation
262 def MLE_speech_enhancement(speech_noisy_freq):
263     """
264     Perform maximum likelihood estimation (MLE) based speech enhancement.
265
266     Parameters:
267     speech_noisy_freq (ndarray): Noisy speech frequency domain representation with shape (D, T), where D is the number of
268     frequency bins, T is the number of frames.
269
270     Returns:
271     speech_enhanced (ndarray): Enhanced speech in the time domain.
272     speech_enhanced_freq (ndarray): Enhanced speech frequency domain representation with shape (D, T), where D is the
273     number of frequency bins and T is the number of frames.
274
275     """
276
277     # parameters
278     max_gamma = 40 # maximum value of gamma
279     mu_VAD = 0.9 # parameter for VAD noise tracking
280     th_VAD = 15 # VAD decision threshold
281
282     phase_noisy = np.angle(speech_noisy_freq)
283     mag_noisy = np.abs(speech_noisy_freq)
284
285     D, T = np.shape(mag_noisy)
286
287     time_index_num = int(1//0.02) # initial noise length: 1s
288     noise_segment = mag_noisy[:, 0:time_index_num]
289     mag_noise = np.mean(np.abs(noise_segment), axis=1)
290     power_noise = mag_noise**2
291
292     mag_enhance = np.zeros([D, T])
293
294     for i in range(T):
295
296         # get energy spectrum and magnitude spectrum for each frame
297         mag_frame = mag_noisy[:, i]
298         power_frame = mag_frame**2
299
300         # get Signal-to-Noise Ratio (SNR) used for VAD calculation
301         SNR_VAD = 10 * np.log10(np.sum(power_frame) / np.sum(power_noise))
302
303         # calculate posterior SNR
304         gamma = np.minimum(power_frame / power_noise, max_gamma)
305
306         # update power_noise based on VAD
307         mu = mu_VAD
308         if SNR_VAD < th_VAD:
309             power_noise = mu * power_noise + (1 - mu) * power_frame
310
311

```

```

312     G = 0.5 + 0.5 * np.sqrt((abs(gamma - 1)) / gamma)
313
314
315     mag_enhance_frame = G * mag_frame
316     mag_enhance[:, i] = mag_enhance_frame
317
318     speech_enhanced_freq = mag_enhance * np.exp(1j * phase_noisy)
319
320     speech_enhanced = librosa.istft(speech_enhanced_freq, n_fft=frame_size, hop_length=hop_size, window='hann')
321
322     return speech_enhanced, speech_enhanced_freq
323
324
325 # Denoise the first noisy speech channel by the MLE speech enhancement
326 enhanced_speech, enhanced_speech_freq = MLE_speech_enhancement(speech_noisy_freq[:, :, 0])
327
328
329 # #### Evaluation
330 print("SNR of the enhanced speech is: ", SNR(enhanced_speech))
331
332 #plot the enhanced speech
333 #subplots
334 plt.figure(figsize=(10,5))
335 plt.subplot(2,1,1)
336 librosa.display.specshow(librosa.amplitude_to_db(abs(speech_noisy_freq[:, :, 0]), ref=np.max), sr=sample_rate, hop_length=
337 hop_size, x_axis='time', y_axis='log', cmap='jet')
338 plt.title('Original noisy Speech in dB')
339 plt.xlabel('Time')
340 plt.ylabel('Frequency')
341 plt.colorbar(format='%+2.0f dB')
342
343 plt.subplot(2,1,2)
344 librosa.display.specshow(librosa.amplitude_to_db(abs(enhanced_speech_freq), ref=np.max), sr=sample_rate, hop_length=
345 hop_size, x_axis='time', y_axis='log', cmap='jet')
346 plt.title('Enhanced Speech in dB by MLE')
347 plt.xlabel('Time')
348 plt.ylabel('Frequency')
349 plt.colorbar(format='%+2.0f dB')
350 plt.subplots_adjust(hspace=0.5, wspace=0.5)
351 plt.show()
352
353 sf.write("enhce_mmse.wav", enhanced_speech, fs)
354
355
356 #.....#
357 # #### 2.1.2. MLE for multi channel
358 #.....#
359
360 # #### Estimation
361 def MLE_multi_channel_speech_enhancement(speech_noisy_freq):
362     """
363     Applies Maximum Likelihood Estimation (MLE) based multi-channel speech enhancement to the given noisy speech frequency
364     domain representation.
365
366     Parameters:
367     speech_noisy_freq (ndarray): Noisy speech frequency domain representation with shape (D, T, L), where D is the number
368     of frequency bins, T is the number of frames, and L is the number of channels.
369
370     Returns:
371     speech_enhanced (ndarray): Enhanced speech in the time domain.
372     speech_enhanced_freq (ndarray): Enhanced speech frequency domain representation with shape (D, T), where D is the
373     number of frequency bins and T is the number of frames.
374     """
375
376     speech_noisy_freq = np.mean(speech_noisy_freq, axis=2) # average over all channels
377
378     # parameters
379     max_gamma = 40 # maximum value of gamma
380     mu_VAD = 0.9 # parameter for VAD noise tracking
381     th_VAD = 5 # VAD decision threshold
382
383     phase_noisy = np.angle(speech_noisy_freq)
384     mag_noisy = np.abs(speech_noisy_freq)
385
386     D, T = np.shape(mag_noisy)
387
388     time_index_num = int(1/0.02) # initial noise length: 1s

```

```

388 noise_segment = mag_noisy[:, 0:time_index_num]
389 mag_noise = np.mean(np.abs(noise_segment), axis=1)
390 power_noise = mag_noise**2
391
392 mag_enhance = np.zeros([D, T])
393
394 for i in range(T):
395
396     # get energy spectrum and magnitude spectrum for each frame
397     mag_frame = mag_noisy[:, i]
398     power_frame = mag_frame**2
399
400     # get Signal-to-Noise Ratio (SNR) used for VAD calculation
401     SNR_VAD = 10 * np.log10(np.sum(power_frame) / np.sum(power_noise))
402
403     # calculate posterior SNR
404     gamma = np.minimum(power_frame / power_noise, max_gamma)
405
406     # update power_noise based on VAD
407     mu = mu_VAD
408     if SNR_VAD < th_VAD:
409         power_noise = mu * power_noise + (1 - mu) * power_frame
410
411
412     G = 0.5 + 0.5 * np.sqrt((abs(gamma - 1)) / gamma)
413     mag_enhance_frame = G * mag_frame
414     mag_enhance[:, i] = mag_enhance_frame
415
416     speech_enhanced_freq = mag_enhance * np.exp(1j * phase_noisy)
417
418     speech_enhanced = librosa.istft(speech_enhanced_freq, n_fft=frame_size, hop_length=hop_size, window='hann')
419
420     return speech_enhanced, speech_enhanced_freq
421
422
423 # Denoise the noisy speech by 16-channel MLE speech enhancement
424 enhanced_speech, enhanced_speech_freq = MLE_multi_channel_speech_enhancement(speech_noisy_freq)
425
426
427
428 # ##### Evaluation
429
430 # var_emp evaluation
431 def Var_emp(enhanced_speech_freq_all):
432     """
433     Calculate the empirical variance for each frame in the enhanced speech frequency.
434
435     Parameters:
436     enhanced_speech_freq_all (ndarray): Array of shape (D, T, L) representing the enhanced speech frequency.
437
438     Returns:
439     var_emp_all (ndarray): Array of shape (L,) containing the empirical variance for each frame.
440     """
441     var_emp_all=np.zeros(16)
442     D,T,L = np.shape(enhanced_speech_freq_all)
443     for i in range(L):
444
445         for j in range(T):
446             var=np.sum((np.abs(clean_freq[:,j]-enhanced_speech_freq_all[:,j,i]))**2)
447
448             var_emp_all[i] = (1/(D*T))*np.sum(var)
449     return var_emp_all
450
451
452 #Enhancement for all channels only by calculating the mean of all channels
453 snr_all_mean=np.zeros(16)
454 speech_enhanced_all_mean=np.zeros((601600,16))
455 speech_enhanced_freq_all_mean=np.zeros((321,3761,16))
456
457 for i in range(16):
458     speech_enhanced_freq_all_mean[:, :, i]=np.mean(speech_noisy_freq[:, :, 0:i+1],axis=2)
459     speech_enhanced_all_mean[:, i] = librosa.istft(speech_enhanced_freq_all_mean[:, :, i], n_fft=frame_size, hop_length=
hop_size, window='hann')
460     snr_all_mean[i] = SNR(speech_enhanced_all_mean[:, i])
461
462
463 print("SNR of the enhanced speech is: ", SNR(enhanced_speech))
464 #plot the enhanced speech
465 #subplots
466 plt.figure(figsize=(10,5))

```

```

467 plt.subplot(2,1,1)
468 librosa.display.specshow(librosa.amplitude_to_db(abs(speech_noisy_freq[:,0]), ref=np.max), sr=sample_rate, hop_length=
hop_size, x_axis='time', y_axis='log', cmap='jet')
469 plt.title('Original noisy Speech in dB')
470 plt.xlabel('Time')
471 plt.ylabel('Frequency')
472 plt.colorbar(format='%+2.0f dB')
473
474 plt.subplot(2,1,2)
475 librosa.display.specshow(librosa.amplitude_to_db(abs(enhanced_speech_freq), ref=np.max), sr=sample_rate, hop_length=
hop_size, x_axis='time', y_axis='log', cmap='jet')
476 plt.title('Enhanced Speech in dB by MLE in multi-channel')
477 plt.xlabel('Time')
478 plt.ylabel('Frequency')
479 plt.colorbar(format='%+2.0f dB')
480 plt.subplots_adjust(hspace=0.5, wspace=0.5)
481 plt.show()
482
483 sf.write("enhce_mmse.wav", enhanced_speech, fs)
484
485
486 #Evaluation for 0-16 channels by SNR and var_emp
487 enhanced_speech_all_mle=np.zeros((len(enhanced_speech),16))
488 enhanced_speech_freq_all_mle=np.zeros((321,3761,16))
489 snr_all_mle=np.zeros(16)
490
491 for i in range(16):
492     enhanced_speech_all_mle[:,i], enhanced_speech_freq_all_mle[:,i] =
MLE_multi_channel_speech_enhancement(speech_noisy_freq[:,0:i+1])
493     snr_all_mle[i] = SNR(enhanced_speech_all_mle[:,i])
494
495
496 #plot the enhanced speech for all 16 channels
497 #subplots
498 plt.figure(figsize=(45,25))
499 for i in range(16):
500     plt.subplot(4,4,i+1)
501     librosa.display.specshow(librosa.amplitude_to_db(abs(enhanced_speech_freq_all_mle[:,i]), ref=np.max), sr=sample_rate,
hop_length=hop_size, x_axis='time', y_axis='log', cmap='jet')
502     plt.title(str(i+1)+' Channels')
503     plt.xlabel('Time')
504     plt.ylabel('Frequency')
505     plt.colorbar(format='%+2.0f dB')
506
507
508 #plot snr for each channel
509 print("SNR of the enhanced speech for all 0-16 channels is: ", snr_all_mle)
510 plt.figure(figsize=(10,5))
511 plt.scatter(range(1, 17), snr_all_mean)
512 plt.scatter(range(1, 17), snr_all_mle)
513 plt.title('SNR for 0-16 channels after multi enhancement by MLE')
514 plt.xticks(range(1, 17))
515 plt.xlabel('Number of Channels')
516 plt.ylabel('SNR')
517 plt.grid()
518 #legend
519 plt.legend(['Enhancement by averaging the channels','Enhancement by MLE'])
520 plt.show()
521
522
523 # calculate the Var_emp
524 var_emp_all_mle=Var_emp(enhanced_speech_freq_all_mle)
525 var_emp_all_mean=Var_emp(speech_enhanced_freq_all_mean)
526 print("var_emp of the enhanced speech for all 0-16 channels is: ", var_emp_all_mle)
527 plt.figure(figsize=(10,5))
528 plt.scatter(range(1, 17), var_emp_all_mean)
529 plt.scatter(range(1, 17), var_emp_all_mle)
530 plt.plot(range(1, 17), CRLB, color='red')
531 plt.title('var_emp for 0-16 channels after multi enhancement by MLE')
532 plt.xticks(range(1, 17))
533 plt.xlabel('Number of Channels')
534 plt.ylabel('var_emp')
535 plt.legend(['Enhancement by averaging the channels','Enhancement by MLE','CRLB'])
536 plt.grid()
537 plt.show()
538
539
540
541
542
543

```

```

544
545
546
547 #-----#
548 # ### 2.2 Denoise with MMSE
549 #-----#
550
551
552 #.....#
553 # #### 2.2.1 MMSE for single channel
554 #.....#
555
556 # #### Estimation
557 #single channel speech enhancement
558 def MMSE_speech_enhancement(speech_noisy_freq):
559     """
560     Enhances the given noisy speech signal using the MMSE (Minimum Mean Square Error) algorithm.
561
562     Parameters:
563     speech_noisy_freq (ndarray): Noisy speech frequency domain representation with shape (D, T), where D is the number of
564     frequency bins, T is the number of frames.
565
566     Returns:
567     speech_enhanced (ndarray): Enhanced speech in the time domain.
568     speech_enhanced_freq (ndarray): Enhanced speech frequency domain representation with shape (D, T), where D is the
569     number of frequency bins and T is the number of frames.
570     """
571
572     # parameters
573     max_gamma = 40 # maximum value of gamma
574     alpha = 0.98 # parameter for updating ksi using decision-direct
575
576     ksi_min = 10 ** (-15 / 10) # minimum value of ksi: -15dB
577                                #recommended in literature as -15dB
578
579     mu_VAD = 0.9 # parameter for VAD noise tracking
580     th_VAD = 5 # VAD decision threshold
581
582     phase_noisy = np.angle(speech_noisy_freq)
583     mag_noisy = np.abs(speech_noisy_freq)
584
585     D, T = np.shape(mag_noisy)
586
587     time_index_num = int(1//0.02) # initial noise length: 1s
588     noise_segment = mag_noisy[:, 0:time_index_num]
589     mag_noise = np.mean(np.abs(noise_segment), axis=1)
590     power_noise = mag_noise**2
591
592     mag_enhance = np.zeros([D, T])
593     alpha = alpha
594
595     for i in range(T):
596
597         # get energy spectrum and magnitude spectrum for each frame
598         mag_frame = mag_noisy[:, i]
599         power_frame = mag_frame**2
600
601         # get Signal-to-Noise Ratio (SNR) used for VAD calculation
602         SNR_VAD = 10 * np.log10(np.sum(power_frame) / np.sum(power_noise))
603
604         # calculate posterior SNR
605         gamma = np.minimum(power_frame / power_noise, max_gamma)
606
607         # calculate prior SNR
608         if i == 0:
609             ksi = alpha + (1 - alpha) * np.maximum(gamma - 1, 0)
610         else:
611             ksi = alpha * power_enhance_frame / power_noise + (1 - alpha) * np.maximum(gamma - 1, 0)
612             # limit the minimum value of ksi
613             ksi = np.maximum(ksi_min, ksi)
614
615         # update power_noise based on VAD
616         mu = mu_VAD
617         if SNR_VAD < th_VAD:
618             power_noise = mu * power_noise + (1 - mu) * power_frame
619         c = np.sqrt(np.pi) / 2
620
621         v = gamma * ksi / (1 + ksi)
622

```



```

623     # standard mmse
624     j_0 = sp.iv(0, v/2)
625     j_1 = sp.iv(1, v/2)
626     C = np.exp(-0.5 * v)
627     H = ((C * (v ** 0.5)) * C) / gamma * ((1 + v) * j_0 + v * j_1)
628
629     mag_enhance_frame = H * mag_frame
630     mag_enhance[:, i] = mag_enhance_frame
631
632     power_enhance_frame = mag_enhance_frame ** 2
633
634     speech_enhanced_freq = mag_enhance * np.exp(1j * phase_noisy)
635
636     speech_enhanced = librosa.istft(speech_enhanced_freq, n_fft=frame_size, hop_length=hop_size, window='hann')
637
638     return speech_enhanced, speech_enhanced_freq
639
640
641
642
643 # Try to use the first channel to do the speech enhancement
644 enhanced_speech, enhanced_speech_freq = MMSE_speech_enhancement(speech_noisy_freq[:, :, 0])
645
646
647 # #### Evaluation
648 print("SNR of the enhanced speech is: ", SNR(enhanced_speech))
649
650 #plot the enhanced speech
651 #subplots
652 plt.figure(figsize=(10,5))
653 plt.subplot(2,1,1)
654 librosa.display.specshow(librosa.amplitude_to_db(abs(speech_noisy_freq[:, :, 0]), ref=np.max), sr=sample_rate, hop_length=
hop_size, x_axis='time', y_axis='log', cmap='jet')
655 plt.title('Original noisy Speech in dB')
656 plt.xlabel('Time')
657 plt.ylabel('Frequency')
658 plt.colorbar(format='%+2.0f dB')
659
660 plt.subplot(2,1,2)
661 librosa.display.specshow(librosa.amplitude_to_db(abs(enhanced_speech_freq), ref=np.max), sr=sample_rate, hop_length=
hop_size, x_axis='time', y_axis='log', cmap='jet')
662 plt.title('Enhanced Speech in dB by MMSE')
663 plt.xlabel('Time')
664 plt.ylabel('Frequency')
665 plt.colorbar(format='%+2.0f dB')
666 plt.subplots_adjust(hspace=0.5, wspace=0.5)
667 plt.show()
668
669 sf.write("enhce_mmse.wav", enhanced_speech, fs)
670
671
672 #.....#
673 # #### 2.2.2 Multi channel MMSE
674 #.....#
675
676
677 # #### Estimation
678 #multi channel speech enhancement
679 def MMSE_MultiChannel_speech_enhancement(speech_noisy_freq):
680     """
681     Enhances the given noisy speech signal using the MMSE (Minimum Mean Square Error) algorithm.
682
683     Parameters:
684     speech_noisy_freq (ndarray): Noisy speech frequency domain representation with shape (D, T), where D is the number of
frequency bins, T is the number of frames.
685
686     Returns:
687     speech_enhanced (ndarray): Enhanced speech in the time domain.
688     speech_enhanced_freq (ndarray): Enhanced speech frequency domain representation with shape (D, T), where D is the
number of frequency bins and T is the number of frames.
689     """
690
691     speech_noisy_freq=np.mean(speech_noisy_freq,axis=2) #average over all channels
692
693     # parameters
694     max_gamma = 40 # maximum value of gamma
695     alpha = 0.98 # parameter for updating ksi using decision-direct
696
697     ksi_min = 10 ** (-15 / 10) # minimum value of ksi: -15dB
698                                #recommended in literature as -15dB
699

```

```

700 mu_VAD = 0.9 # parameter for VAD noise tracking
701 th_VAD = 5 # VAD decision threshold
702
703 phase_noisy = np.angle(speech_noisy_freq)
704 mag_noisy = np.abs(speech_noisy_freq)
705
706 D, T = np.shape(mag_noisy)
707
708 time_index_num = int(1//0.02) # initial noise length: 1s
709 noise_segment = mag_noisy[:, 0:time_index_num]
710 mag_noise = np.mean(np.abs(noise_segment), axis=1)
711 power_noise = mag_noise**2
712
713 mag_enhance = np.zeros([D, T])
714 alpha = alpha
715
716 for i in range(T):
717
718     # get energy spectrum and magnitude spectrum for each frame
719     mag_frame = mag_noisy[:, i]
720     power_frame = mag_frame**2
721
722     # get Signal-to-Noise Ratio (SNR) used for VAD calculation
723     SNR_VAD = 10 * np.log10(np.sum(power_frame) / np.sum(power_noise))
724
725     # calculate posterior SNR
726     gamma = np.minimum(power_frame / power_noise, max_gamma)
727
728     # calculate prior SNR
729     if i == 0:
730         ksi = alpha + (1 - alpha) * np.maximum(gamma - 1, 0)
731     else:
732         ksi = alpha * power_enhance_frame / power_noise + (1 - alpha) * np.maximum(gamma - 1, 0)
733         # limit the minimum value of ksi
734         ksi = np.maximum(ksi_min, ksi)
735
736     # update power_noise based on VAD
737     mu = mu_VAD
738     if SNR_VAD < th_VAD:
739         power_noise = mu * power_noise + (1 - mu) * power_frame
740
741     c = np.sqrt(np.pi) / 2
742
743     v = gamma * ksi / (1 + ksi)
744
745     # standard mmse
746     j_0 = sp.iv(0, v/2)
747     j_1 = sp.iv(1, v/2)
748     C = np.exp(-0.5 * v)
749     H = (((c * (v ** 0.5)) * C) / gamma) * ((1 + v) * j_0 + v * j_1)
750
751     mag_enhance_frame = H * mag_frame
752     mag_enhance[:, i] = mag_enhance_frame
753
754     power_enhance_frame = mag_enhance_frame ** 2
755
756     speech_enhanced_freq = mag_enhance * np.exp(1j * phase_noisy)
757
758     speech_enhanced = librosa.istft(speech_enhanced_freq, n_fft=frame_size, hop_length=hop_size, window='hann')
759
760     return speech_enhanced, speech_enhanced_freq
761
762
763
764
765 enhanced_speech, enhanced_speech_freq = MMSE_MultiChannel_speech_enhancement(speech_noisy_freq)
766
767
768
769
770 # #### Evaluation
771 print("SNR of the enhanced speech is: ", SNR(enhanced_speech))
772
773
774 #plot the enhanced speech
775 plt.figure(figsize=(10,5))
776 plt.subplot(2,1,1)
777 librosa.display.specshow(librosa.amplitude_to_db(abs(speech_noisy_freq[:, :, 0]), ref=np.max), sr=sample_rate, hop_length=
hop_size, x_axis='time', y_axis='log', cmap='jet')
778 plt.title('Original noisy Speech in dB')

```

```

779 plt.xlabel('Time')
780 plt.ylabel('Frequency')
781 plt.colorbar(format='%+2.0f dB')
782
783 plt.subplot(2,1,2)
784 librosa.display.specshow(librosa.amplitude_to_db(abs(enhanced_speech_freq), ref=np.max), sr=sample_rate, hop_length=
hop_size, x_axis='time', y_axis='log', cmap='jet')
785 plt.title('Enhanced Speech in dB by MMSE in multi-channel')
786 plt.xlabel('Time')
787 plt.ylabel('Frequency')
788 plt.colorbar(format='%+2.0f dB')
789 plt.subplots_adjust(hspace=0.5, wspace=0.5)
790 plt.show()
791
792 sf.write("enhce_mmse.wav", enhanced_speech, fs)
793
794
795
796 #Evaluation for 0-16 channels by SNR and var_emp
797 enhanced_speech_all_mmse=np.zeros((len(enhanced_speech),16))
798 enhanced_speech_freq_all_mmse=np.zeros((321,3761,16))
799 snr_all_mmse=np.zeros(16)
800
801 for i in range(16):
802     enhanced_speech_all_mmse[:,i], enhanced_speech_freq_all_mmse[:, :,i] =
MMSE_MultiChannel_speech_enhancement(speech_noisy_freq[:, :,0:i+1])
803     snr_all_mmse[i] = SNR(enhanced_speech_all_mmse[:,i])
804
805 #plot the enhanced speech for all 16 channels
806 #subplots
807 plt.figure(figsize=(45,25))
808 for i in range(16):
809     plt.subplot(4,4,i+1)
810     librosa.display.specshow(librosa.amplitude_to_db(abs(enhanced_speech_freq_all_mmse[:, :,i]), ref=np.max), sr=
sample_rate, hop_length=hop_size, x_axis='time', y_axis='log', cmap='jet')
811     plt.title(str(i+1)+' Channels')
812     plt.xlabel('Time')
813     plt.ylabel('Frequency')
814     plt.colorbar(format='%+2.0f dB')
815
816
817 #plot snr for each channel
818 print("SNR of the enhanced speech for all 0-16 channels is: ", snr_all_mmse)
819 plt.figure(figsize=(10,5))
820 plt.scatter(range(1, 17), snr_all_mean)
821 plt.scatter(range(1, 17), snr_all_mle)
822 plt.scatter(range(1, 17), snr_all_mmse)
823
824 plt.title('SNR for 0-16 channels after multi enhancement by MMSE')
825 plt.xticks(range(1, 17))
826 plt.xlabel('Number of Channels')
827 plt.ylabel('SNR')
828 plt.grid()
829 #legend
830 plt.legend(['Enhancement by averaging the channels', 'Enhancement by MLE', 'Enhancement by MMSE'])
831 plt.show()
832
833
834 # calculate the Var_emp
835 var_emp_all_mmse=Var_emp(enhanced_speech_freq_all_mmse)
836 var_emp_all_mean=Var_emp(speech_enhanced_freq_all_mean)
837 print("var_emp of the enhanced speech for all 0-16 channels is: ", var_emp_all_mmse)
838 plt.figure(figsize=(10,5))
839 plt.scatter(range(1, 17), var_emp_all_mean)
840 plt.scatter(range(1, 17), var_emp_all_mle)
841 plt.scatter(range(1, 17), var_emp_all_mmse)
842 plt.plot(range(1, 17), CRLB, color='red')
843 plt.title('var_emp for 0-16 channels after multi enhancement by MMSE')
844 plt.xticks(range(1, 17))
845 plt.xlabel('Number of Channels')
846 plt.ylabel('var_emp')
847 plt.legend(['Enhancement by averaging the channels', 'Enhancement by MLE', 'Enhancement by MMSE', 'CRLB'])
848 plt.grid()
849 plt.show()
850
851 # calculate the Var_emp with only MMSE
852 plt.figure(figsize=(10,5))
853 plt.scatter(range(1, 17), var_emp_all_mmse, color='green')
854 plt.plot(range(1, 17), CRLB, color='red')
855 plt.title('var_emp for 0-16 channels after multi enhancement by MMSE')
856 plt.xticks(range(1, 17))

```

```
857 plt.xlabel('Number of Channels')
858 plt.ylabel('var_emp')
859 plt.legend(['Enhancement by MMSE', 'CRLB'])
860 plt.grid()
861 plt.show()
862
863
864
865 #-----END-----#
```