



IMPROVEMENT IN SIGNAL QUALITY THROUGH MEDIAN BASE FILTERING

Junaid Masood¹, Sheeraz Ahmed², Asim Ali³, Ubaid Ullah⁴
Said-ul-Abrar⁵, Muhammad Tayyab⁶, Samhita Priyadarsini Gundala⁷

^{1,2}IQRA National University, Peshawar, Pakistan

^{3,4}Brains Institute, Peshawar, Pakistan

⁵Institute of Computer Science and Information Technology,
Agriculture University Peshawar, Pakistan

⁶Career Dynamics Research Centre, Peshawar, Pakistan

⁷Vignn's Foundation for Science, Tech and Research, Guntur, India

Email: ¹junaidkttk@gmail.com, ²sheerazahmed306@gmail.com,
³asim.aries.664@gmail.com, ⁴ubaid.pk49@gmail.com, ⁵saidulabrar@aup.edu.pk,
⁶tayyab.uetpeshawar@gmail.com, ⁷spd4india@gmail.com

<https://doi.org/10.26782/jmcms.2021.09.00003>

(Received: July 5, 2021; Accepted: September 1, 2021)

Abstract

Speech signal segmental framing and the scaling factor is basis for the speech recognition process as first step. The next followed step is existing noise reduction in the recognized speech signal for quality improvement. In this work, the noise reduction is done using newly proposed adaptive median based filtering. Comparison of the observations based on adaptive median filtering with Minimum Mean-Square Error Short-time Spectral Amplitude (MMSE-STSA) and Minimum Mean-Square Error (MMSE) based noise reduction reveal a list of worthy to mention relevant observations. The drawn conclusion also accumulates possible contributions by the proposed adaptive median based filtering technique. Lastly is mentioning of Signal-to-noise ratio (SNR) as the primary metric for observations collection for the newly proposed adaptive median based filtering technique analysis.

Keywords : Filters, Speech Signal, Signal to Noise Ratio, Mean Square Error, Scaling Factor

I. Introduction

Intelligibility is the ability of listener to understand the content of speech signals. The intelligibility varies from person to person but the minimum speech reception threshold is fifty (50) percent [I]. Whereas quality is the comparison of system output speech signal with the original input signal. Quality may vary from person to person. Some people mark low-quality speech as a high quality and other people mark it as low quality. So quality is a measure of speech signal better effect on

listener's ear. Both of them are independent attributes of the speech signal. Mostly both are inversely related to with each other [II], [III].

Speech signal comprises of both useful and unwanted signals. A common definition, the unwanted signal part of the speech signal is called noise. In presence of background noise, the speech signal becomes degraded by noise. An important worthy to note relevant discrepancy is domination of speech signal by the noise signal and this occurs with occurrence of negative SNR which results in zero intelligibility. However, the noise can be reduced easily if attributes of noise or speech signal are known. Achieving limited speech signal distortion and noise quantity, different speech enhancement techniques are used. The list includes Speech enhancement using MMSE filter in wavelet domain [IV], Speech enhancement using Empirical Mode Decomposition method [V], adaptive median based filtering for quality improvement [VI]. Enhancement technique works when the attributes of noise are known. List of different types of noise includes pink noise, babble noise, grey noise and additive noise [VII]. Additive noise is the background noise which has a resemblance with speech signal. The best example of additive noise is additive white Gaussian noise. Additive white Gaussian (AWGN) noise is used as a basis noise for adaptive thresholding using median base filtering for quality improvement of signals [VIII]. Improve signal has many application such as hearing impaired devices. Real time analysis of speech enhancement is too much expensive and not in range.

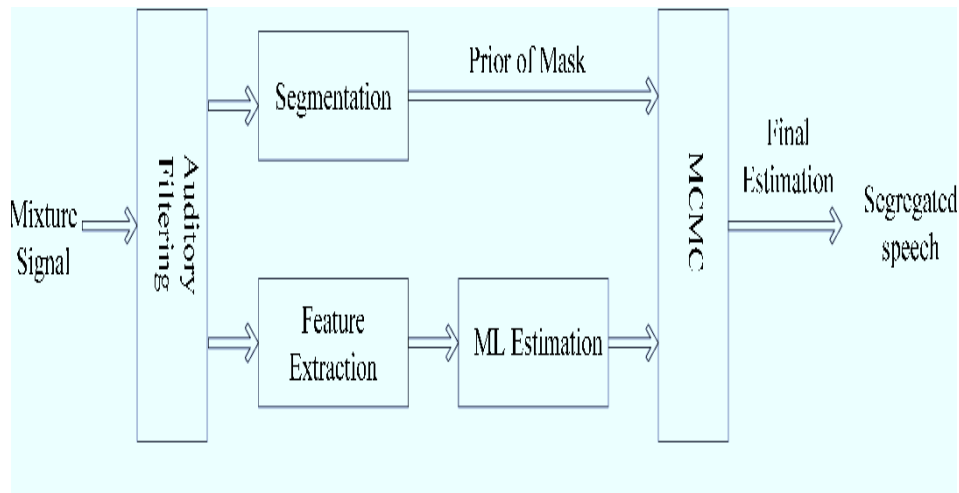


Figure 1. Block diagram of binary mask method

Alternatively the mathematical analysis has solved the purpose with significant quality which can be seen in subsequent sections of this research. The speech enhancement can be done using different proposed method. These methods comprises of parametric and non-parametric method [IX]. The parametric model method does not need the posterior information [X]. The priori model information helps in the process of speech enhancement with reduced noise and better quality. For single channel the non-parametric method use the silent frame as frame of reference of speech signal and estimate the noisy observation. However, for multiple channels

one source is assumed to be reference source for noise estimation. Single channel is more effective because of its clear implementation and least cost [XI], [XII]. Also, it is important to note that non-parametric approach is flexible and therefore allow changes in the used parameters over a cartage specified range therefore not restricted to fixed number of parameters as the case in parametric approach.

Here in this case, adaptive thresholding using median base filtering algorithm is based on non-parametric and single channel speech enhancement. The algorithm needs the noise observation of silent frame of speech signal which is achieved using single microphone. In the following text the survey of few speech enhancement algorithms provides the basis for proposed algorithm.

II. Review of Literature

Enhancement of Speech comprises of 3 different algorithm stages [XIII]. The algorithms categorization depends upon pros and cons in their operation. Four (04) different speech enhancement categorizations are binary mask, spectral subtraction, subspace algorithm, statistical based algorithms.

Spectral Subtraction method

Spectral Subtraction needs the continuous output for reduction of noise. The spectral subtraction method for the first time induces by the Weiss using in the correlation domain. Later Boll uses it for in the Fourier transform domain. The spectral subtraction method is based on subtraction of noise from the corrupted signal [XIV]. During speech enhancement using spectral subtraction method first determine the noise. The noise can be determined from the pause region of the speech signal. And subtraction of the highlighted noise is carried out from the original signal. However, this algorithm can distort the useful information from the speech signal.

Binary mask method

The binary mask method use the binary value for noise reduction. In binary masking method original signal is added with selected frequencies of corrupted signal. After addition of noise the time domain signal is converted into the time-frequency domain. The frequency domain analysis of signal provides the frequency information of speech signal. The frequency domain not provides the change of energy of speech signal at specific point. So it needs the use of time domain. As the speech signal is non-stationary so it need the use of both frequency and time domain. Classical method of frequency domain is Fourier analysis but for the binary masking method represented by Short time Fourier transform [XVI]. In time and frequency domain the binary masking & signal-to-noise ratio thresholding criterion is created. The difference between the target energy signal and energy of mask signal is greater than local threshold criterion than assign value 1. If it is not greater than the local threshold criterion assign zero.

$$BM(t, f) = \begin{cases} 1 & \text{if } TE(t, f) - m(t, f) > LTC \\ 0 & \text{otherwise} \end{cases} \quad (1)$$

The binary masking value is multiplied with the magnitude of FFT. The output from the FFT is converted into inverse FFT and the result is weighted by using

overlap add method and apply windowing and get output frames in row frame matrix. The block diagram is shown in figure 2.

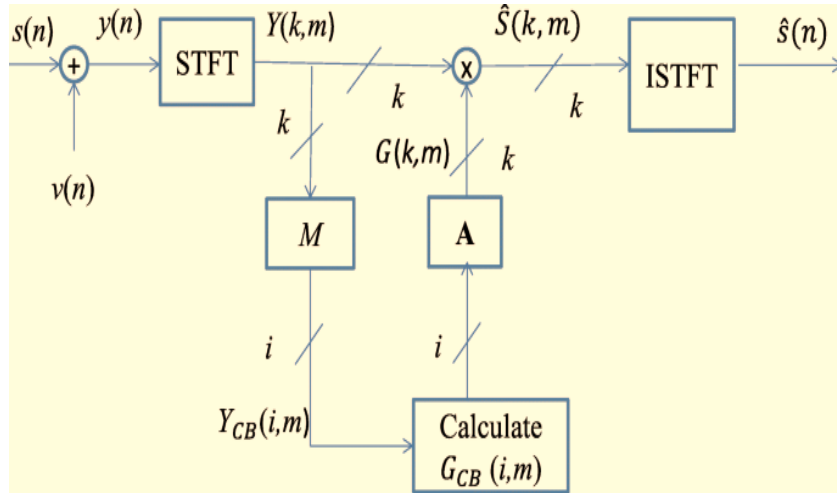


Figure 2. Flow Chart for Adaptive binary mask method

Subspace method

In speech enhancement using subspace method noisy speech signal is the sum of speech signal subspace and the additive noise subspace [XVII]. As a first step speech signal is converted into subspace covariance matrix for the noisy speech signal. The noisy speech signal is the sum of the covariance matrix of speech signal and covariance matrix of additive noise signal. Then estimation of the speech signal subspace covariance matrix is next followed step. The Eigen value decomposition method is used for finding the diagonal element greater than zero of covariance matrix of signal subspace and noisy subspace. Then get the signal subspace and noisy subspace spanning the Eigen vector against the Eigen values. The diagonal element of covariance matrix of speech signal provides the dimension of covariance matrix speech signal. After this conversion of the product of Eigen value and Eigen vector into frequency domain occurs. In frequency domain the power spectral density of Eigen Product divided by the number of samples of the frames is found. Calculation of the tonal and non-tonal component using the thresholding masking is the next in-line step of the process. For minimizing the noise in the enhance speech signal the noise frame with the maximum value of the noise is normalized in the frame. Now after normalization is the time for back into Eigen domain conversion from converted frequency domain. Calculation of the gain and the product of noise frames and filter provide the estimation of speech enhancement.

Statistical based Method

Initially for statistical base speech enhancement method the Ephraim used the MMSE STSA method [XVIII]. MMSE STSA is primarily used for the noise estimation. The noisy speech signal is divided into magnitude and phase part. The estimation of speech signal depends on the product of response of the MMSE STSA

Junaid Masood et al

estimator and frame index with sampling frequency. The gain or response of MMSE depends on the priori SNR. The priori SNR depend on the ratio of speech variance to variance of noise known as posterior SNR or instantaneous SNR. Since it is here that conversion of non-stationary speech signal into stationary occurs by estimating the one fixed frequency and time required in the stationary signal. Point worthy to note is improved quality signal results by the conversion into stationary signal using the newly proposed clean adaptive thresholding that is based on median based filtering.

III. MMSE method details

For better explanation of the proposed median based filtering technique in this paper, the relevant description of the preliminary used MMSE method needs to be considered in context of stationary and non-stationary speech signal.

The conventional Minimum mean square error (MMSE) method which is not as useful for non-stationary noise as it is usually effective for stationary noise. One of the main reasons for the mentioned inflexibility is tracking the noise variance along time dimension and typical focus on statistical characteristics of the noise [XVIII].

Role of Noise

The additive white Gaussian noise is added as a noise with the clean speech signal, mathematically noisy speech signal can expressed as

$$n(t) = c(t) + r(t) \quad (2)$$

Hamming of Speech Signal

The noisy speech signal is further divided into frames using hamming windowing.

$$n(t) = \sum_{j=1}^c Q_j(t) + l_m(t) \quad (3)$$

MMSE Noise Estimation

For different frame of $Q_j(t)$, the window noise estimation can be done using adaptive parameter of H_j . Successful preprocessing lead to adaptive median base filtering. For the nose estimation using different frame u at sampling frequency F_s in presence of adaptive parameter H_j is mathematically expressed using gamma function as basis function given below

$$\widetilde{g_j} = \Gamma[Q_j(t); H_j] \quad (4)$$

MMSE Estimation

Posterior signal-to-noise ratio is the ratio of square of frames of noisy signal and square of noisy signal using sampling frequency at instantaneous frame u .

$$SNR_{\text{post}} = \frac{Q^2(F_s, u)}{D^2(F_s, u)} \quad (5)$$

Priori SNR can be obtain from posterior SNR, while prior SNR is the sum of maximum value between posterior SNR and zero weight with $(1-\beta)$ and square of estimation noisy spectral proceeding frame $(u-1)$ at sampling frequency F_s with square of estimation spectral proceeding frame $(u-1)$ at sampling frequency F_s .

$$\text{SNR}_{\text{priori}} = \alpha + (1 - \alpha) \max(\text{SNR}_{\text{post}}(f_s, u), 0) \quad (6)$$

MMSE filter response is the ratio of priori signal-to-noise of frame u at sampling frequency f_s to priori signal-to-noise of frame u at sampling frequency plus one. Can be expressed mathematically

$$L(f_s, u) = \frac{\text{SNR}_{\text{priori}}}{\text{SNR}_{\text{post}} + 1} \quad (7)$$

Estimated noise spectral is the product of MSSE filter system response $L(f_s, u)$ and $Q_j(f_s, u)$.

$$\widetilde{R}_j = L(f_s, u) Q_j(f_s, u) \quad (8)$$

IV. Overview of Proposed Method

Adapting filtering using median based filtering

The proposed adaptive median based filtering minimizes the effect of original speech signal distortion. And show more smooth effect during reconstruction of estimated speech signal.

The noisy speech signal is divided into three parts. First comprises of start region with fewer information of speech, so that a noise can be easily subtracted. Due to low distortion ratio of original signal, it is scale with 0.4. The difference between the frames of region 1 and median of frames 1 region and overall median product with scaling factor of 0.4 gives the value of noise region of σ_1 . So for the first region of noise estimation value can be expressed mathematically

$$\sigma_1 = 0.4 \times \text{Median}^*(Q_1(t) - \text{Median}^* Q_1(t)) \quad (9)$$

And this is same for region 2 but requires more information of speech. To get smoothing effect its scaling is weighted with 0.8. If any low amplitude noise signal present it can be amplified with larger weight.

$$\sigma_2 = 0.8 \times \text{Median}^*(Q_2(t) - \text{Median}^* Q_2(t)) \text{ same for the third region so the estimated noise level is}$$

$$\sigma_3 = 0.4 \times \text{Median}^*(Q_3(t) - \text{Median}^* Q_3(t))$$

Take the average of these three regions

$$\sigma = \frac{\sigma_1 + \sigma_2 + \sigma_3}{3} \quad (10)$$

And finally the value of σ is use in the equation of Donoho

$$A = \sqrt{2 \log L} \sigma \quad (11)$$

In above equation L is the length of original speech signal. This whole process is shown in form of block diagram in figure 2 and subsequently in figure 3 together.

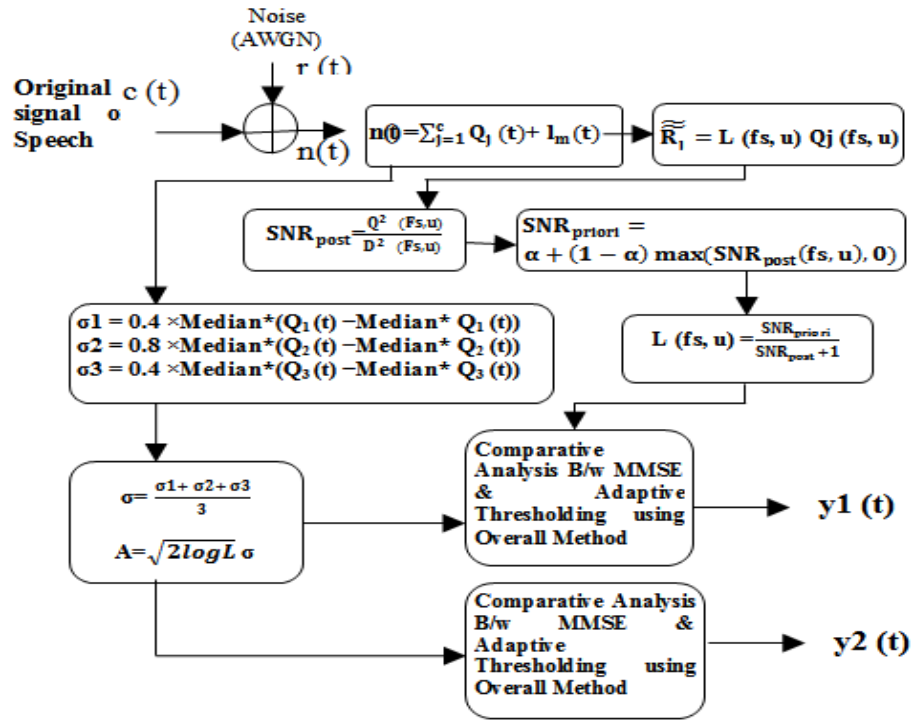


Figure 3. Adaptive thresholding for median based filtering

V. Experimental Results

Random test are perform for comparative analysis of MMSE and adaptive median based filtering. For initial SNR the output of change in SNR is calculated for the SNR range between -2 dB to -18 dB shown in figure 4 and figure 5.

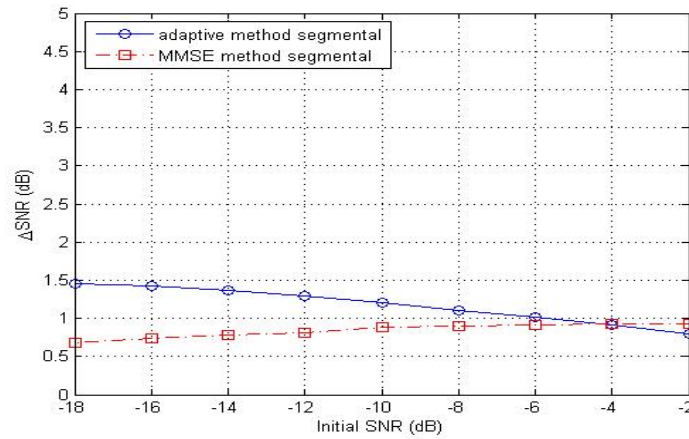


Figure 4. Segmental result of Δ SNR comparison with MMSE for SNR range between -2dB to -18 dB.

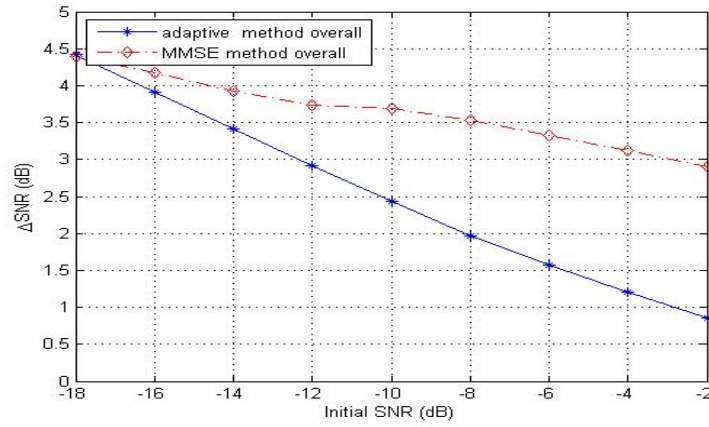


Figure 5. Overall result of Δ SNR comparison with MMSE for SNR range between -2dB to -18 dB.

The adaptive thresholding using median based filtering show better performance between the ranges of SNR -14 dB -8dB for the input SNR range -2dB to -18dB. As compare to segmental result the proposed method show better performance result for overall analysis. Greater the input range of SNR the better the performance result of proposed method.

Another set of experiment for SNR range 4dB to -16dB. Change in SNR for the input value of initial SNR is calculated shown in figure 6 [XVIX] for various types of MMSE speeches.

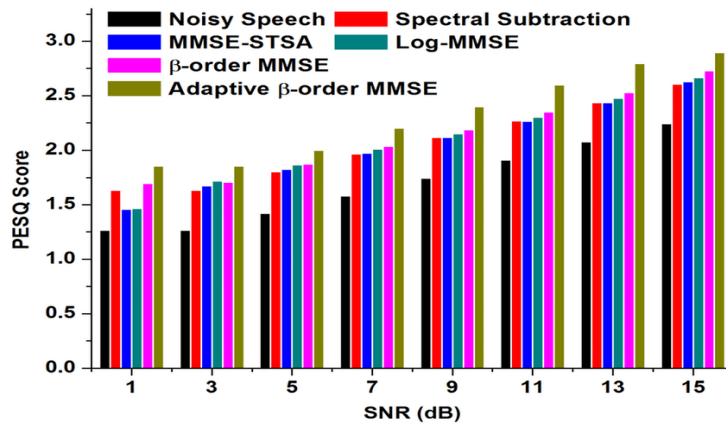


Figure 6: Segmental result of Δ SNR comparison with MMSE for SNR range between 4dB to -16dB.

VI. Conclusion

The research shows that the noise reduction is done using newly proposed adaptive median based filtering. Comparison of the observations based on adaptive median filtering with Minimum Mean-Square Error Short-time Spectral Amplitude (MMSE-STSA) and Minimum Mean-Square Error (MMSE) based noise reduction

Junaid Masood et al

reveal a list of worthy to mention relevant observations. In terms of overall signal Δ SNR as compared to adaptive thresholding MMSE filtering technique offers improved performance. The overall performance of adaptive thresholding is better in wide range of SNR and show low performance on short range of SNR. On short range of SNR the MMSE performance is better than adaptive thresholding.

.

Conflict of Interest

The authors declare that they have no conflicts of interest to report regarding the present study.

References

- I. Brown, A., S. Garg, and J. Montgomery, Automatic and Efficient Denoising of Bioacoustics Recordings Using MMSE STSA. IEEE Access, 2018. 6: p. 5010-5022.
- II. Di Liberto, G.M., et al., Atypical cortical entrainment to speech in the right hemisphere underpins phonemic deficits in dyslexia. NeuroImage, 2018.
- III. Djaziri-Larbi, S., et al., Watermark-Driven Acoustic Echo Cancellation. IEEE/ACM Transactions on Audio, Speech, and Language Processing, 2018. 26(2): p. 367-378.
- IV. Fu, J., L. Zhang, and Z. Ye, Supervised monaural speech enhancement using two-level complementary joint sparse representations. Applied Acoustics, 2018. 132: p. 1-7.
- V. Heitkaemper, Jens, Joerg Schmalenstroer, Joerg Ullmann, Valentin Ion, and Reinhold Haeb-Umbach. "A Database for Research on Detection and Enhancement of Speech Transmitted over HF links." arXiv preprint arXiv:2106.02472 (2021).
- VI. Heese, F., et al. Selflearning codebook speech enhancement. in Speech Communication; 11. ITG Symposium; Proceedings of. 2014. VDE.
- VII. Kandagatla, R.K. and P. Subbaiah, Speech enhancement using MMSE estimation of amplitude and complex speech spectral coefficients under phase-uncertainty. Speech Communication, 2018. 96: p. 10-27.
- VIII. Khaldi, K., A.-O. Boudraa, and A. Komaty, Speech enhancement using empirical mode decomposition and the Teager–Kaiser energy operator. The Journal of the Acoustical Society of America, 2014. 135(1): p. 451-459.

- IX. Khaldi, K., et al., Speech enhancement via EMD. EURASIP Journal on Advances in Signal Processing, 2008. 2008(1): p. 873204.
- X. Kuortti, J., J. Malinen, and A. Ojalampi, Post-processing speech recordings during MRI. Biomedical Signal Processing and Control, 2018.
- XI. Masood, J., Shahzad, M., Khan, Z.A., Akre, V., Rajan, A., Ahmed, S. and Masood, F., 2020, November. Effective Classification Algorithms and Feature Selection for Bio-Medical Data using IoT. In 2020 Seventh International Conference on Information Technology Trends (ITT) (pp. 42-47). IEEE.
- XII. Michelsanti, Daniel, Zheng-Hua Tan, Shi-Xiong Zhang, Yong Xu, Meng Yu, Dong Yu, and Jesper Jensen. "An overview of deep-learning-based audio-visual speech enhancement and separation." IEEE/ACM Transactions on Audio, Speech, and Language Processing (2021).
- XIII. Nabi, W., et al., A dual-channel noise reduction algorithm based on the coherence function and the bionic wavelet. Applied Acoustics, 2018.
- XIV. Rao, C.V.R., M.R. Murthy, and K.S. Rao. Speech enhancement using perceptual Wiener filter combined with unvoiced speech—A new scheme. in Recent Advances in Intelligent Computational Systems (RAICS), 2011 IEEE. 2011. IEEE.
- XV. Tabassum Feroz, Uzma Nawaz. : ‘SUPPRESSION OF WHITE NOISE FROM THE MIXTURE OF SPEECH AND IMAGE FOR QUALITY ENHANCEMENT’. *J. Mech. Cont. & Math. Sci., Vol.-16, No.-7, July (2021) pp 67-78*. DOI : 10.26782/jmcms.2021.07.00006
- XVI. Wang, X., et al., A comparison of recent waveform generation and acoustic modeling methods for neural-network-based speech synthesis. preprint arXiv:1804.02549, 2018.
- XVII. Wiem, B., P. Mowlae, and B. Aicha, Unsupervised single channel speech separation based on optimized subspace separation. Speech Communication, 2018. 96: p. 93-101.
- XVIII. Yang, Fan, Ziteng Wang, Junfeng Li, Risheng Xia, and Yonghong Yan. "Improving generative adversarial networks for speech enhancement through regularization of latent representations." Speech Communication 118 (2020): 1-9.
- XIX. Yilmaz, O. and S. Rickard, Blind separation of speech mixtures via time-frequency masking. IEEE Transactions on signal processing, 2004. 52(7): p. 1830-1847.