

Noise Power Spectral Density Estimation based on Maximum a Posteriori and Generalized Gamma Distribution

Xin Dang, Takayoshi Nakai and Md. Iqbal Aziz Khan

*Information Science and Technology, Graduate School of Science and Technology,
Shizuoka University, Hamamatsu, 432-8561, Japan*
Email: f5045013@ipc.shizuoka.ac.jp, tdtnaka@ipc.shizuoka.ac.jp

Abstract

Noise power spectral density (PSD) estimation is a crucial part of speech enhancement system due to its contributory effect on the quality of the noise reduced speech. A novel estimation method for color noise PSD on the basis of an assumption of generalized Gamma distribution and maximum a posteriori (MAP) criterion is proposed. In the experiment, generalized Gamma PDF which is a natural extension of the Gaussian modeling of a non-white components distribution is found best fitting in four types of color noises compared with Laplace, Rayleigh distributions. After that MAP noise estimators based on the reported generalized Gamma PDF models are competed with Minimum Statistics (MS), minimum mean square error (MMSE) based PSD estimation and Maximum Likelihood estimation (MLE) noise tracking methods in evaluations. The performance of the proposed noise estimations are good as demonstrated by log error, segmental SNR and PESQ measures when they are integrated with the speech enhancement technique.

Keywords: Speech enhancement, noise PSD estimation, MAP, MLE **1**.

Introduction

An estimation of the power spectral density (PSD) of noise is a crucial part to retrieve speech in a noisy environment. Thus, the performance of the man-machine communication system in a noisy environment highly depends on the accurate estimation of the unknown noise model. In the recent years, many methods have been proposed and the estimation of noise from a noisy speech remains a challenging task, especially due to the wide variety of non-stationary and non-white nature of environmental noises.

The PSD of the most real world noises change rapidly over time due to its non-stationary nature, for that, an efficient noise estimation method must have real time update capability and estimation accuracy. The worse accuracy of the noise (over or under) estimation might lead to a reduced intelligibility of original speech or generate an unnecessary amount of residual noise due to an inherent mismatch between the original and the estimated noise.

The most common approach for the estimation of noise PSD is to exploit speech presence by means of a Voice Activity Detector (VAD) [1] and the speech/pause detection plays the major role in the performance of the whole system. This approach cannot update the noise estimation promptly. However, these systems can perform well for voiced speech and high Signal-to-Noise Ratio (SNR), but their performance degrades with unvoiced speech at low SNR.

To improve the estimation of noise PSD, several real time approaches have been proposed during the last decade. One of the most famous methods is noise PSD estimation based on Minimum Statistics (MS) [2], which can update the noise power spectrum directly from the noisy speech, even in a non-stationary noisy environment. The noise PSD estimation using

MS is based on the assumption that within the observed time-span, there is a silent part that is at least a small fraction of the total time-span. The spectral noise power is then obtained from the minimum values of the estimated power spectrum of the noisy signal. But this method may attenuate low energy phonemes occasionally, and the minimal search length should be set at least to the length of the smoothed speech components to avoid over-estimation, which restricts the tracking capability of the noise estimator in case of varying noise spectrum[4].

Recently, noise power estimation methods based on minimum mean square error (MMSE) [4-5] and Maximum Likelihood estimation (MLE) [8] have been proposed. These methods are based on the assumption of noise DFT components following a Gaussian distribution, and achieve good online noise-tracking capability for non-stationary noises with low complexity. However, the DFT components of color noise hardly follow Gaussian distribution. The accuracy of such estimation degrades in color noisy environments.

To estimate the PSD of unknown noise, a study on the real distribution of speech and noise spectrum is essential because the optimal spectral estimator is based on assumed appropriate statistical model and criterion. A distribution of speech spectrum is discussed by Yariv Ephraim [1]. It concludes real statistical model of speech spectrum seems to be inaccessible, the validity of speech model can be judged a posteriori based on the results obtained. In his paper, he used a Gaussian model for both of the speech and noise. An advanced study is presented by Martin and Lotter [6]. These studies explored speech distribution data with 1-hour duration, in which only speech with a high and narrow band-wide SNR interval (for example 19–21 dB) was selected. The results show the gamma model is better to fit the real statistical model in comparison with Laplace and Gaussian model. However, since the speech is a non-stationary signal, the distributions of voiced and unvoiced sound were found to be quite different. Therefore it is difficult to improve estimation of speech using a single distribution theory. Since the distributions of noise spectral amplitudes change slightly. Therefore, an improved estimation methods based on noise distribution is expected.

In this paper, we investigate a new noise PSD estimation method for color noise based on a generalized gamma model rather than assuming a Gaussian distribution. The estimation of noise PSD can be derived by MAP criterions based on the generalized gamma probability density function (PDF). In addition, the parameters of the underlying PDF can be optimally fitted to the real distribution of four different typical color noise spectral amplitudes at each frequency bin. Using this statistical model, an accurate and computationally efficient noise estimator can be established. By integrating with the Wiener filter, the proposed algorithm shows noticeable efficiency compared with other recently developed speech-enhancement methods.

This paper is organized as follows. In Section 2, we analyze the real distributions of noise spectral amplitudes by making a comparison with Rayleigh, the gamma and Laplace CDFs. From this comparison, we propose two new noise PSD estimation methods based on the generalized gamma PDFs. In Section 3, we test the proposed methods integrated with an improved Wiener filter using highly non-stationary noisy speech and compare it with some popular noise PSD-tracking algorithms, such as MS [2], MMSE [4], and MLE [8]. The results in terms of objective measurements are described in subsection 3.3. Then we give a conclusion in Section 4.

2. Noise Magnitude Models and Statistical Analysis

2.1. The Distribution of Noise Spectral Amplitude

We assumed that the mixing noise is additive and that speech $s(l)$ and noise $n(l)$ signals are uncorrelated. The noisy speech $y(l)$ becomes

$$y(l) = s(l) + n(l) \tag{1}$$

After discrete Fourier transform (DFT) of the noisy speech, it can be written as

$$Y(\lambda, k) = S(\lambda, k) + N(\lambda, k) \tag{2}$$

where $Y(\lambda, k)$ denotes the DFT coefficient of noisy speech at the frequency index k and frame index λ .

For white noise, both real and imaginary parts of the Fourier coefficients are distributed independently and identically and has a Gaussian distribution. This allows for Rayleigh-distributed noise spectral amplitudes. The PDF with parameter δ_{Ray} can be written as

$$f_{Ray}(x) = \frac{2x}{\delta_{Ray}^2} \exp\left(-\frac{x^2}{\delta_{Ray}^2}\right) \tag{3}$$

and $\delta_{Ray} = N_{mean}(k)$, where $N_{mean}(k)$ denotes the mean of noise spectral amplitudes of the frequency bin k . For non-white noise, the spectral amplitudes are not uniform. Therefore, a more accurate PDF is required. The real PDF of noise spectral is near to Rayleigh PDF in shape but need proper adjustment to fit different kinds of color noises. Therefore gamma and Laplace PDF are taken into account. If we assume that the noise spectrum follows a gamma distribution, then the PDF can be written as

$$f_{gamma}(x) = \frac{1}{\delta_{gam} \Gamma(\nu)} \left(\frac{x}{\delta_{gam}}\right)^{\nu-1} \exp\left(-\frac{x}{\delta_{gam}}\right) \tag{4}$$

where $\delta_{gam} = N_{mean}(k)$, $\nu = N^2_{mean}(k)/V_N(k)$, according to the definition of the underlying gamma distribution, $V_N(k)$ denotes the variance of the noise spectral amplitudes of k -th frequency bin. If we assume that the noise spectrum follows Laplace distribution, the PDF can be written as

$$f_{Laplace}(x) = \frac{1}{\delta_{Lap}} \exp\left(-\frac{|x|}{\delta_{Lap}}\right) \tag{5}$$

where $\delta_{Lap} = N_{m}$

Figure 1 shows an example of the comparison of the analytical CDFs and real CDF of noise spectral amplitudes at 1 kHz and the SNR is set as 5dB. It is difficult to determine which curve is the best. The Rayleigh and gamma CDFs are close to that of real CDF of noise spectral, the Laplace CDF is far away from the real noise CDF. The Rayleigh CDF is good fitted in the fan noise situation, which is close to the white noise. However at other non-white noise environments the Rayleigh model show worse matching. This suggests that the assumption of the Gaussian distribution model is not suiting for the color noise any more.

Furthermore, the distribution of the color noise spectral amplitude changes in the frequency domain. In order to show the detailed noise spectral distribution in the frequency domain, Kolmogorov–Smirnov (KS) statistics is introduced, the KS statistics is defined as

$$KS = \max_{1 \leq k \leq N} |F(N(k)) - F_{Ana}(k)| \tag{6}$$

where $F(N(k))$ and $F_{Ana(k)}$ denotes the Cumulative distribution function (CDF) of real noise spectral amplitude and analytical distribution, which includes $f_{Ray}(x)$, $f_{gamma}(x)$ and $f_{Laplace}(x)$. In this experiment, we use 50 second noise signals at a SNR of 5 dB to calculate $F(N(k))$ with a frame length of 16 ms, achieving a 50% overlap between adjacent frames. The parameters of the distribution models are computed using $V_N(k)$ and $N_{mean}(k)$ as shown in Eqs. (3) (4) and (5).

Figure 2 shows the KS statistics of four types of noises. The KS statistics of Laplace is the largest for all noises. In Figure 2(a) and (c), the KS statistics of Rayleigh model are smaller than Laplace but larger than the gamma model. Gamma model is superior for the train, traffic and factory noises at most frequency. However, in Figure 2(c) and (d), the KS statistics suggest that the Rayleigh model is close even better for the factory and fan noises at some frequency.

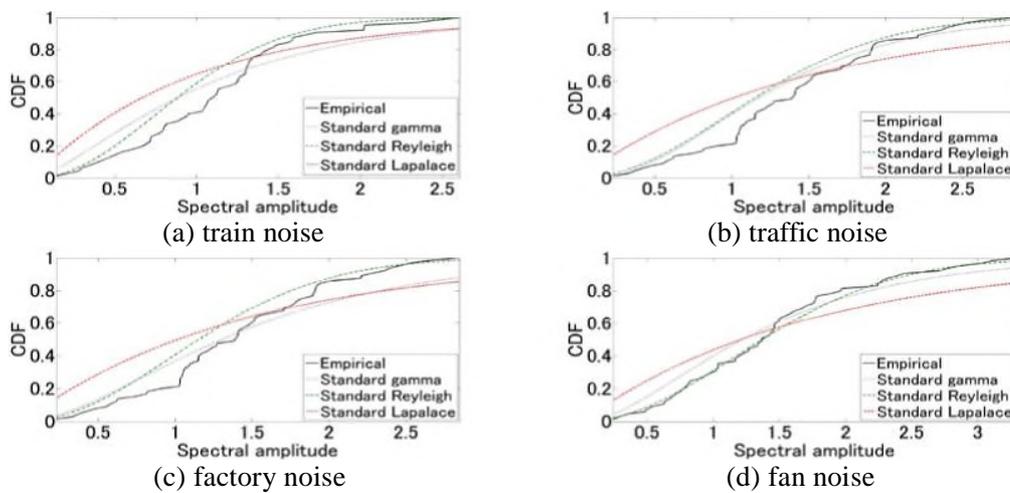


Figure 1. The CDFs of the Spectral Amplitudes of 1kHz at SNR of 5dB

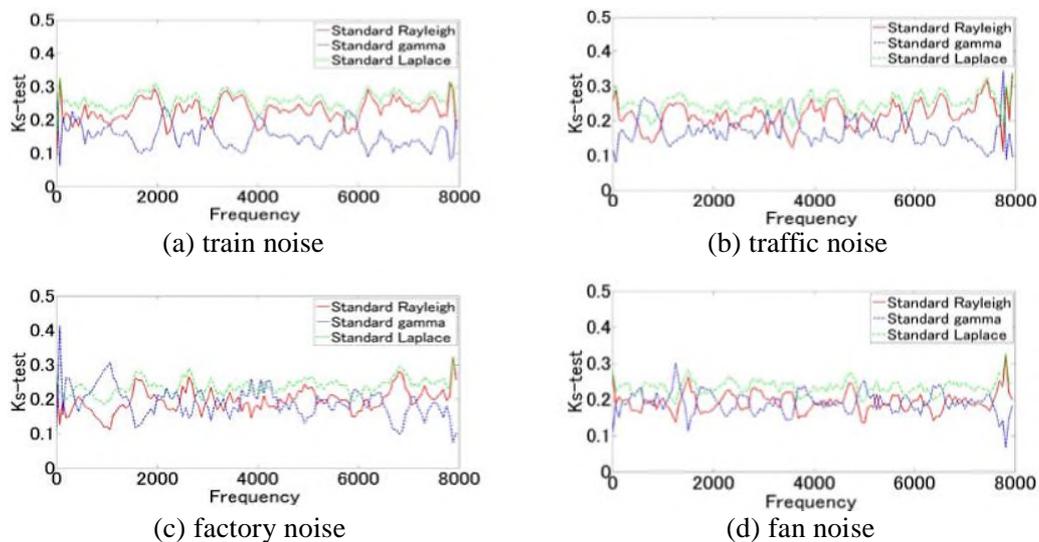


Figure 2. The Kolmogorov–Smirnov (KS) Statistics between Real Noise Spectral CDF and the Standard CDFs

2.2. The Approximate Generalized Gamma PDF

For the standard Rayleigh or gamma distribution model, it is impossible to fit nonwhite noise with frequency bin k . In order to settle this problem, an approximate probability density function is introduced. Similar to the speech DFT component estimation in [6] [7], the generalized Gamma PDF is also suitable for the noise DFT component distribution. The approximating function [6] [7] with shape parameter v and scale parameter p is defined as

$$f_N(n) = \frac{n^{\alpha-1} \exp(-\alpha n)}{\Gamma(\alpha)} \quad (7)$$

Where $\hat{\sigma}_N = N_{mean}(k)$ is the training data of the noise PSD, which makes an accurate function to approximate the real noise distribution. After choosing suitable p and v for the Gamma approximate PDF, the generalized gamma approximate PDF will be close to the standard gamma distribution when $\alpha=1$, and close to the Rayleigh distribution when $\alpha=2$.

2.3. The Parameters of the Normalized Approximate Probability Density Function

In this subsection, we develop a parameter estimation based on the moment matching method, to find out a parameter set (p, v) that best approximates the real noise distribution at each frequency bin.

We use a fourth order moment of the spectral amplitudes of observed noisy speech and it can be denoted by [7]

$$Y^4 = S^4 + N^4 + 4 S^2 N^2 \cos 2(\phi_S - \phi_N) \quad (8)$$

where \oplus denotes expectation. Because the phase difference $(\phi_S - \phi_N)$ is a uniform distribution [6], then the Eq. (8) is given by

$$Y^4 = S^4 + N^4 + 4 S^2 N^2 \quad (9)$$

where $S^2(L, k)$ and $N^2(L, k)$ are the moments powers of the speech and noise.

When $\alpha=2$, according to equation (7),

$$f_N(n) = \frac{2n \exp(-2n)}{\Gamma(2)} \quad (10)$$

and Gaussian assumption of the speech DFT component we can

$$\frac{Y^4}{4 S^2 N^2} = \frac{S^2 + N^2}{2 S N} \quad (11)$$

Then the Eq. (8) becomes

$$8 S^2 N^2 = Y^4 - 4 S^2 N^2 \cos 2(\phi_S - \phi_N)$$

Finally, we can define

$$\frac{Y^4}{4 S^2 N^2} = \frac{S^2 + N^2}{2 S N} \quad (12)$$

where $S^2 = S^2$, $N^2 = N^2$, and the normalization condition implies the constraint

$$\alpha = \gamma = q.$$

(13)

When $\alpha=1$, the Eq. (8) becomes

$$\left(\frac{4}{\Gamma M} \right)^{\frac{1}{2}} + \frac{1}{2} + \frac{1}{2} + 8 \frac{\Gamma M}{\Gamma M} \quad (14)$$

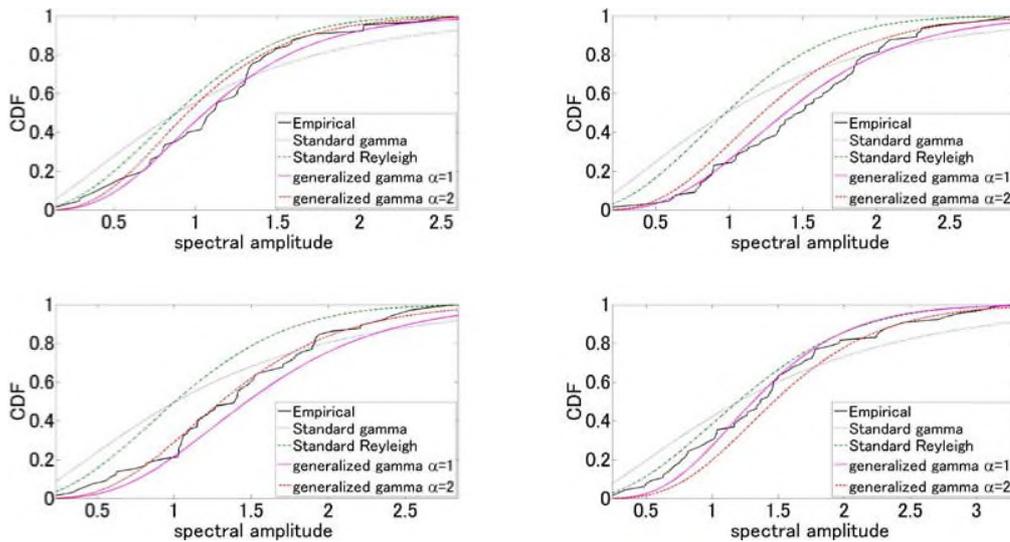
Same as above, we obtain

$$\left\{ = \frac{4}{\Gamma} + \frac{4}{\Gamma} \right\} q$$

(15)

and $\alpha = (\gamma + 1)$.

In the Figure 3, we show the CDFs of noise spectral amplitudes at 1kHz. It is very clear that the generalized gamma PDFs with adaptive parameter are fitting much better



to the real noise CDFs. The KS statistics results are shown in Figure 4. Two generalized gamma models have less KS distance with the real noise spectral CDFs. But from Figure 3 or Figure 4, it is very difficult to find out which generalized gamma model is better.

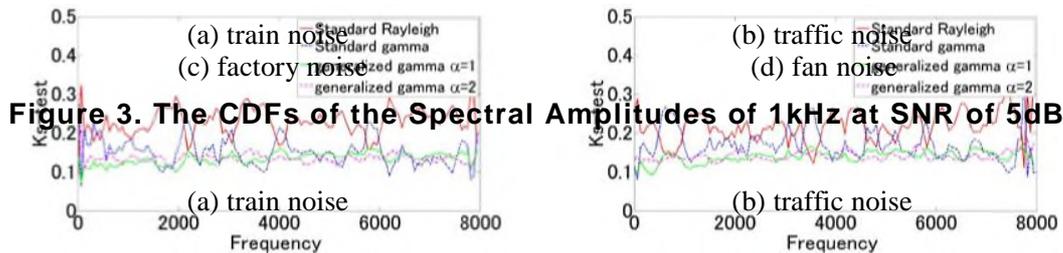


Figure 3. The CDFs of the Spectral Amplitudes of 1kHz at SNR of 5dB

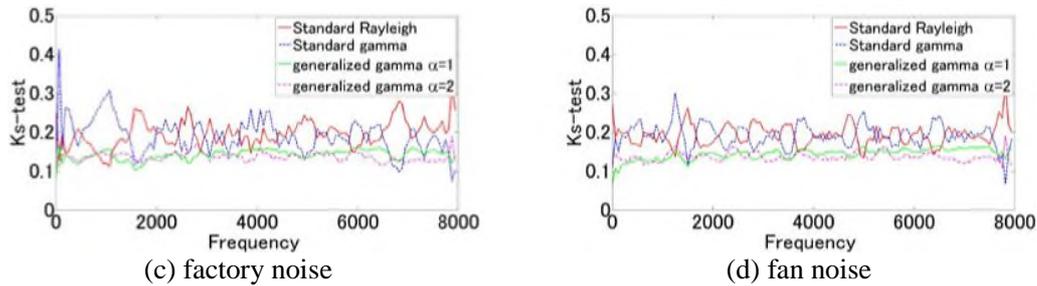


Figure 4. The Kolmogorov–Smirnov (KS) Statistics between Real Noise Spectral CDF and the Generalized Gamma CDFs

2.5. Noise PSD Estimation based on MAP

2.5.1. MAP Estimation when $\alpha=1$

At first, according to the Gaussian assumption, we can define the $p(Y/N)$ as

$$p_s(Y^N) = \frac{1}{2^N} \exp \left\{ -\frac{1}{2} \sum_{s=1}^N \frac{Y_s^2}{\sigma^2} \right\} \quad (16)$$

Instead of differentiating, the maximization can perform better after applying the natural logarithm, because the product of the polynomial and exponential converts into a sum

$$\frac{d \log [p(Y/N)p(N)]}{dN} \propto \sum_{s=1}^N \frac{Y_s^2}{\sigma^2} - \frac{1}{\sigma^2} = 0 \quad (16)$$

After multiplication with N , one reasonable solution $\hat{N} = GY$ to the quadratic equation is found, because the second solution delivers spectral amplitudes $N < 0$ at least for $\gamma > 0$. The second derivative at N is negative, thus a local maximum is guaranteed.

$$G = u + \frac{1}{2} \left[\frac{2v + u}{2} \pm \sqrt{\left(\frac{2v + u}{2} \right)^2 - \frac{1}{4}} \right] \quad (18)$$

2.5.2. MAP Estimation when $\alpha=2$

$$\frac{d \log [p(Y/N)p(N)]}{dN} \propto \sum_{s=1}^N \frac{Y_s^2}{\sigma^2} - \frac{2}{\sigma^2} = 0 \quad (19)$$

Considering the above calculations, when $\alpha=2$ the MAP estimation can be written as

$$\frac{d \log [p(Y/N)p(N)]}{dN} \propto \sum_{s=1}^N \frac{Y_s^2}{\sigma^2} - \frac{2}{\sigma^2} = 0 \quad (19)$$

The weight G will be

$$G = u + \frac{1}{2} \left[\frac{2v + u}{2} \pm \sqrt{\left(\frac{2v + u}{2} \right)^2 - \frac{1}{4}} \right], \quad u = 2 + 4 \left(\frac{1}{\sigma^2} \right) \quad (20)$$

In Figure 5 the weights G are shown, the slope of MAP-gamma $\alpha=1$ shows a peak between 1 and 2 of posteriori SNR $\gamma = \frac{N \sigma^2}{\sigma^2}$. MAP-gamma $\alpha=2$ makes the fastest delay. The MMSE and MLE methods based on Gaussian model show less dynamic range and makes delay slowly, which means a worse resolution.

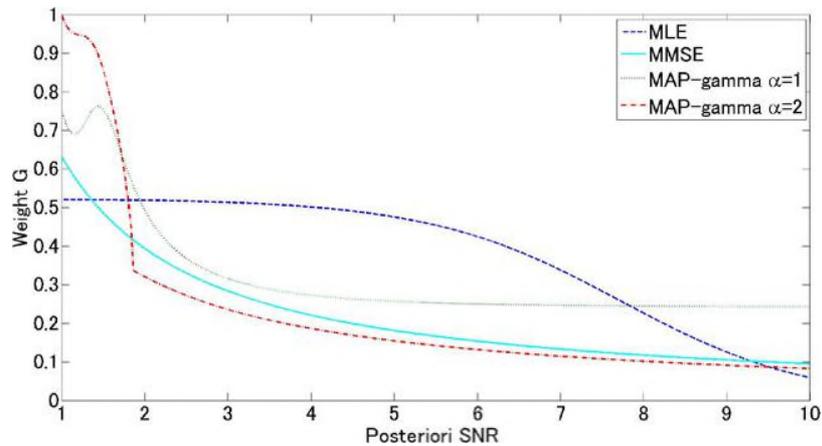


Figure 5. The Weight C of the Proposed and Reference Methods

3. Experiments and Simulated Results

3.1. Modified Wiener Filter

The traditional Wiener filter can be formulated as

$$\hat{S}(\lambda, k) = \frac{P_s(\lambda, k)}{P_s(\lambda, k) + P_d(\lambda, k)} Y(\lambda, k) \quad (21)$$

where $P_s(\lambda, k)$ is the clean speech power spectrum, and $P_d(\lambda, k)$ is the noise power spectrum. Thus, we have

$$\hat{S}(\lambda, k) = H(\lambda, k) Y(\lambda, k), \quad \hat{S}(\lambda, k) = \frac{P_s(\lambda, k)}{P_s(\lambda, k) + P_d(\lambda, k)} S(\lambda, k) \quad (22)$$

where $S(\lambda, k)$ is the estimation of the Fourier transform of the clean speech. Then the modified Wiener Filter [9] can be simplified into

$$H(\lambda, k) = \frac{Y(\lambda, k)^2}{Y(\lambda, k)^2 + \frac{Y(\lambda, k)^2}{N(\lambda, k)^2}} \quad (23)$$

The modified Wiener Filter improves the performance of the noise reduction system at a satisfactory level depending on the accuracy of the noise spectral estimation, and this method produces less musical noise as discussed in [9].

3.2. Experiment Setup

The experiment is performed using the 3 male and 3 female speeches. The data used were taken from the University of Tsukuba Multilingual Speech Corpus (UT-ML). The speeches were degraded by noise sources with input SNRs of 0, 5, 10 and 15dB. The four additive acoustic noises were taken from JEIDA-NOISE database. All signals and noises are sampled at a frequency of 16 kHz. All the frames have a length of $N=256$ with an overlap of 50%, and 84 are windowed using Hanning window.

3.3. Evaluation Measures

In order to measure the objective evaluation of the quality, we employed the Perceptual Estimation of Speech Quality (PESQ ITU-T P.862) [10], log-error distortion and segmental SNR improvement. The symmetric segmental log-error in Figure 6 is defined

$$LogErr = \frac{1}{L} \sum_{k=1}^K \sum_{l=1}^L \left| 10 \oplus \log_{10} \left(\frac{\hat{s}(l, k)}{s(l, k)} \right)^2 \right| \quad (24)$$

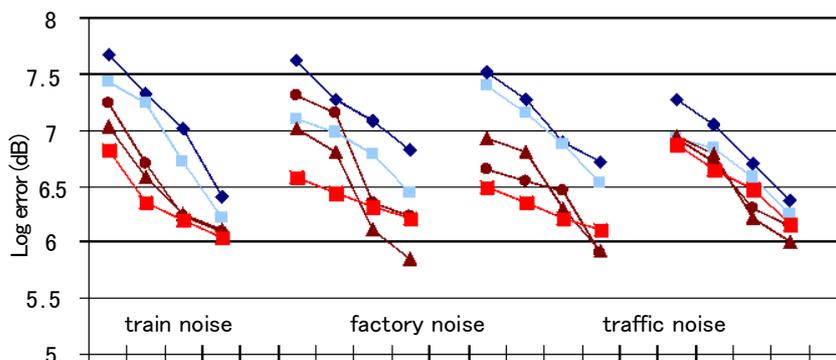
and segmental SNR is defined as

$$SNR_{seg} = \frac{1}{L} \sum_{l=1}^L \left(\frac{\sum_{k=1}^K \hat{s}(l, k)^2}{\sum_{k=1}^K s(l, k)^2} \right) \quad (25)$$

3.4. Performance Evaluation

Figure 6 shows the log error results of the noise PSD estimation method. The results show MAP gamma $\alpha=1$ and $\alpha=2$ have least log errors at high SNRs (10 and 15 dB) and low SNRs (0 and 5 dB) respectively. The MS has less log error than MLE in most of the situations except at the SNR of 0 and 5 dBs in factory noise. At the high SNRs the MS shows less log error than MAP gamma $\alpha=2$. The log errors of MMSE are largest in all noise environments. The log errors of MLE show similar trend and less than MMSE.

In Figure 7 the segmental SNRs of the enhanced speech are given. The MLE shows best results at 0 and 5 dBs and better results at 10 and 15 dBs. MAP gamma $\alpha=1$ and $\alpha=2$ show almost same results in this test. That is worst at low SNRs (0 and 5 dBs) but best at high



SNRs. The results of MS improve slowly and better than MMSE, which has similar trend.

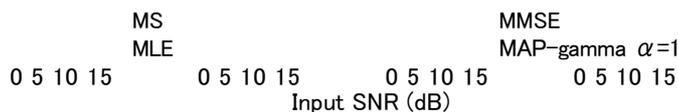


Figure 6. The Logerror of Noise PSD Estimation

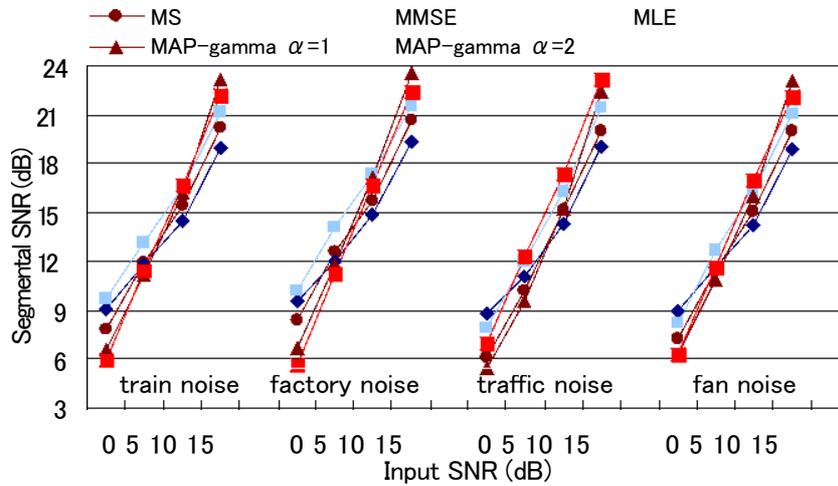


Figure 7. The Segmental SNR of Enhanced Speech

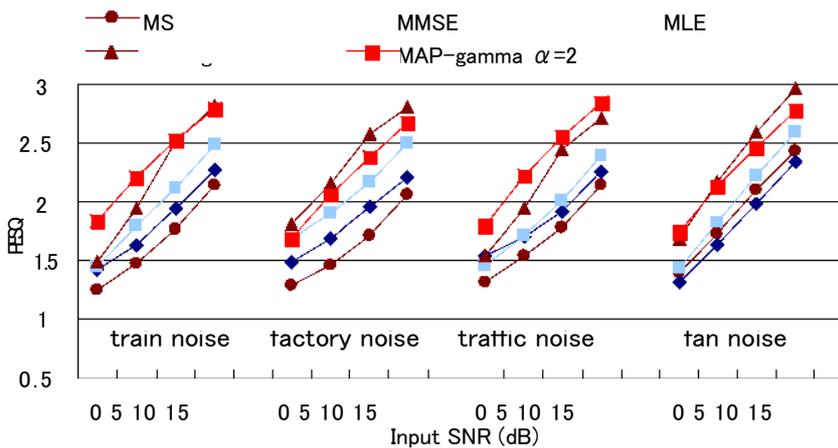


Figure 8. The PESQ of Enhanced Speech

The PESQs of enhanced speech are shown in Figure 8. The MAP gamma $\alpha=2$ shows the best speech quality in train and traffic noises and MAP gamma $\alpha=1$ shows the best performance in factory and fan noises. The MLE is better than MMSE and MS. The MS shows worst results in all environments except in the fan noise environment.

According to the Figure 5, the slope of the MAP gamma $\alpha=2$ has a largest dynamic range when posteriori SNR ranges from 1 to 2, which suggests it has better resolution at low SNR. However its resolution degrades when posteriori SNR is larger than 2 because in this field the dynamic range of this method is short. In contrast, the slope of the MAP gamma $\alpha=1$ makes fast delay during the posteriori SNR between 2 and 5, which suggests better resolution at high SNRs. The MLE makes delay slowly at low posteriori SNRs and makes fastest delay at high posteriori SNRs above 6, which indicates that this method will show better performance at high SNRs. The MMSE method shows an inverse proportionality character, and kept in middle level of the dynamic range in all posteriori fields.

4. Conclusions

In this paper we focused on the noise power spectral density estimation of noisy speech for speech enhancement system. Two MAP estimators for color noise PSD on the basis of the assumption of generalized Gamma distribution are proposed. The parameters of generalized Gamma PDFs are estimated based on moment matching method and above assumption. The results show those estimators have good online tracking capability. Compared with reference methods such as MMSE, MAP and MLE, MAP gamma $\alpha=1$ has best performance at high SNRs, while MAP gamma $\alpha=2$ has best performance at low SNRS.

It is confirmed that the adaptive generalized Gamma models are suitable for describing the real spectral distribution of color noises. Evaluations show that proposed method has better performance compared with Gaussian noise assumption based MMSE and MLE noise PSD estimation algorithms.

References

- [1] Y. Ephraim and D. Malah, "Speech enhancement using a minimum mean-square error short-time spectral amplitude estimator", *IEEE Trans. Acoustic, Speech, Signal Processing*, vol. ASSP-32, no. 6, (1984), pp. 1109-1121.
- [2] R. Martin, "Noise power spectral density estimation based on optimal smoothing and minimum statistics", *IEEE Trans. Speech, Audio Processing*, vol. 9, no. 5, (2001), pp. 504-512.
- [3] N. S. Kim and J.-H. Chang, "Spectral enhancement based on global soft decision", *IEEE Signal Processing Letter*, vol. 7, (2000), pp. 108-110.
- [4] R. C. Hendriks, R. Heusdens and J. Jensen, "MMSE based noise PSD tracking with low complexity", *IEEE International Conference on Acoustics, Speech, and Signal Processing (ICASSP 2010)*. IEEE, (2010), pp. 4266-4269.
- [5] R. Yu, "A low-complexity noise estimation algorithm based on smoothing of noise power estimation and estimation bias correction", *IEEE International Conference on Acoustics, Speech, and Signal Processing (ICASSP 2009)*, (2009), pp. 4421-4424.
- [6] T. Lotter and P. Vary, "Speech Enhancement by MAP Spectral Amplitude Estimation using a Super-Gaussian Speech Modeling", *EURASIP Journal on Applied Signal Processing*, (2005), pp. 1110-1126.
- [7] T. Huy Dat and K. Takeda "Generalized gamma modeling of speech and its online estimation for speech enhancement", *IEEE International Conference on Acoustics, Speech, and Signal Processing (ICASSP 2005)*, vol. 4, (2005), pp. iv/181-iv/184.
- [8] M. Delcroix, K. Kinoshita, "Noise Power Spectral Density Tracking: A Maximum Likelihood Perspective", *IEEE Signal Processing Letters*, vol. 19, no. 8, (2012), pp. 495-498.
- [9] L. M. Arslan, "Modified Wiener filtering", *Signal Processing*, vol. 86, no. 2, (2006), pp. 267-272.
- [10] A. Rix, J. Beerends, M. Hollier and A. Hekstra, "Perceptual Evaluation of Speech Quality (PESQ)-a new method for speech quality assessment of telephone networks and codecs", *IEEE International Conference on Acoustics, Speech, and Signal Processing (ICASSP 2001)*, vol. 2, (2001), pp. 749-752.

Authors



Xin Dang received the M.E. degrees in Electrical Engineering from Shizuoka University in 2010. He is currently in the doctor's course at Shizuoka University in Japan.



Takayoshi Nakai received the B.E., M.E., and Dr.E. degrees from Shizuoka University, Hamamatsu, Japan, in 1974, 1976, and 1994, respectively. From 1976 to 1996 he was a Research Associate, from 1996 to 2004 he was an Associate Professor, and since 2004 he has been a Professor at the Faculty of Eng., Shizuoka University.



Md. Iqbal Aziz Khan received his B.Sc. (Hons.) and M.Sc. Degrees from the Department of Computer Science and Engineering, University of Rajshahi, Bangladesh. Currently he has been working as a PhD student at Shizuoka University, Japan.