JAIST Repository

https://dspace.jaist.ac.jp/

Title	Psychoacoustically-motivated adaptive -order generalized spectral subtraction based on data- driven optimization
Author(s)	Li, Junfeng; Jiang, Hui; Akagi, Masato
Citation	Proceedings of INTERSPEECH 2008: 171-174
Issue Date	2008-09-23
Туре	Conference Paper
Text version	publisher
URL	http://hdl.handle.net/10119/9984
Rights	Copyright (C) 2008 International Speech Communication Association. Junfeng Li, Hui Jiang, Masato Akagi, Proceedings of INTERSPEECH 2008, pp.171–174.
Description	



Japan Advanced Institute of Science and Technology





Psychoacoustically-motivated Adaptive β -order Generalized Spectral Subtraction Based on Data-driven Optimization

Junfeng Li¹, Hui Jiang², Masato Akagi¹

¹ School of Information Science, Japan Advanced Institute of Science and Technology, Japan ² Department of Computer Science and Engineering, York University, Canada

Email: junfeng@jaist.ac.jp, hj@cse.yorku.ca, akagi@jaist.ac.jp

Abstract

To mitigate the performance limitations caused by the constant spectral order β in the traditional spectral subtraction methods, we previously presented an adaptive β -order generalized spectral subtraction (GSS) in which the spectral order β is updated in a heuristic way [10]. In this paper, we propose a psychoacoustically-motivated adaptive β -order GSS, by considering that different frequency bands contribute different amounts to speech intelligibility (i.e., the bandimportance function). Specifically, in this proposed adaptive β -order GSS, the tendency of spectral order β to change with the input local signal-to-noise ratio (SNR) is quantitatively approximated by a sigmoid function, which is derived through a data-driven optimization procedure by minimizing the intelligibility-weighted distance between the desired speech spectrum and its estimate. The inherent parameters of the sigmoid function are further optimized with the data-driven optimization procedure. Experimental results indicate that the proposed psychoacoustically-motivated adaptive β -order GSS yields great improvements over the traditional spectral subtraction methods with the intelligibility-weighted measures.

Index Terms: Psychoacoustically-motivated adaptive β -order GSS, Data-driven optimization, Importance function, Speech intelligibility, Sigmoid function.

1. Introduction

Background noises severely degrade performance of state-ofthe-art speech applications, for example, in the reduced intelligibility of perceived speech for cochlear implant users. As a pre-processor, noise reduction has been shown to be effective in improving the quality and/or intelligibility of noisy signal, and feeding the enhanced signal to implant users results in significant benefits in sentence/word recognition [1].

Research studies on speech intelligibility show that various factors affect speech intelligibility to different degrees [2]. These factors include the audibility in each subband derived from the corresponding subband SNR and the listener's hearing threshold, and the band-importance functions that indicate to which degree each frequency band contributes to intelligibility [3, 4]. Speech intelligibility was traditionally computed from the long-term speech and noise spectra in stationary noise [2, 3]. This calculation approach was recently extended for non-stationary noise conditions by first calculating the speech intelligibility in each short-term frame followed by averaging these intelligibility values across frames, yielding the overall intelligibility for that particular condition [4].

To improve the speech intelligibility of the received signal in noisy conditions, a number of noise reduction algo-

rithms have been proposed in the past several decades [5]. As a well-known noise reduction method, spectral subtraction (SS) has been used widely due to its simplicity in implementation and improved to overcome its shortcomings in different ways [5, 6, 7, 8]. Boll applied several secondary procedures to the processed signal after SS to further attenuate the residual noise [5]. To enhance noise reduction and mitigate "musical" noise, Berouti et al. introduced two additional parameters, an oversubtraction factor that controls the amount of noise to be subtracted, and a spectral flooring factor that mitigates "musical" noise [6]. Schless et al. suggested to set both oversubtraction factor and spectral flooring factor based on the signal-tonoise ratio (SNR) [7]. Kamath et al. proposed to empirically set different oversubtraction factors in different subbands [8]. Sim et al. derived a short-time spectral amplitude estimator of the speech signal based on a parametric formulation of generalized spectral subtraction (GSS) [9]. In all SS methods mentioned above, however, the spectral order β is always fixed to some constants which result in performance limitation to a certain degree. More recently, therefore, Li et al. suggested an adaptive β -order GSS in which the spectral order β is updated in each subband according to the input local SNRs in a heuristic way [10].

In this paper, we propose a psychoacoustically-motivated adaptive β -order GSS based on the research results on speech intelligibility and a data-driven optimization procedure. Considering the different contributions of different frequency bands to speech intelligibility which is quantified by the band-importance function [2, 3], we propose to quantitatively determine the tendency of spectral order β to change with the input local SNR in MMSE sense, and introduce a data-driven optimization procedure that quantifies this change tendency as a sigmoid function and optimizes the inherent parameters. Experimental results in various noise conditions show that the proposed method outperforms the traditional SS algorithms in terms of the intelligibility-weighted measures.

2. β -order generalized spectral subtraction

2.1. Signal model

Suppose the observed noisy signal to be the sum of the clean speech signal and the uncorrelated additive noise signal. Applying the *short-time Fourier transform* (STFT), the observed signal in the time-frequency domain is represented as

$$X(k,\ell) = S(k,\ell) + N(k,\ell), \tag{1}$$

where k and ℓ are the frequency bin index and the time frame index, respectively; $X(k, \ell)$, $S(k, \ell)$ and $N(k, \ell)$ are the STFTs of the noisy signal, the clean signal and the noise signal.

2.2. β -order generalized spectral subtraction

The β -order generalized spectral subtraction is defined as [9]

$$\left|\hat{S}_{\beta}\left(k,\ell\right)\right|^{\beta} = a_{\beta}\left(k,\ell\right) \left|X(k,\ell)\right|^{\beta} - b_{\beta}(k,\ell) E\left[\left|N(k,\ell)\right|^{\beta}\right], \quad (2)$$

where β denotes the spectral order, $\hat{S}_{\beta}(k, \ell)$ is the spectral estimate of the speech enhanced by the β -order GSS, and $a_{\beta}(k, \ell)$ and $b_{\beta}(k, \ell)$ are two parameters. Note that the speech spectrum estimate $\hat{S}_{\beta}(k, \ell)$ is dependent not only on time and frequency, but also on the spectral order β .

The parameters $a_{\beta}(k, \ell)$ and $b_{\beta}(k, \ell)$ in the β -order GSS are determined and optimized by minimizing the MSE of the β -order speech spectrum amplitude $|S_{\beta}(k, \ell)|^{\beta}$ and its estimate $|\hat{S}_{\beta}(k, \ell)|^{\beta}$. Under the complex Gaussian assumption and substituting the two newly-derived optimal coefficients into Eq. (2), the gain function of the β -order GSS is derived as [9]

$$\hat{G}_{\beta}(k,\ell) = \left\{ \frac{\left[\xi_{\beta}(k,\ell)\right]^{\beta}}{1 + \left[\xi_{\beta}(k,\ell)\right]^{\beta}} \right\}^{\frac{1}{\beta}} \left\{ 1 - \left(1 - \left[\xi_{\beta}(k,\ell)\right]^{\frac{-\beta}{2}}\right) \\ \Gamma\left(\frac{\beta}{2} + 1\right) \left(\frac{1}{\gamma(k,\ell)}\right)^{\frac{\beta}{2}} \right\}^{\frac{1}{\beta}},$$
(3)

where $\xi_{\beta}(k, \ell)$ and $\gamma(k, \ell)$ are the *a priori* SNR and the *a posteriori* SNR as defined in [11], and $\Gamma(\cdot)$ denotes the Gamma function. The estimate of $\xi_{\beta}(k, \ell)$ is updated in a decision-directed scheme, greatly decreasing the residual "musical" noise [11].

3. Psychoacoustically-motivated adaptive β -order GSS using data-driven optimization

Considering the research findings on speech intelligibility, in this section, we propose a psychoacoustically-motivated adaptive β -order GSS in which the dependence of β on the local SNR is derived and optimized through a data-driven optimization procedure.

3.1. Subband-processing-based determination of spectral order β

The value of spectral order β has been shown to be largely dependent on the SNR in the current condition [10]. Furthermore, the SNRs vary greatly with time due to the time-varying characteristics of speech and noise signals, and also significantly vary in different subbands because of the colorness of noise signals and the non-uniform spectral energy distribution of speech signals. As a result, speech signal corrupted by real-world noises is characterized by different local SNRs in different partitions in the time-frequency domain. Further taking the strong correlations of spectral components between adjacent frequency bins into account, the appropriate value of β should be determined depending not only on knowledge of the current frequency bin under consideration, but also on knowledge of the neighboring bins. As a result, the appropriate value of spectral order β must be adaptively determined according to the local SNRs calculated in subbands instead of the instantaneous SNR in each individual frequency bin in the time-frequency domain.

3.2. Psychoacoustically-motivated derivation of spectral order β

Improving speech intelligibility motivates us to derive the appropriate value for spectral order β by integrating the psychoacoustical research findings on speech intelligibility. There are

a number of aspects which play an important role in speech intelligibility enhancement. For example, it is understood that the short-term spectrum is of primary importance in the perception of speech, and different frequency bands contribute different amounts to speech intelligibility [2, 3]. Speech intelligibility index (SII) in ANSI S.35-1997 standard quantifies the degree of speech intelligibility in the presence of background noise [3]. Audibility is measured as the ratio of time-averaged speech power and time-averaged noise power in a set of frequency bands, followed by weighting with the so-called bandimportance function. The band-importance function indicates to which degree each frequency band contributes to intelligibility. Finally, the SII is determined by accumulation of the audibility across the different frequency bands. As a result, the band-importance functions should be integrated when determining the appropriate value of spectral order β , as shown in the following section.

3.3. Data-driven optimization of spectral order β

In our previous study, the change tendency of the spectral order β along with the input SNR was qualitatively analyzed and intuitively described by a sigmoid function, where no theoretical/experimental evidences were provided [10]. With the consideration of the non-uniform effect of noise on speech in Section 3.1 and the psychoacoustical findings on speech intelligibility in Section 3.2, in this section, we aim to quantify the dependency of the spectral order β on the input local SNR using a data-driven optimization approach.

In our data-driven optimization procedure, ten speech sentences were randomly selected from the NTT database [12], and two noise signals ("car" and "babble") were taken from the NOISEX-92 database [13]. The speech and noise signals were first downsampled to 8kHz and then mixed with the global SNR ranging from -40 to 40 dB. We assume that the noise spectrum is a known prior in this optimization procedure. For a given value of β , the gain function of the β -order GSS is calculated using Eq. (3), and then used to enhance the target speech signal. Furthermore, based on the discussions in Sections 3.1 and 3.2, we propose to optimize the spectral order β by minimizing the overall intelligibility-weighted distance between the spectral amplitude $|S(k, \ell)|$ of the clean signal and that of its estimate $|\hat{S}_{\beta}(k, \ell)|$ summed across all subbands, that is,

$$\beta^{\text{opt}} = \underset{0.1 \le \beta \le 3.0}{\arg\min} \left(\sum_{m=1}^{M} \sum_{k=\omega_m}^{\omega_{m+1}} I_m \left[\left| S(k,\ell) \right| - \left| \hat{S}_{\beta}(k,\ell) \right| \right]^2 \right), \quad (4)$$

where I_m is the band-importance function in the *m*-th subband [3], *M* is the number of subbands, ω_m denotes the boundary frequency of the *m*-th subband, and the range of β is empirically confined to [0.1, 3.0]. Note that the importance functions are normally given as constant values [3] and each term in the overall intelligibility-weighted spectral amplitude distance in Eq. (4) are positive, therefore, the minimization problem performed in each independent subband. That is, the minimization of the overall intelligibility-weighted spectral amplitude distance can be achieved by the value that yields the minimal intelligibility-weighted spectral amplitude distance of the clean signal and its estimate in each individual subband, given by

$$\beta_m^{\text{opt}} = \underset{0.1 \le \beta \le 3.0}{\operatorname{arg\,min}} \left(\sum_{k=\omega_m}^{\omega_{m+1}} \left[\left| S(k,\ell) \right| - \left| \hat{S}_\beta(k,\ell) \right| \right]^2 \right).$$
(5)



Figure 1: Scatter plots of the optimized β value with respect to the input local SNR, the mean of the scattered data (solid line) and the fitted sigmoid function (dashed line) with the parameters A = 0.1, B = 2.0, D = 7.0, in the car noise condition (upper panel) and in the babble noise condition (lower panel).

Though only ten speech sentences are used in this optimization procedure, we should note that the optimization is performed in the following scenarios: in each frames (each speech sentence is divided into $220 \sim 350$ overlapping frames by windowing before Fourier transform.), in each critical subband (e.g., 18 subbands) and at the different global SNR conditions (that is, -40 to 40 dB in 10 dB increments). It is therefore believed that the optimization of spectral order β is sufficient in the statistical sense, and the parameters obtained from this optimization procedure might also be applicable in other different conditions.

Fig. 1 shows the scatter plots of the optimized β value against the local SNR (defined in Eq. (6)), the mean curve as well as the fitted sigmoid function in the "car" and "babble" noise conditions. The results shown in Fig. 1 indicate that: (1) The optimal β value increases (decreases) as the local input SNR increases (decreases), which proves our previous qualitative analysis results [10]; (2) Most importantly, the change tendency of the appropriate value of β with the change of the local SNR can be approximated by a sigmoid function (defined in Eq. (7)); (3) Moreover, the inherent parameters are also determined through this data-driven optimization.

3.4. Adaptive scheme for an appropriate value of spectral order β

As stated in Sections 3.1 and 3.2, the appropriate value of spectral order β vary with the change of the local SNR in each subband. Further considering the mechanism of human perception, the whole spectrum is first divided into subbands according to

the critical-band scale [14]. Then, the local SNR $\rho(m,\ell)$ in the m-th critical band and the $\ell\text{-th}$ frame is calculated as

$$\rho(m,\ell) = 10 \log_{10} \left(\frac{\sum_{k=\omega_m}^{\omega_{m+1}} \left| |X(k,\ell)| - |\hat{N}(k,\ell)| \right|^2}{\sum_{k=\omega_m}^{\omega_{m+1}} \left| \hat{N}(k,\ell) \right|^2} \right), \quad (6)$$

where ω_m is specified as the boundary frequency of the *m*-th critical band.

Based on the data-driven optimization results in Section 3.3, we propose to determine the optimized value of the spectral order $\tilde{\beta}(m, \ell)$ according to the local SNR $\rho(m, \ell)$ in each critical band, frame by frame, by the use of the sigmoid function, given by

$$\tilde{\beta}(m,\ell) = \frac{B}{1 + e^{-A\left[\rho(m,\ell) - D\right]}},\tag{7}$$

where the parameter A controls the changing speed of the value of $\tilde{\beta}(m, \ell)$ with respect to the local SNR $\rho(m, \ell)$, B determines the range of the value of $\tilde{\beta}$, and D denotes the shift along the SNR axis. In order to avoid severe speech distortion caused by an extremely low β value, we further confine the value of $\tilde{\beta}(k, \ell)$ with a minimum value β_{min} . As a result, the optimized value of the spectral order $\hat{\beta}(m, \ell)$ is finally determined as

$$\hat{\beta}(m,\ell) = \max\left[\tilde{\beta}(m,\ell), \beta_{min}\right].$$
(8)

4. Experiments and results

4.1. Experimental configuration

Performance of the proposed psychoacoustically-motivated adaptive β -order GSS was assessed by the following experiments. We randomly selected 40 clean continuous speech sentences produced by two females and two males from NTT speech database [12]. Two types of noise source, "car" and "babble", were chosen from NOISEX-92 database [13]. The clean speech and noise signals were first downsampled to 8kHz. We generated noisy speech signals artificially by adding various noise signals to the clean signals at different SNRs ranging from 0 to 15 dB in 5 dB increments. Note that the car noise was a stationary signal, whereas the babble noise was a nonstationary signal. Based on the optimization results shown in Fig. 1, the parameters in our experiments were set as follows: A = 0.1, B = 2.0, D = 7.0 and $\beta_{min} = 0.1$. Note that the spectral order β in the current experiments ranges from 0.1 to 2.0. Performance of the proposed adaptive β -order GSS was further compared with that of the traditional SS methods (power SS, amplitude SS and SS when $\beta = 0.1$) that were implemented by setting the spectral order β to 2.0, 1.0 and 0.1 in Eq. (3) derived in [9], respectively.

4.2. Experimental results

Performance of the studied algorithms was objectively examined in terms of *intelligibility-weighted SNR* (SNR_{int}) and *intelligibility-weighted log-spectral distance* (LSD_{int}), as defined in [15]. The experimental results of SNR_{int} and LSD_{int} , averaged across all sentences in two noise conditions, are shown in Tables 1 and 2, respectively.

Table 1 illustrates that the proposed psychoacousticallymotivated adaptive β -order GSS consistently provides the highest SNR_{int} improvements, compared to the traditional SS methods, for all conditions at all SNRs. While the SS method with

Table 1: Intelligibility-weighted SNR [dB] of the noisy signal, the power SS ($\beta = 2.0$) output, the amplitude SS ($\beta = 1.0$) output, the SS output when $\beta = 0.1$, and the proposed adaptive β -order GSS output.

Global SNR	0	5	10	15
Condition	car noise			
Noisy	-11.51	-6.26	-1.15	3.85
$\beta = 2.0$ (Power SS)	-6.02	-2.35	1.06	4.02
$\beta = 1.0$ (Amplitude SS)	-3.56	-0.54	2.50	5.32
$\beta = 0.1$	-5.33	-5.54	-5.35	-5.33
Adaptive- β	-2.12	0.89	3.76	6.04
Condition	babble noise			
Noisy	-13.96	-8.61	-3.41	1.65
$\beta = 2.0$ (Power SS)	-9.54	-5.31	-1.42	1.95
$\beta = 1.0$ (Amplitude SS)	-6.56	-3.05	0.28	3.32
$\beta = 0.1$	-5.22	-5.28	-5.31	-5.32
Adaptive- <i>β</i>	-4.38	-1.33	1.58	4.22

Table 2: Intelligibility-weighted LSD [dB] of the noisy signal, the power SS ($\beta = 2.0$) output, the amplitude SS ($\beta = 1.0$) output, the SS output when $\beta = 0.1$, and the proposed adaptive β -order GSS output.

Global SNR	0	5	10	15
Condition	car noise			
Noisy	18.93	15.96	13.43	11.30
$\beta = 2.0$ (Power SS)	15.29	13.22	11.46	10.00
$\beta = 1.0$ (Amplitude SS)	13.97	12.15	10.57	9.29
$\beta = 0.1$	52.78	54.33	55.24	55.69
Adaptive- β	13.69	11.63	10.05	8.97
Condition	babble noise			
Noisy	20.47	17.25	14.51	12.21
$\beta = 2.0$ (Power SS)	17.41	14.89	12.76	11.04
$\beta = 1.0$ (Amplitude SS)	15.75	13.62	11.79	10.27
$\beta = 0.1$	50.37	52.43	53.84	54.69
Adaptive- β	14.78	12.82	11.13	9.81

 $\beta = 0.1$ yields greatly reduced SNR_{int} results because it introduces severe speech distortion due to the too small value of β (i.e., 0.1). The highest SNR_{int} by our proposed algorithm indicates high noise reduction ability corresponding to high speech intelligibility in some sense [15]. This might be attributed to the use of low gains in speech-absence periods due to the low values of the spectral order β .

Concerning the results of LSD_{int} shown in Table 2, we can observe that all tested algorithms decrease the LSD_{int} in all conditions, except for the SS algorithm with $\beta = 0.1$ that markedly increases LSD_{int} (i.e., high speech distortion and low intelligibility). The enhanced speech signal processed by the proposed algorithm involves the lowest speech distortion, which corresponds to high speech intelligibility to a certain degree. This achievement can specifically be attributed to the use of high gains in speech-presence periods, due to the high values of the spectral order β .

As a result, the proposed psychoacoustically-motivated adaptive β -order GSS outperforms the traditional SS algorithms in speech intelligibility enhancement, due to the time-varying frequency-dependent spectral order β that is derived and optimized through a data-driven optimization procedure based on the consideration of psychoacoustic research findings on speech intelligibility (e.g., the band-importance function).

5. Conclusion

In this paper, we proposed a psychoacoustically-motivated β order generalized spectral subtraction for speech intelligibility enhancement, in which the dependence of the spectral order β on the local SNRs is derived and optimized through a datadriven optimization procedure based on consideration of the different contributions of different frequency bands to speech intelligibility. Experimental results illustrate that the proposed algorithm outperforms the traditional SS algorithms in various noise conditions in terms of intelligibility-weighted measures. As a future work, the effectiveness of our proposed algorithm in improving speech intelligibility for cochlear implant users (using state-of-the-art cochlear implants) should be assessed in daily-life noise conditions.

6. Acknowledgement

This study was supported by a Grant-in-Aid for Young Scientists (B) (No. 19700156) from the Ministry of Education, Science, Sports and Culture of Japan.

7. References

- P.C. Loizou, "Introduction to cochlear implants," *IEEE Engineering in Medicine and Biology Magazine*, vol. 18, pp. 32-42, 1999.
- [2] C.V. Pavlovic, "Derivation of primary parameters and procedures for use in speech intelligibility prediction," J. Acoust. Soc. Am., vol. 82, no. 2, pp. 413-422, 1987.
- [3] ANSI S3.5-1997, "American National Standard Methods for Calculation of the Speech Intelligibility Index," 1997.
- [4] K.S. Rhebergen and N.J. Versfeld, "A speech intelligility indexbased approach to predict the speech reception threshold for sentences in fluctuating noise for normal hearing listeners," J. Acoust. Soc. Am. vol. 117, no. 4, pp. 2181-2191, 2005.
- [5] S. F. Boll, "Suppression of acoustic noise in speech using spectral subraction," *IEEE Tans. ASSP*, vol. ASSP-27, no. 2, pp. 113-120, 1979.
- [6] M. Berouti, R. Schwartz, and J. Makhoul, "Enhancement of speech corrupted by acoustic noise," In *Proc. ICASSP*, pp. 208-211, 1979.
- [7] V. Schless and F. Class, "SNR-dependent flooring and noise overestimation for joint application of spectral subtraction and model combination," In *Proc. ICSLP*, pp. 721-725, 1998.
- [8] S.D. Kamath and P.C. Loizou, "A multi-band spectral subtraction method for enhancing speech corrupted by colored noise," In *Proc. ICASSP*, pp. 4164-4167, 2002.
- [9] B.L. Sim, et al. "A parametric formulation of the generalized spectral subtraction," *IEEE Trans. SAP*, vol. 6, no. 4, pp. 328-337, 1998.
- [10] J. Li, et al., "Noise reduction based on adaptive β-order generalized spectral subtraction for speech enhancement," In Proc. Interspeech2007, pp. 802-805, 2007.
- [11] Y. Ephraim and D. Malah, "Speech enhancement using a minimum mean-square error short-time spectral amplitude estimator," *IEEE Trans. ASSP*, vol. 32, no. 6, pp. 1109-1121, 1984.
- [12] http://www.ntt-at.com/products_e/speech2002/.
- [13] A. Varga and H.J.M. Steeneken, "Assessment for automatic speech recognition: II. NOISEX-92: A database and an experiment to study the effect of additive noise on speech recognition systems," *Speech Communication*, vol. 12, pp. 247-251, 1993.
- [14] E. Zwicker and E. Terhardt, "Analytical expressions for critical band rate and critical bandwidth as a funtion of frequency," J. Acoust. Sco. Am., vol. 68, pp. 1523-1525, 1980.
- [15] J.E. Greenberg, P.M. Peterson and P.M. Zurek, "Intelligibilityweighted measures of speech-to-interference ratio and speech system performance," *J. Acoust. Soc. Am.*, vol. 94, no. 5, pp. 3009-3010, 1993.