

Speech enhancement based on leakage constraints DF-GSC

Zou Cairong^{1,2} Chen Guoming¹ Zhao Li¹

(¹ School of Information Science and Engineering, Southeast University, Nanjing 210096, China)

(² Foshan University, Foshan 528000, China)

Abstract: In order to improve the performance of general sidelobe canceller (GSC) based speech enhancement, a leakage constraints decision feedback generalized sidelobe canceller(LCDFF-GSC) algorithm is proposed. The method adopts DF-GSC against signal mismatch, and introduces a leakage factor in the cost function to deal with the speech leakage problem which is caused by the part of the speech signal in the noise reference signal. Simulation results show that although the signal-to-noise ratio (SNR) of the speech signal through LCDF-GSC is slightly less than that of DF-GSC, the IS measurements show that the distortion of the former is less than that of the latter. MOS (mean opinion score) scores also indicate that the LCDF-GSC algorithm is better than DF-GSC and the Wiener filter algorithm.

Key words: speech enhancement; general sidelobe canceller (GSC); speech leakage

In speech communication applications, such as mobile phones, hands-free telephones and hearing aids, speech signals are often corrupted by acoustic background noise. Many speech enhancement methods have been introduced to solve these problems^[1-2]. Since a multi-microphone system exploits spatial information in addition to temporal and spectral information of the desired signal and noise signal, it achieves much better performance than traditional single microphone-based speech enhancement algorithms such as Wiener filtering^[1], subspace-based enhancement^[2] etc.. Frost^[3] provided an algorithm which deals with the problem of a broadband signal received by an array. The algorithm is capable of satisfying some desired signal in the look direction by using constrained minimization of the total output power. Griffiths and Jim^[4] reconsidered Frost's algorithm and introduced the generalized sidelobe canceller (GSC) solution which is a widely used noise reduction algorithm for a multi-microphone system by adjusting the weights of a sensor array with adaptive filters.

The GSC algorithm achieves better performance assuming that the desired speaker location, the microphone characteristics and positions are known beforehand. However, in reality, these assumptions are often violated, resulting in speech leakage into the noise references which causes speech distortion^[5]. Although

some techniques^[6-7] were introduced to reduce the amount of speech leakage, it can never be completely avoided. We consider speech leakage with leakage constraints in the GSC cost function and adopt the DF-GSC method against signal mismatch.

1 GSC-Based Optimal Filter

A GSC-based optimal filter is illustrated in Fig. 1. It consists of three parts. Part A is a fixed beamformer filter which is responsible for providing a speech reference signal. Part B is a blocking matrix which creates a noise reference signal by blocking the direction of the speech source, and part C is an adaptive filter which uses the noise reference signal as an input signal and the speech reference signal as a desired signal. The GSC attempts to recover the speech signal by constraining the array response to unity in the direction of the speech source and minimizes the total energy which is coming from all other directions.

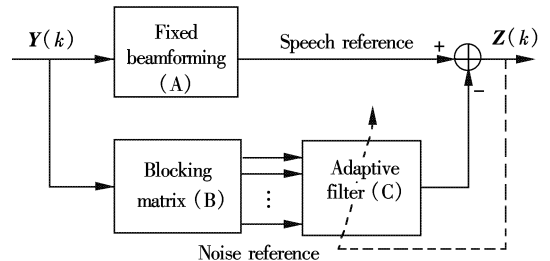


Fig. 1 Basic structure of GSC

Consider a uniform linear array(ULA) of N microphones, where each microphone signal $Y_n(k)$, $n = 0, 1, \dots, N - 1$, at tap time k , consists of a filtered version of the clean speech signal and additive noise:

$$Y_n(k) = x_n(k) + v_n(k) \quad (1)$$

Received 2006-11-28.

Foundation items: The National Natural Science Foundation of China (No. 60472058), the Ph. D. Programs Foundation of Ministry of Education of China(No. 20050286001), Program for New Century Excellent Talents in University(No. NCET-04-0483).

Biography: Zou Cairong (1963—), male, doctor, professor, zoucairong@seu.edu.cn.

where $x_n(k)$ and $v_n(k)$ are the speech component and the noise component received at the n -th microphone, respectively. The additive noise can be colored and is assumed to be uncorrelated with the speech signal.

In Fig. 2, a detailed fixed beamformer is provided. The speech signal from far field impinges on the array from a known DOA of θ_0 along with $M-1$ uncorrelated noise from unknown DOAs $\{\theta_1, \theta_2, \dots, \theta_{M-1}\}$, so the received signal can be written as

$$Y(k) = a(\theta_0)x_0(k) + \sum_{m=1}^{M-1} a(\theta_m)x_m(k) + v(k) = x(k) + v(k) \quad (2)$$

where $a(\theta_0) = \{1, \exp(i\tau_{\theta_0}), \exp(i2\tau_{\theta_0}), \dots, \exp(i(N-1)\tau_{\theta_0})\}^T$, with $\tau_{\theta_0} = \frac{2\pi d}{\lambda} \sin\theta_0$, and λ is the signal wavelength. Similarly, $a(\theta_m) = \{1, \exp(i\tau_{\theta_m}), \exp(i2\tau_{\theta_m}), \dots, \exp(i(N-1)\tau_{\theta_m})\}^T$, with $\tau_{\theta_m} = \frac{2\pi d}{\lambda} \sin\theta_m$. The output signal is $z(k) = \sum_{n=0}^{N-1} w_n^T(k)y_n(k)$.

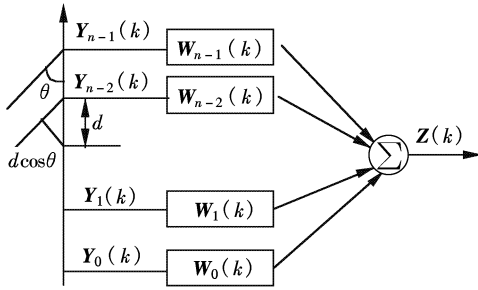


Fig. 2 Structure of fixed beamformer

The matrix B with dimension $N \times (N-1)$, blocking the direction of the speech source should satisfy the condition given as

$$B^H a(\theta_0) = 0 \quad (3)$$

Part C is depicted in Fig. 3. The output signal $z(k)$ can be written as $z(k) = \sum_{n=0}^{N-1} w_n^T(k)y_n(k) = W^T(k)y(k)$.

In Fig. 3, $x(k)$ is the desired response vector and $e(k) = x(k) - z(k)$ is the estimation error vector. The mean square error (MSE) cost function leads to the well-known multidimensional Wiener filter $W(k) = R_{yy}^{-1}(k) R_{yx}(k)$.

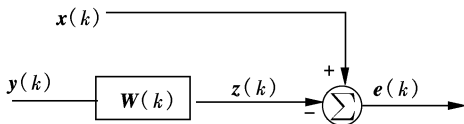


Fig. 3 Structure of MWF

2 Leakage Constraints Decision Feedback GSC (LCDF-GSC)

2.1 GSC solution

An effective approach to determining the weight

$W(k)$ in part C is based on the LCMV (linearly constrained minimum variance) criterion, which can be expressed as

$$\min_W W^H R_y W \quad \text{subject to } C^H W = f \quad (4)$$

where $R_y(k) = E\{y(k)y^T(k)\}$ is the input correlation matrix; C is an $N \times p$ constraint matrix and can be expressed as $C = I_N \otimes C^k$ with $C^k = \{C_0, C_1, \dots, C_{p-1}\}$; and $f = \{f^k, 0_p^T, \dots, 0_p^T\}^T$ is a $p \times 1$ response vector, in which 0_p^T is a $p \times 1$ zero vector and I_N is an $N \times N$ identity matrix. By decomposing $W = W_q - BW_a$, where $W_q = C(C^H C)^{-1}f$, $B = I_N \otimes B^k$, in which \otimes denotes the Kronecker product, the LCMV beamformer can be implemented via the GSC structure and can be formulated as an unconstrained optimization problem:

$$\min_{W_a} J = \min_{W_a} (W_q - BW_a)^H R_y (W_q - BW_a) \quad (5)$$

And the optimum solution of W_a can be calculated as

$$W_{a, \text{opt}} = (B^H R_y B)^{-1} B^H R_y W_q \quad (6)$$

2.2 Decision feedback GSC

In Refs. [8–9], a new scheme to improve the performance of the traditional GSC algorithm was proposed. Fig. 4 depicts the whole structure of the proposed DF-GSC. Accordingly, the cost function (5) is changed to

$$\min_{W_a} J = \min_{W_a} E\{|X(k) - W_b^* \hat{S}_0(k)|^2\} = \min_{W_a} E\{|(W_q - BW_a)^H Y(k) - W_b^* \hat{S}_0(k)|^2\} \quad (7)$$

where W_b is the feedback tap weight and $\hat{S}_0(k)$ is the detected desired signal. The optimum solution of (7) is

$$W_{a, \text{opt}} = (B^H R_y B)^{-1} B^H R_y W_q, \quad W_{b, \text{opt}} = a(\theta_0) W_q \quad (8)$$

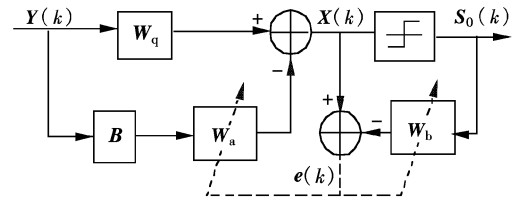


Fig. 4 Structure of DF-GSC

2.3 Leakage constraints decision feedback GSC

Theoretically, the output of the blocking matrix B is a noise reference signal; i. e., $d_i(k) = v_i(k)$, $k = 1, 2, \dots, M-1$. In practical application, the vector d often contains part of the speech signal which is a so-called speech leakage; i. e., $d_i(k) = x_i(k) + v_i(k)$, $k = 1, 2, \dots, M-1$; hence, the adaptive filter will also remove part of the speech signal from the speech reference signal. In this case speech distortion is introduced. In this section we consider a new beamformer which incorporates the new speech leakage constraints with the decision feedback GSC technique. The cost function is changed to

$$J = \min_{\mathbf{W}_a, \mathbf{W}_b} E\{ |[(\mathbf{W}_q - \mathbf{B}\mathbf{W}_a)^H \mathbf{Y}(k) - \mathbf{W}_b^* S_0(k)]|^2 \} + \mu E\{(\mathbf{W}_a^H \mathbf{X}(k))^2\} \quad (9)$$

This cost function consists of a term $E\{(\mathbf{W}_a^H \mathbf{X}(k))^2\}$ related to the speech leakage, then

$$J = \min_{\mathbf{W}_a, \mathbf{W}_b} E\{ |[(\mathbf{W}_q - \mathbf{B}\mathbf{W}_a)^H \mathbf{Y}(k) - \mathbf{W}_b^* S_0(k)]|^2 \} + \mu E\{(\mathbf{W}_a^H \mathbf{X}(k))^2\} = \min_{\mathbf{W}_a, \mathbf{W}_b} E\{ |[(\mathbf{W}_q - \mathbf{B}\mathbf{W}_a)^H \mathbf{Y}(k) - \mathbf{W}_b^* S_0(k)]|^2 \} + \mu E\{(\mathbf{W}_a^H \mathbf{X}(k))^2\} = \min_{\mathbf{W}_a, \mathbf{W}_b} \mathbf{W}_q^H \mathbf{R}_Y \mathbf{W}_q -$$

$$\mathbf{W}_q^H [\mathbf{R}_Y \mathbf{B} \mathbf{P}] \mathbf{W}_c - \mathbf{W}_c^H \left[\begin{array}{c} \mathbf{B}^H \mathbf{R}_Y \\ \mathbf{P}^H \end{array} \right] \mathbf{W}_q + \mathbf{W}_c^H \mathbf{R}_Y \mathbf{W}_c +$$

$$\mu E\{(\mathbf{W}_a^H \mathbf{X}(k))^2\}$$

where

$$\mathbf{W}_c = \begin{bmatrix} \mathbf{W}_a \\ \mathbf{W}_b \end{bmatrix}$$

$$\mathbf{R}_c = E\left\{ \left[\begin{array}{c} \mathbf{B}^H \mathbf{Y}(k) \\ S_0^*(k) \end{array} \right] \left[\mathbf{Y}^H(k) \mathbf{B} \quad S_0^*(k) \right] \right\} = \begin{bmatrix} \mathbf{B}^H \mathbf{R}_Y \mathbf{B} & 0 \\ \mathbf{0}^H & \sigma_{S_0}^2 \end{bmatrix}$$

$$\mathbf{P} = E\{\mathbf{Y}(k) S_0^*(k)\} = \sigma_{S_0}^2 \mathbf{a}(\theta_0)$$

in which $\sigma_{S_0}^2$ is the power of the transmitted desired signal.

$$\text{From } \frac{\partial J}{\partial \mathbf{W}_a^*} = -\mathbf{B}^H \mathbf{R}_Y \mathbf{W}_q + \mathbf{B}^H \mathbf{R}_Y \mathbf{B} \mathbf{W}_a + \mu \mathbf{R}_X \mathbf{W}_a =$$

$\mathbf{0}$, we obtain

$$\mathbf{W}_{a, \text{opt}} = [\mu \mathbf{R}_X + \mathbf{B}^H \mathbf{R}_Y \mathbf{B}]^{-1} \mathbf{B}^H \mathbf{R}_Y \mathbf{W}_q \quad (10)$$

where assuming that the speech and noise signals are uncorrelated, \mathbf{R}_X can be estimated as

$$\mathbf{R}_X = \mathbf{R}_Y - \mathbf{R}_V$$

where \mathbf{R}_Y is estimated during the periods of speech + noise and \mathbf{R}_V during the periods of noise only.

$$\text{From } \frac{\partial J}{\partial \mathbf{W}_b^*} = -\mathbf{W}_q \mathbf{P} + \sigma_{S_0}^2 \mathbf{W}_b = \mathbf{0}, \text{ we obtain}$$

$$\mathbf{W}_{b, \text{opt}} = \frac{\mathbf{P}^H \mathbf{W}_q}{\sigma_{S_0}^2} = \mathbf{a}^H(\theta_0) \mathbf{W}_q \quad (11)$$

The parameter μ gives the trade off between distortion of the speech reference and noise reduction. The greater the amount of speech leakage, the more attention is paid to speech distortion. For $\mu = 0$, all emphasis is placed on the noise reduction and the speech leakage is not taken into account, which corresponds to the GSC-solution. Hence, the LCDF-GSC encompasses the GSC algorithm as a special case. For $\mu = \infty$, all emphasis is placed on speech leakage, $\mathbf{W}_a = \mathbf{0}$, so $\mathbf{X}(k)$ is equal to the output of the fixed beamformer and no speech leakage is calculated, which corresponds to the delay-and-sum(DS) beamformer solution.

2.4 Analysis of output SNR

Let $J(\infty)$ denote the MSE in steady state of the adaptive algorithm:

$$J(\infty) = J_{\min} + J_{\text{ex}}(\infty) \quad (12)$$

where $J_{\text{ex}}(\infty)$ is the excess MSE of the adaptation.

$$J_{\min} = \min_{\mathbf{W}_a, \mathbf{W}_b} E\{ |[(\mathbf{W}_q - \mathbf{B}\mathbf{W}_{a, \text{opt}})^H \mathbf{Y}(k) - \mathbf{W}_b^* S_0(k)]|^2 \} + \mu E\{(\mathbf{W}_{a, \text{opt}}^H \mathbf{X}(k))^2\} = \mathbf{W}_q^H \mathbf{R}_Y \mathbf{W}_{\text{opt}} - \sigma_{S_0}^2 | \mathbf{W}_{\text{opt}}^H \mathbf{a}(\theta_0) |^2 + \mu \mathbf{W}_{a, \text{opt}}^H \mathbf{R}_X \mathbf{W}_{a, \text{opt}} = \mathbf{W}_q^H \mathbf{R}_Y \mathbf{W}_{\text{opt}} - \mathbf{P}_S + \mathbf{P}_L \quad (13)$$

where $\mathbf{P}_S = \sigma_{S_0}^2 | \mathbf{W}_{\text{opt}}^H \mathbf{a}(\theta_0) |^2$, $\mathbf{P}_L = \mu \mathbf{W}_{a, \text{opt}}^H \mathbf{R}_X \mathbf{W}_{a, \text{opt}}$.

The output signal-to-noise ratio (SNR) can be written as

$$\text{SNR}(k) = \frac{E\{ | \mathbf{W}^H(k) \mathbf{a}(\theta_0) S_0(k) |^2 \}}{E\{ | \mathbf{W}^H(k) \mathbf{Y}(k) - \mathbf{W}^H(k) \mathbf{a}(\theta_0) S_0(k) |^2 \} + \mu E\{(\mathbf{W}_a^H \mathbf{X}(k))^2\}} = \frac{\mathbf{P}_S}{\mathbf{W}_q^H \mathbf{R}_Y \mathbf{W}_{\text{opt}} - \mathbf{P}_S + \mathbf{P}_L + J(\infty)}$$

such that

$$\text{SNR}_{\text{LCDF-GSC}} = \lim_{k \rightarrow \infty} \text{SNR}(k) = \frac{\mathbf{P}_S}{\mathbf{W}_q^H \mathbf{R}_Y \mathbf{W}_{\text{opt}} - \mathbf{P}_S + \mathbf{P}_L} \quad (14)$$

The output signal-to-noise ratio of DF-GSC and GSC can be respectively written as

$$\text{SNR}_{\text{DF-GSC}} = \frac{\mathbf{P}_S}{\mathbf{W}_q^H \mathbf{R}_Y \mathbf{W}_{\text{opt}} - \mathbf{P}_S} \quad (15)$$

$$\text{SNR}_{\text{GSC}} = \frac{\mathbf{P}_S}{J_{\min(\text{GSC})} - \mathbf{P}_S} = \frac{\mathbf{P}_S}{\mathbf{f}^H (\mathbf{C}^H \mathbf{R}_X^{-1} \mathbf{C})^{-1} \mathbf{f} - \mathbf{P}_S} \quad (16)$$

Since in DF-GSC, $\mathbf{W}_q^H \mathbf{R}_Y \mathbf{W}_{\text{opt}} - \mathbf{P}_S \cong \mathbf{0}$, the MMSE of DF-GSC is smaller than that of the GSC. As a result, the output SNR of DF-GSC is higher than that of the GSC. For $\mu = 0$, Eq. (10) shows that the solution of $\mathbf{W}_{a, \text{opt}}$ is the same as that of DF-GSC, and so is the output SNR of LCDF-GSC. For $\mu = 0$, $\mathbf{W}_a = \mathbf{0}$, $\mathbf{X}(k)$ is equal to the output of the fixed beamformer, and the output SNR corresponds to the weight \mathbf{W}_q .

3 Results

This section discusses the performance (SNR improvement and speech distortion) of the LCDF-GSC algorithm.

Segmental SNR is defined as

$$G_{\text{SEGSNR}} = \frac{1}{L} \sum_{l=1}^{L-1} 10 \log \frac{\frac{1}{N} \sum_{n=0}^{N-1} s^2(n + Nl)}{\frac{1}{N} \sum_{n=0}^{N-1} [s(n + Nl) - \hat{s}(n + Nl)]^2} \quad (17)$$

where L is the frame length, $s(n)$ is the original signal, and $\hat{s}(n)$ is the estimated signal.

Speech distortion will be analyzed by considering

Itakura-saito distance between the speech component of the first microphone signal and the speech component of the considered signal. Itakura-saito distance is defined as

$$d_{IS}(\bar{\alpha}_d, \bar{\alpha}_\phi) = \left[\frac{\sigma_\phi^2}{\sigma_d^2} \right] \left[\frac{\bar{\alpha}_d \mathbf{R}_\phi \bar{\alpha}_d}{\bar{\alpha}_\phi \mathbf{R}_\phi \bar{\alpha}_\phi} \right] + \log \left(\frac{\sigma_d^2}{\sigma_\phi^2} \right) - 1 \quad (18)$$

where $\bar{\alpha}_\phi$ is an original clean frame of speech with linear prediction (LP) coefficient vector, and $\bar{\alpha}_d$ is the processed speech coefficient vector. σ_d^2 and σ_ϕ^2 represent the all-pole gains for the processed and clean speech frames, respectively. We have calculated this distance with an LPC-order of 12.

The white noise from the Noisex-92 database was used for the objective evaluation of the proposed algorithm. The speech signals were four sentences taken from two female speakers and two male speakers, respectively. We compare the performance of the LCDF-GSC algorithm with DF-GSC and the single channel speech enhancement algorithm of the Wiener filter. In the experiment, we choose $N=4, \mu=0.01$ in Eq. (10).

Fig. 5 depicts the SNR improvement in white noise. It can be seen that the SNR improvement of multi-channel speech enhancement is higher than that of single-channel speech enhancement, and the SNR improvement of DF-GSC is higher than that of LCDF-GSC due to the introduction of speech leakage. However, from Fig. 6, it can be seen that the IS distance of DF-GSC is higher than the proposed algorithm and lower than that of the Wiener filter. So the proposed algorithm has lower distortion than the other two algorithms. Finally, we assume that there is a 2° difference between the estimated and the actual desired signal's DOA^[9]. From Fig. 7 we can see that the mismatch affects conventional GSC more than LCDF-GSC in SNR improvement.

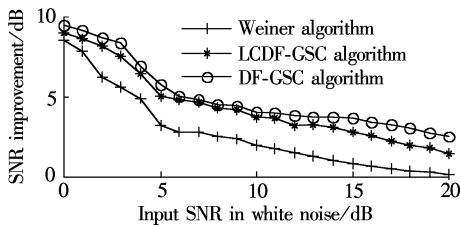


Fig. 5 SNR improvement in white noise

We adopted the MOS score to demonstrate the informal listening test results. Everyone listened to the four enhanced speech signals two times and gave eight scores. The results of the average score show that the speech enhanced by LCDF-GSC has more quality improvement than the other two algorithms.

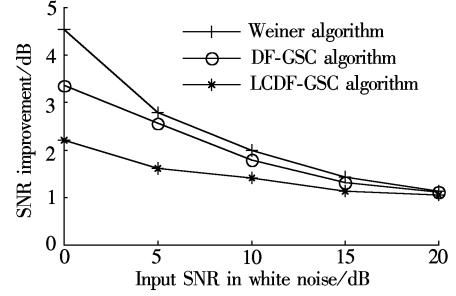


Fig. 6 IS distance measurement

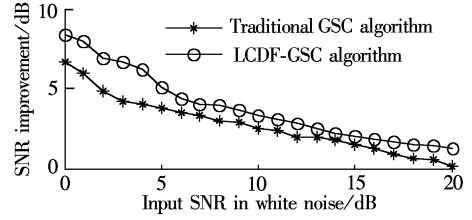


Fig. 7 SNR improvement in white noise with mismatch

4 Conclusion

In this paper, a new LCDF-GSC algorithm has been proposed for speech enhancement. With the introduction of leakage constraints, LCDF-GSC achieves less speech distortion than the DF-GSC algorithm and Wiener filtering, and simulation results also show that LCDF-GSC can maintain high SNR values even in mismatch than the traditional GSC algorithm.

References

- [1] Sreenivas T V, Kimapure Pradeep. Codebook constrained Wiener filtering for speech enhancement [J]. *IEEE Transactions on Speech and Audio Processing*, 1996, **4**(5): 383 – 389.
- [2] Hu Yi, Loizou Philipos C. A generalized subspace approach for enhancing speech corrupted by colored noise[J]. *IEEE Transactions on Speech and Audio Processing*, 2003, **11** (4): 334 – 340.
- [3] Frost O L III. An algorithm for linearly constrained adaptive array processing[J]. *Proceedings of IEEE*, 1972, **60**(8): 926 – 935.
- [4] Griffiths L J, Jim C W. An alternative approach to linearly constrained adaptive beamforming[J]. *IEEE Trans Antennas Propag*, 1982, **30**(1): 27 – 34.
- [5] Spriet Ann, Moonen Marc, Wouters Jan. Stochastic gradient-based implementation of spatially preprocessed speech distortion weighted multichannel Wiener filtering for noise reduction in hearing aids [J]. *IEEE Trans Signal Process*, 2005, **53**(3): 911 – 925.
- [6] Hoshuyama O, Sugiyama A, Hirano A. A robust adaptive beamformer for microphone arrays with a blocking matrix using constrained adaptive filters [J] . *IEEE Trans Signal*

- Process*, 1999, **47**(10): 2677 – 2683.
- [7] Nordebo S, Claesson I, Nordholm S. Adaptive beamforming: spatial filter designed blocking matrix[J]. *IEEE Journal of Oceanic Engineering*, 1994, **19**(4): 583 – 590.
- [8] Lee Yinman, Wu Wen-Rong. An LMS-based adaptive generalized sidelobe canceller with decision feedback[C]// *IEEE International Conference on Communications*. Seoul, Korea, 2005, **3**: 2047 – 2051.
- [9] Lee Yinman, Wu Wen-Rong. A robust adaptive generalized sidelobe canceller with decision feedback[J]. *IEEE Trans Antennas Propag*, 2005, **53**(11): 3822 – 3832.
- [8] Lee Yinman, Wu Wen-Rong. An LMS-based adaptive generalized sidelobe canceller with decision feedback[C]//

基于泄漏约束的 DF-GSC 语音增强

邹采荣^{1,2} 陈国明¹ 赵 力¹

(¹ 东南大学信息科学与工程学院, 南京 210096)

(² 佛山科学技术学院, 佛山 528000)

摘要: 为了改善广义旁瓣抵消(GSC)语音增强方法的性能,提出了一种带有泄漏约束的判决反馈旁瓣抵消(LCDF-GSC)方法. 采用 DF-GSC 方法以解决 GSC 对波达方向敏感的问题,在代价函数中引入泄漏因子,以此改善语音失真的问题,而这种问题是由于噪声参考信号中含有语音成份造成的. 试验结果表明,尽管经过 LCDF-GSC 处理后的语音信号信噪比要略低于 DF-GSC,IS 测度表明这时前者的语音信号失真度要小于后者. MOS 分也表明 LCDF-GSC 方法要优于 DF-GSC 和单通道 Weiner 滤波算法.

关键词: 语音增强; 广义旁瓣抵消; 语音泄漏

中图分类号: TN912.35