

基于阵列抗串扰自适应噪声抵消的语音增强

曾庆宁^{1,2}, 欧阳缮², 廖桂生¹

(1. 西安电子科技大学通信工程学院, 陕西西安 710071; 2. 桂林电子工业学院通信与信息工程系, 广西桂林 541004)

摘要: 本文提出了一种针对阵列交叉串扰信号的自适应噪声抵消方法, 并将其用于麦克风阵语音增强. 该方法仅使用两级滤波系统, 计算量小, 稳定性好, 且对麦克风阵的几何结构及噪声类型均没有严格限制. 试验表明, 该方法消噪量大, 对语音损伤小, 语音增强效果显著, 适用于多种噪声环境并易于实时实现.

关键词: 语音增强; 串扰; 自适应噪声抵消; 麦克风阵

中图分类号: TN957. 52 **文献标识码:** A **文章编号:** 0372-2112 (2005) 02-0241-04

Speech Enhancement Based on Array Cross-Talk Resistant Adaptive Noise Cancellation

ZENG Qing-ning^{1,2}, OUYANG Shan², LIAO Gui-sheng¹

(1. College of Telecommunications Engineering, Xidian University, Xi'an, Shaanxi 710071, China;

2. Department of Telecommunications and Information Engineering, Guilin Institute of Electronic Technology, Guilin, Guangxi 541004, China)

Abstract: An adaptive noise cancellation method is presented to extract desired signal from array cross-talk signals, and is used to microphone array speech enhancement. Employing a two-stage filtering system, the method not only has a character of small amount of computations and high stability, but also has no strict limitations to the geometric structure of the microphone array and even to the noise types. Experiment results show that the method is quite effective in speech enhancement as it can cancel a large amount of noise with little distortion to speech signals, and that the method is suitable to many kinds of noisy environments and easy for real-time implementation.

Key words: speech enhancement; cross-talk; adaptive noise cancellation; microphone array

1 引言

众所周知, 语音通讯及语音识别等都常受到背景噪声的干扰, 严重影响通话质量及语音识别系统的识别率. 以消除噪声并提取尽可能纯净的原始语音为主要目的的语音增强研究, 具有重要的应用价值.

总结国内外的研究成果, 容易发现自适应噪声抵消法 (Adaptive Noise Cancellation, 简为 ANC) 是语音增强的一种主要方法. 它可以适应多种噪声环境, 对语音损伤小, 计算量不大, 是一种实用的语音增强方法. 近年来, 研究者们认识到: 由于噪声传播及声场的复杂性, 基于一个或两个麦克风的消噪系统的实际消噪性能是有限的. 要提高消噪效果, 需要增加麦克风数目以组成麦克风阵来获得, 因此, 麦克风阵消噪系统已成为一个研究热点, 并在该领域取得了一系列成果^[1~3].

实际应用中, 许多场合要求麦克风阵总体尺寸不宜过大^[2,3]. 因此, 在阵列中各麦克风拾取的信号之间, 语音与噪声的交叉串扰 (cross-talk) 现象是十分严重的, 如果要用噪声抵消法进行语音增强, 一般的抵消方法效果将会很差. 语音信号往往将会一起被抵消掉. 文献[4~7]针对仅两路交叉串扰信

号的场合, 分别讨论了几种方案, 然而, 仅两路信号的噪声抵消法, 其语音增强效果有限, 而且我们发现文[4~6]的方法容易导致算法不稳定, 甚至经常发散. 文[7]主要探讨生物医学信号处理问题, 而且使用了三级滤波系统进行处理. 本文中我们提出一种针对阵列交叉串扰信号的噪声抵消方法, 即阵列抗交叉串扰自适应噪声抵消法 (Array Cross-talk Resistant Adaptive Noise Cancellation, 简为 ACRANC), 并将其用于麦克风阵语音增强. 该方法仅使用两级滤波系统, 计算量不大, 稳定性好, 非常利于实时处理, 而且对麦克风阵的几何结构、麦克风与麦克风之间的距离以及噪声的类型均没有严格要求. 试验表明, 该方法消噪量大, 对语音的损伤小, 语音增强效果显著, 而且适用于多种噪声环境, 并完全可以用目前普通的 DSP 器件实时实现.

2 阵列抗串扰自适应噪声抵消法

如图1所示, 设麦克风阵的 $N+1$ 个麦克风为 M_0, M_1, \dots, M_N , 语音信号 $s(t)$ 通过冲击响应为 $h_{s_i}(t)$ 的空间传播 (包括直射、反射等) 到达麦克风 M_i 并转换为语音信号 $s_i(t)$, 噪声信号 $n(t)$ (可以是一个或多个噪声源而成) 通过 $h_{n_i}(t)$ (包

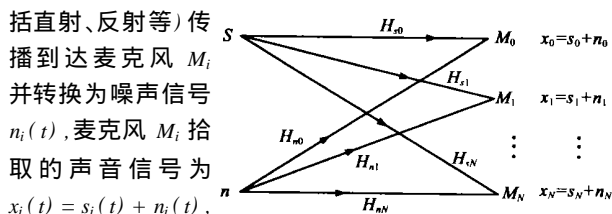


图 1 语音与噪声信号交叉串扰示意图

括直射、反射等)传播到达麦克风 M_i 并转换为噪声信号 $n_i(t)$, 麦克风 M_i 拾取的声音信号为 $x_i(t) = s_i(t) + n_i(t)$, $i = 0, 1, \dots, N$. 每个麦克风 M_i 实际获得的信号 $x_i(t)$ 其实都包括语音信号和噪声信号, 此即语音与噪声信号的交叉串扰^[4-7] (cross-talk). 对 $s(t)$ 、 $s_i(t)$ 、 $h_{s_i}(t)$ 、 $n(t)$ 、 $n_i(t)$ 、 $h_{n_i}(t)$ 及 $x_i(t)$ 抽样后的离散信号分别仍记为 $s(k)$ 、 $s_i(k)$ 、 $h_{s_i}(k)$ 、 $n(k)$ 、 $n_i(k)$ 、 $h_{n_i}(k)$ 及 $x_i(k)$, 显然

$$x_i(k) = s_i(k) + n_i(k), \quad i = 0, 1, \dots, N \quad (1)$$

$$s_i(k) = h_{s_i}(k) * s(k), \quad i = 0, 1, \dots, N \quad (2)$$

$$n_i(k) = h_{n_i}(k) * n(k), \quad i = 0, 1, \dots, N \quad (3)$$

图 1 中麦克风阵可以是线、面或体阵, 将其中的一个麦克风视为主麦克风(通常选为中心位置的麦克风), 其余的视为参考麦克风, 为方便计, 不妨将主麦克风记为 M_0 , 参考麦克风记为 M_1, \dots, M_N . 令 s_i 到 s_j 的冲击响应为 $h_{s_{ij}}(k)$, n_i 到 n_j 的冲击响应为 $h_{n_{ij}}(k)$, 即

$$s_j(k) = h_{s_{ij}}(k) * s_i(k), \quad i, j = 0, 1, \dots, N \quad (4)$$

$$n_j(k) = h_{n_{ij}}(k) * n_i(k), \quad i, j = 0, 1, \dots, N \quad (5)$$

则由(2)~(5)可见

$$H_{s_{ij}}(z) = \frac{H_{s_i}(z)}{H_{s_i}(z)}, \quad i, j = 0, 1, \dots, N \quad (6)$$

$$H_{n_{ij}}(z) = \frac{H_{n_i}(z)}{H_{n_i}(z)}, \quad i, j = 0, 1, \dots, N \quad (7)$$

其中 $H_{s_{ij}}(z)$ 是 $h_{s_{ij}}(k)$ 的 z 变换, 其余类同.

图 2 中使用两级滤波系统来完成交叉串扰自适应噪声抵消而提取语音 s_0 , 且噪声抵消是基于“语音信号可分为‘有’和‘无’

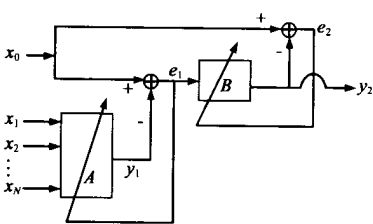


图 2 阵列抗串扰自适应噪声抵消原理图

首先, 在无语音信号阶段, 麦克风 M_0, M_1, \dots, M_N 拾取纯粹的噪声信号 n_0, n_1, \dots, n_N , 视 n_0 为主信号, 并将 n_1, \dots, n_N 用作参考信号输入滤波器 A, 目的是抵消 n_0 .

我们知道, 在两通道自适应噪声抵消系统中, 除去一路主信号外, 只有一路参考信号, 即只能用某一路 n_i 输入滤波器 A 去抵消 n_0 , 而且实际中 A 往往取为 FIR 滤波器, 即

$$n_0(k) = w_i n_i(k) + err_i(k) \quad (8)$$

$$w_i = (w_{i0}, w_{i1}, \dots, w_{iL}) \quad (9)$$

$$\bar{n}_i(k) = [n_i(k), n_i(k-1), \dots, n_i(k-L)]^T \quad (10)$$

w_i 为滤波器系数, L 为滤波器阶数, $err_i(k)$ 为使用 $\bar{n}_i(k)$ 对 $n_0(k)$ 进行预测的预测误差, 误差功率 $P[err_i(k)]$ 的大小将直接

影响到语音增强的效果. 但在实际环境下, 一般说来, 无论如何优化 w_i , 甚至进一步增大或优选 L , 式(8)的最小误差功率 $P[err_i(k)]$ 往往远不够小, 这可能主要是由于噪声在空间传播的直射、反射、回响等复杂因素造成的, 致使 M_i 未能获得与 M_0 完全相关的噪声信号, 而且已经证实麦克风 M_i 离 M_0 越远, M_0 与 M_i 的信号相关性越弱, $P[err_i(k)]$ 越大.

如果将多路参考信号 n_1, \dots, n_N 输入 FIR 滤波器 A 来抵消 n_0 , 则

$$n_0(k) = w \bar{n}(k) + err(k) \quad (11)$$

$$w = (w_1, w_2, \dots, w_N) \quad (12)$$

$$w_i = (w_{i0}, w_{i1}, \dots, w_{iL}) \quad (13)$$

w 为滤波器系数, 它是一个 $N(L+1)$ 维的行向量, 此时滤波器 A 的阶数为 $N(L+1) - 1$; 而

$$\bar{n}(k) = [\bar{n}_1(k), \bar{n}_2(k), \dots, \bar{n}_N(k)]^T \quad (14)$$

$$\bar{n}_i(k) = [n_i(k), n_i(k-1), \dots, n_i(k-L)] \quad (15)$$

$\bar{n}(k)$ 为 $N(L+1)$ 维的列向量; $err(k)$ 为预测误差.

实际中发现, 如果适当地选择 L 和 N , 通过优化 w , 式(11)的最小误差功率常常满足:

$$P[err^*(k)] < P[err_i^*(k)] \quad (16)$$

而且对宽带或多噪声源等噪声, $P[err^*(k)]$ 往往比 $P[err_i^*(k)]$ 小得多. 这意味着主麦克风 M_0 中的噪声信号 n_0 通过参考麦克风 M_1, \dots, M_N 中的噪声信号 n_1, \dots, n_N 进行抵消后剩余的误差噪声比仅通过单个参考麦克风 M_i 中的噪声信号 n_i 进行抵消后剩余的误差噪声往往更小, 其实这也正是使用麦克风阵的优越性之所在. 这一现象说明增加参考麦克风能更全面接收与主麦克风 M_0 相关的噪声信号.

记式(11)中达到最小预测误差功率 $P[err^*(k)]$ 的最优滤波器系数为

$$w^* = (w_1^*, w_2^*, \dots, w_N^*) = (w_{10}^*, w_{11}^*, \dots, w_{1L}^*, w_{20}^*, w_{21}^*, \dots, w_{2L}^*, \dots, w_{N0}^*, w_{N1}^*, \dots, w_{NL}^*) \quad (17)$$

显然, 为使误差噪声功率 $P[err(k)]$ 最小, 调整图 2 中滤波器 A 的系数 w 以使 e_1 的平方和最小即可.

然后, 在有语音信号阶段, 由于麦克风阵中各个麦克风的相对位置不变, 所在的声场环境短时间内也基本不变, 因此在前段无语音信号时得到的噪声传播的冲击响应 $h_{n_{i0}}(k)$ 可认为不会马上改变(即假设噪声环境是短时不变或缓变的), 因此有:

$$y_1(k) = w^* \bar{x} = w^* (\bar{s} + \bar{n}) = w^* \bar{s} + w^* \bar{n} = w^* \bar{s} + [n_0(k) - err^*(k)] \quad (18)$$

这里 \bar{s} 、 \bar{x} 类似式(14)、(15)的 \bar{n} , 比如

$$\bar{s}(k) = [s_1(k), s_2(k), \dots, s_N(k)]^T$$

$$\bar{s}_i(k) = [s_i(k), s_i(k-1), \dots, s_i(k-L)]$$

于是

$$\begin{aligned} e_1(k) &= x_0(k) - y_1(k) \\ &= [s_0(k) + n_0(k)] - [w^* \bar{s} + n_0(k) - err^*(k)] \\ &= s_0(k) - w^* \bar{s} + err^*(k) = p(k) + err^*(k) \end{aligned} \quad (19)$$

其中

$$p(k) = s_0(k) - w^* \bar{s} \quad (20)$$

对式(20)两边取 z 变换有

$$\begin{aligned} P(z) &= S_0(z) - Z \left[\sum_{i=1}^N \sum_{j=0}^L W_{ij}^* s_i(k-j) \right] \\ &= S_0(z) - \sum_{i=1}^N \sum_{j=0}^L w_{ij}^* z^{-2j} H_{s_i s_0}(z) S_0(z) \\ &= \left[1 - \sum_{i=1}^N \sum_{j=0}^L w_{ij}^* z^{-2j} H_{s_i s_0}(z) \right] S_0(z) = \tilde{H}(z) S_0(z) \quad (21) \end{aligned}$$

其中

$$\tilde{H}(z) = 1 - \sum_{i=1}^N \sum_{j=0}^L w_{ij}^* z^{-2j} H_{s_i s_0}(z) \quad (22)$$

由此可见 $p(k)$ 是由语音信号 $s_0(k)$ 所决定的信号。

进而, $p(k)$ 的功率 $P[p(k)]$ 通常不会“很小”,即对含噪语音,经过滤波,虽然噪声功率被大幅降低了,但语音的功率却没有降低或相对大幅降低的噪声功率而言降低得不多。这其实是因为语音与噪声传播到麦克风的系统函数通常不同而保障的(严格说来应与空间环境及麦克风阵的结构有关),实际上这也是阵列信号处理的零点形成技术,即通过滤波,在噪声的所有直射、反射、折射等方向形成了零点,但由于语音方向通常与噪声方向不同,其所有直射、反射、折射等方向一般不会全都落到这些噪声方向的零点上。

所以,在有语音阶段,含噪语音通过第一级滤波,其输出 $e_1(k)$ 的信噪比(此时信号为 $p(k)$,噪声为 $err^*(k)$),见式(19)通常比滤波前含噪语音 $s_0(k)$ 的信噪比更大,也通常比通过只有一个参考麦克风的滤波系统滤波后输出的信噪比更大。但是, $p(k)$ 并非语音信号 $s_0(k)$,而是 $s_0(k)$ 的畸变信号,多麦克风时畸变尤其严重。由滤波器 B 组成的第二级滤波系统,其目的正是为了将畸变的语音信号 $p(k)$ 恢复成 $s_0(k)$ 。

为此,我们不妨假定主麦克风 M_0 拾取的噪声信号 n_0 与语音信号 x_0 不相关,而且 n_0 的均值为零。于是只要调整滤波器 B 的系数使图 2 中 e_2 的平方和最小,那么滤波器 B 的输出 y_2 即为估计的语音 s_0 ,这是因为

$$\begin{aligned} e_2(k)^2 &= x_0(k) - y_2(k) = s_0(k) + n_0(k) - y_2(k) \\ &= n_0(k) + s_0(k) - y_2(k) \\ &\quad + 2n_0(k)(s_0(k) - y_2(k)) \quad (23) \end{aligned}$$

$$E(e_2^2) = E(n_0^2) + E[(s_0 - y_2)^2] \quad (24)$$

让 e_2 的平方和最小即让 $E[(s_0 - y_2)^2]$ 最小,等价于让 y_2 与 s_0 的误差最小,鉴于滤波器 B 的输入 e_1 由式(19)~(22)可见与 s_0 有关,故用它来抵消 x_0 中的 s_0 是可行的。亦即让滤波器 B 的输出 y_2 最接近 s_0 是可行的。 e_1 的信噪比(信号为 $p(k)$,噪声为 $err^*(k)$)越大, y_2 与 s_0 的误差越小。

实际中,滤波器 B 通常也用 FIR 滤波器作为近似。不难看出,滤波器 B 的系统函数实际上就是 $\tilde{H}(z)$ 的逆系统 $\tilde{H}^{-1}(z)$ 的近似。

总而言之,对图 2 所示的阵列抗串扰自适应噪声抵消系统,只要调节其中 FIR 滤波器 A 、 B 的系数分别使 e_1 、 e_2 的平方和最小,则滤波器 B 的输出 y_2 即为语音 s_0 的估计,其估计误差一般比使用两个麦克风时的估计误差更小。

注意:滤波器 A 的系数仅在无语音信号的时间段内进行自适应调整,滤波器 B 的系数则仅在有语音信号的时间段内进行自适应调整。由于滤波器 A 、 B 系数的自适应调整其实是为了跟踪语音及噪声传播过程的冲击响应的变化,所以,当噪声源和语音源的位置不变,而且它们所处的空间环境也不变时,这些冲击响应就基本不变,因此,此时滤波器 A 、 B 的系数便不必更新。同样,由于这一原因,当上述语音消噪方法应用于噪声源、语音源及空间环境不变或变化较缓的场合时,对含噪语音信号的有声/无声判决的要求就可以不高,只需在肯定为无声时才调整滤波器 A 的系数,并在肯定为有声时才调整滤波器 B 的系数,而在不好确定的阶段,则不妨保持滤波器 A 、 B 的系数不变,只作滤波即可。

对滤波器 A 、 B 系数的优化有 LMS、NLMS、RLS、BTFI、LSLL 等自适应算法^[7-9],可根据所用 DSP 器件的性能及实际问题的需要适当选择。若采用较复杂的 LSLL 的算法,则不难算出上述语音消噪法在每抽样时间间隔内需要的浮点运算次数不超过 $23(L_{\max} + 1)$ (其中 L_{\max} 为滤波器 A 与 B 的最高阶数),这里未计通常按帧处理的有声/无声判决所需要的运算次数,因它平均到每抽样间隔后相对较少,尤其是在上述语音消噪法对有声/无声判决要求并不很高的情况下更是如此。实际应用中, L_{\max} 一般可取为小于 100,因此,如果使用 $8k$ 的常用采样频率,则每秒运算量至多为 19.4MHOPS。假如我们采用的是最简单的 LMS 算法,相应的每秒的运算量则至多为 2.4MHOPS。可见用目前普通的 DSP 即可轻易实现本文方法的实时处理(例如 TI 公司的 TMS320VC33 具有每秒 150MHOPS 的运算能力)。

3 试验结果

我们所做的多个仿真和具体试验说明了本文方法的有效性,这里仅给出其中一个具体试验的结果。该试验中,麦克风阵由四个小麦克风组成,各麦克风之间距离很近,阵列孔径颇小。阵列信号快拍采样频率为 8kHz,语音为真实的语音“本段语音用于测试”,噪声为收音机失谐时的噪音,语音开始前有约 1.5 秒的纯噪声,语音源与噪声源不在同一位置,试验在一普通的 18 平方米的室内进行。

图 3(a) 是阵列中麦克风 M_0 采集到的含噪语音信号 x_0 ,图中只画出了从 0.5 秒至第 3.625 秒的图形,它在有语音阶段的语音与噪声的信噪比为 2.753dB,其他麦克风 M_1 、 M_2 、 M_3 拾取的含噪语音信号的相应信噪比与 x_0 的大致相同。通过试听,这些混合信号中的语音都被噪声严重污染,听者很难听清说话人所说的内容。

图 3 的 (b)~(f) 展示了本文阵列抗串扰自适应噪声抵消法的优越性,也显示了使用阵列比仅使用两个麦克风的优越性。图中同样只画出了从 0.5 秒至 3.625 秒的有关信号图,而从 0 至 0.5 秒是使用 LSLL 自适应滤波算法时的初始收敛过程。

图 3(b) 为使用一般的两路自适应噪声抵消法所增强的语音,即通过主麦克风 M_0 和参考麦克风 M_1 拾取的混合信号进行增强的结果,处理过程中滤波器阶数取为 32。由图可见

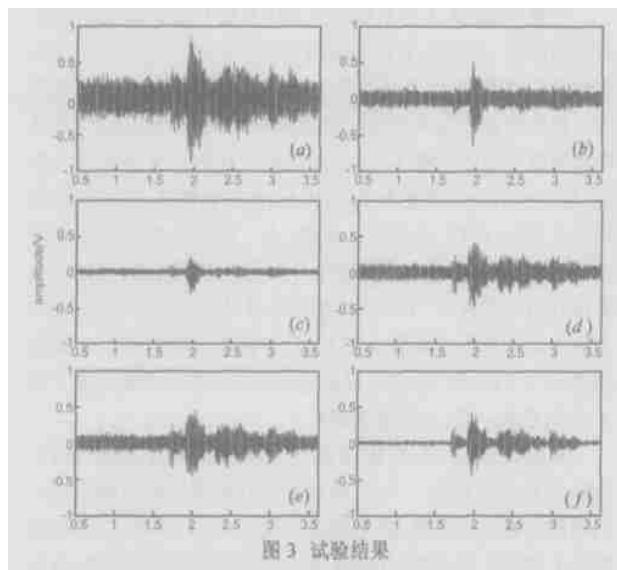


图3 试验结果

残余噪声依然较大,尤其是语音很大部分同时被抵消掉了,相应的信噪比不仅没有改善,反而下降到 1.968dB,试听发现噪声虽有所降低,但语音也有所降低,并且语音有所损伤,这是由于参考信号中串有语音信号,而一般的噪声抵消法并无抗串扰的功能。

图 3(c) 是使用一般的阵列(多通道)自适应噪声抵消法增强的语音,即通过主麦克风 M_0 与参考麦克风 M_1 、 M_2 、 M_3 进行噪声抵消的结果,处理过程中每一路输入滤波器的信号延迟数均为 32,即滤波器的阶数为 98。由图可见残余噪声很小,但语音也几乎被完全抵消,相应的信噪比为 3.039dB,试听发现噪声确实很小,但语音也几乎听不到了,并且语音损伤十分严重。

图 3(d) 是通过 M_0 与 M_1 两路信号使用三级抗串扰自适应噪声抵消法增强的语音,其中三个滤波器的阶数分别取为 32、48、48,所得增强后的语音信噪比为 4.326。试听发现语音有所增强,损伤小,已可以听懂所说内容。

图 3(e) 是通过 M_0 与 M_1 两路信号使用本文的两级抗串扰自适应噪声抵消法增强的语音,其中两个滤波器 A 与 B 的阶数分别为 32、48,所得增强后的语音信噪比为 4.594。注意:(e)需要的计算量比(d)减少约一半,但效果却比(d)的效果略好。试听发现语音有所增强,损伤小,已可以听懂所说内容。

图 3(f) 是通过全部四个麦克风并使用本文的两级阵列抗串扰自适应噪声抵消法增强的语音,处理过程中三路输入滤波器 A 的参考信号延迟数均取为 32(即滤波器 A 的阶数为 98),滤波器 B 的阶数同样取为 48,所得增强后的语音信噪比提高到 13.949dB。试听发现,噪声很小,语音损伤小,所说内容清晰好懂。

如果将(e)中滤波器 A 的阶数提高到 98,则(e)中方法所增强的语音信噪比可提高到 5.133,此时其计算量与(f)的计算量完全相同,可信噪比却远不及(f)的,足见四路比两路效果更佳。

4 结束语

本文提出了一种阵列抗交叉串扰自适应噪声抵消(ACRANC)的方法并用来进行语音增强,该方法仅使用两级滤波系统,稳定性好,易于实时实现,且对阵列结构及噪声类型没有严格限制,对语音损伤小,语音增强效果显著。

参考文献:

- [1] M Dahl, I Claesson, S Nordebo. Simultaneous echo cancellation and car noise suppression employing a microphone array[A]. ICASSP-97[C]. Munich, Germany: ICASSP, 1997. 239 - 242.
- [2] M Brandstein, D Ward. Microphone Arrays[M]. Verlag: Springer, 2001.
- [3] R L Bouquin. Enhancement of noise speech signals: application to mobile radio communications[J]. Speech Communication, 1996, 18(1): 3 - 19.
- [4] R L Zinser, G Mirchandani, J B Evans. Some experimental and theoretical results using a new adaptive filter structure for noise cancellation in the presence of cross-talk [A]. Proc ICASSP[C]. Tampa, 1985, 3: 1253 - 1256.
- [5] G Mirchandani, R L Zinser, J B Evans. A new adaptive noise cancellation scheme in the presence of crosstalk[J]. IEEE Transaction on Circuits and System, 1992, 39(10): 681 - 694.
- [6] S M Kuo, W M Peng. Asymmetric crosstalk-resistant adaptive noise canceller[J]. Proc IEEE Workshop on Signal Processing Systems, 1999. 605 - 614.
- [7] G Madhavan, H D Bruin. Crosstalk resistant adaptive noise cancellation[J]. Annals of Biomedical Engineering, 1990, 18(1): 57 - 67.
- [8] F Ling. Numerically robust LS lattice-ladder algorithm with direct updating of reflection coefficients[J]. IEEE Transaction on Acoustics, Speech and Signal Processing, 1986, 39(7): 837 - 845.
- [9] S Haykin. Adaptive Filter Theory. 北京: 电子工业出版社, 1998. 365 - 696.

作者简介:



曾庆宁 1963年4月生于广西桂林,西安电子科技大学通信工程学院博士生,桂林电子工业学院通信与信息工程系教授,已发表论文30余篇,主要研究方向为语音信号处理、数字图像处理、Markov 决策规划及模糊规划。E-mail: qingning@hotmail.com.



欧阳健 1960年9月生于江西安福,桂林电子工业学院通信与信息工程系教授,博士生导师,02年获全国优秀博士学位论文奖,01年在香港中文大学及03年在美国加州大学河滨分校作研究工作,主要研究方向为自适应信号处理、通信信号处理、神经网络等。