

# 语音信号处理

--前端处理（信号的获取与降噪处理）

<http://blog.csdn.net/shichaog/article/details/52514816>

<http://blog.csdn.net/shichaog/article/details/52403107>

<http://blog.csdn.net/shichaog/article/details/52399354>

Shichaog@126.com

## 目录

第 0 章 信号处理	4
0.1 语音信号特点	4
0.2 信号响应的意义	4
0.3 卷积推导	4
0.4 时域离散系统的输入输出描述法	5
0.4.1 时域离散信号傅里叶变换 (DTFT, Discrete-Time Fourier Transform)	5
0.4.2 周期信号由傅里叶级数表示	5
0.5 离散傅里叶变换 (DFT)	5
0.6 短时傅里叶变换 (STFT, Short-Time Fourier Transform)	7
0.6.1 Overlap-and-Add (OA,OLA) 计算 STFT	7
0.6 FIR 数字滤波器	7
0.7 窗函数设计 FIRDF	8
第一章 重采样算法	8
1.1 信号重采样	8
1.2 音频重采样	10
1.3 sinc 重采样	11
1.4 MATLAB 实现	14
1.5 开源重采样 C 语言实现	16
1.6 基于 FPGA 的多相滤波	16
1.7 重采样性能评估	16
第 2 章 回声消除 (AEC) 原理及实现	17
2.1 回声消除原理	17
2.2 维纳滤波	18
2.3 LMS 算法	19
2.3.1 NLMS	21
2.3.2 SE-LMS	22
2.3.2 SD-LMS	22
2.3.2 SS-LMS	22
2.3.2 LLMS	22
2.3.2 LNLMS	22
2.4 块自适应滤波	23
2.4.1 块自适应滤波器	23
2.4.2 块 LMS	24
2.4.3 块 LMS 算法收敛性	24
2.4.4 块长的选择	25
2.5 FLMS	25
2.6 MDF 自适应权值调整	26
时域解	26
频域解	28
第三章 语音阵列信号处理	29
3.1 阵列模型	31
3.1.1 线阵模型	31
3.1.2 面阵模型	37
3.1.3 圆阵模型	37

3.2 阵列波束形成技术 .....	37
3.2.1 DS 模型 .....	37
3.2.2 最大信噪比 .....	39
3.2.3 最小方差无失真响应滤波器 .....	40
3.2.4 线性约束最小方差 .....	40
3.2.3 广义旁瓣相消 .....	42
3.3 基于阵列定位和跟踪技术 .....	43
3.3.1 互相关方法 .....	43
3.3.2 广义互相关 (GCC Generalized Cross-Correlation) .....	45
3.3.3 基于特征向量的方法 .....	47
3.3.4 最小熵法 .....	50
3.3.5 自适应特征向量分解法 .....	51
3.3.6 自适应盲信号分离 (BSS, Blind Source Separation) .....	52
3.3.1 TDOA .....	53
3.3.3 空域线性预测法 .....	55
3.3.2 SRP-PHAT .....	56
语音信号预加重 .....	58
VAD 算法 .....	58
第五章 ASR (automatic speech recognition) .....	58
5.1 ASR 模型 .....	59
5.1.1 高斯混合模型 GMM(Gaussian mixture model) .....	59
5.1.2 隐马尔可夫模型 HMM (hidden Markov model) .....	61
5.2 MFCC .....	61
5.2.1 Mel 频率分析 (Mel-Frequency Analysis) .....	64
DNN (Deep neutral network) 深度神经网络 .....	67
NLP (native language processing) .....	74
TTS (text to speech) .....	74

## 第 0 章 信号处理

### 0.1 语音信号特点

在一段时间内（10ms~30ms），人的声带和声道形状是相对稳定的，可认为其特征是不变的。语音可以分为周期性的浊音和非周期的清音。浊音和清音经常在一个音节中同时出现。浊音部分和音质关系密切，在时域上呈现出明显的周期性，在频域上有共振峰结构，而且大部分能量集中在较低频段内，是语音中大幅度高能量的部分；清音则具有明显的时域和频域特征，类似于白噪声，能量较小，在强噪声中容易被掩盖，但在较高信噪比时能提供较多的信息。在语音增强中，可以利用浊音的周期性特征，采用梳状滤波器提取语音分量或者抑制非语音信号，而清音则难以与宽带噪声区分。

加性噪声大致上有：周期性噪声、脉冲噪声、宽带噪声和同声道的其它语音干扰等。

周期性噪声主要来源于发动机等周期性运转的机械，电气干扰，特别是电源交流声也会引起周期性噪声，其特点是有许多离散的窄谱峰。脉冲噪声来源于爆炸、撞击和放电等，表现为时域波形中突然出现的窄脉冲。宽带噪声的来源很多，包括热噪声、气流（风、呼吸）噪声及各种随机噪声源，量化噪声也可视为宽带噪声。平稳的宽带噪声，通常也可以视为宽带噪声。平稳的宽带噪声，通常也可以视为高斯白噪声。

语音增强算法大致分为四种：参数法、非参数法、统计法和其它方法。

### 0.2 信号响应的意义

对于任何一个信号均可以使用冲击函数来表示，即：

$$x(n) = \sum_{m=-\infty}^{+\infty} x(m) \delta(n-m)$$

数字信号处理的意义就是通过运算来达到处理的目的，设这种运算关系为：

$T[\bullet]$ ，则输出信号  $y(n)$  和输入信号  $x(n)$  之间的关系描述为：

$$y(n) = T[x(n)]。$$

### 0.3 卷积推导

设系统输入  $x(n) = \delta(n)$ ，系统的输出  $y(n)$  的初始状态为零，这时系统输出用  $h(n)$  表示为：

$$h(n) = T[\delta(n)]。$$

则称  $h(n)$  为系统的单位脉冲响应。则对任意输入信号  $x(n)$ ，系统输出为：

$$y(n) = T \left[ \sum_{-\infty}^{+\infty} x(n) \delta(n-m) \right]$$

,根据叠加原理可得:

$$y(n) = T \left[ \sum_{-\infty}^{+\infty} x(n) \delta(n-m) \right] = \sum_{-\infty}^{+\infty} x(n) T [\delta(n-m)] \quad \text{利用系统时不变性, 可得下式:}$$

$$T [\delta(n-m)] = h(n-m), \quad \text{因此可得:}$$

$$y(n) = \sum_{-\infty}^{+\infty} x(n) \delta(n-m) = x(n) * h(n)$$

上述就是卷积公式的推导。

## 0.4 时域离散系统的输入输出描述法

描述一个系统可以不管系统内部的结构如何, 将系统看成一个黑盒子, 只描述系统的输出和输入之间的关系, 这种描述法被成为输入输出描述法。在模拟系统中使微分方程描述系统的输入和输出之间的关系, 在时域离散系统中使用差分方程描述系统的输入和输出关系。

点评: 微分方程重在描述变化的趋势, 差分方程的过程可以套用卷积的方法。

### 0.4.1 时域离散信号傅里叶变换 (DTFT, Discrete-Time Fourier Transform)

定义:

$$X(\omega) = x(e^{j\omega}) = \sum_{-\infty}^{+\infty} x(n) e^{-j\omega n}$$

上述  $\omega$  的单位是弧度, 范围是  $2\pi$ 。

其傅里叶反变换由如下公式得到:

$$x(n) = \frac{1}{2\pi} \int_{-\pi}^{\pi} X(e^{j\omega}) e^{j\omega n} d\omega。$$

### 0.4.2 周期信号由傅里叶级数表示

傅里叶变换的一些性质:

时域卷积, 频域相乘; 时域相乘, 频域卷积。

$$\sum_{-\infty}^{+\infty} |x(n)|^2 = \frac{1}{2\pi} \int_{-\pi}^{\pi} |x(n)|^2 d\omega \quad (\text{巴塞伐尔 parseval 定理}) \text{-----信号的功率也可以在频域求。}$$

## 0.5 离散傅里叶变换 (DFT)

将有限长时域离散信号变换到频域的变换, 但变换的结果是对时域离散信号的频谱的等

间隔采样。

### DFT 定义

设序列  $x(n)$  的长度为  $M$ ，定义  $x(n)$  的  $N$  点 DFT 为

$$X(k) = DFT[x(n)]_N = \sum_{n=0}^{N-1} x(n) e^{-j\frac{2\pi}{N}kn}, 0 \leq k \leq N-1$$

式中， $N$  成为离散傅里叶变换区间长度，要求  $N \geq M$ 。DTFT 中  $\omega_k = 2\pi k / N$  即可得

DFT。为书写简单，令  $W_N = e^{-j\frac{2\pi}{N}}$ ，则可以简写为：

$$X(k) = DFT[x(n)]_N = \sum_{n=0}^{N-1} x(n) W_N^{kn}, 0 \leq k \leq N-1$$

其反变换如下：

$$x(n) = IDFT[X(k)]_N = \frac{1}{N} \sum_{k=0}^{N-1} X(k) W_N^{-kn}, 0 \leq k \leq N-1$$

DFT 和 FT 之间的关系：

$$X(k) = X(e^{j\omega})|_{\omega=\frac{2\pi}{N}k}, k=0,1,\dots,N-1$$

DFT 的主要性质：

- 1) 线性性质
- 2) 隐含周期性
- 3) 循环移位性质

1) 有限长序列的循环移位

设序列  $x(n)$  的长度为  $M$ ，对  $x(n)$  以  $N(N \geq M)$  为周期进行周期延拓，得到：

$$\tilde{x}_N(n) = x((n))_N$$

定义  $x(n)$  的循环移位序列为：

$$y(n) = \tilde{x}_N(n+m) R_N(n) = x((n+m))_N R_N(n)$$

上式表示将序列  $x(n)$  以  $N$  为周期进行周期延拓，再左移  $m$  个单位取主值序列，就得到  $x(n)$  的循环移位序列  $y(n)$ 。

则有如下结论：

设序列  $x(n)$  的长度为  $M$ ，其循环移位序列为

$$y_N(n) = x((n))_N R_N(n), N \geq M$$

$$X(k) = DFT[x(n)]_N$$

则

$$Y(k) = DFT[y(n)]_N = W_{-km}^N X(k)$$

## 0.6 短时傅里叶变换 (STFT, Short-Time Fourier Transform)

DTFT 针对平稳信号的变换, 语音信号在长时间跨度上不平稳, 但其每个 10-30ms 时间段内可看成是平稳的。

定义:

$$X(n, \omega) = \sum_{m=-\infty}^{\infty} x(m) w(n-m) e^{-j\omega m}, \quad x(m) \text{ 是输入信号, } w(m) \text{ 是分析窗口,}$$

$w(n-m)$  是经过时域翻转并右移  $n$  个采样点。类似于 DFT, 离散 STFT 定义如下:

$$X(n, \omega_k) \triangleq X(n, k) = \sum_{m=-\infty}^{\infty} x(m) w(n-m) e^{-j\frac{2\pi}{N} km}, \quad \omega_k = \frac{2\pi}{N} k$$

其含义是在时域用窗函数截取信号, 对截取部分的信号进行傅里叶变换, 即在  $n$  时刻得到  $n$  时刻该段信号的傅里叶变换, 不断移动  $n$ , 即可得到不同的傅里叶变换, 将这些傅里叶变换组合起来即得  $X(n, \omega_k)$ 。

### 0.6.1 Overlap-and-Add (OA,OLA) 计算 STFT

OA 在计算  $x(n)$  和 FIR 滤波器  $h(n)$  卷积效率较高。

$$y(n) = x(n) * h(n) \stackrel{\text{def}}{=} \sum_{m=-\infty}^{\infty} h(m) x(n-m) = \sum_{m=-1}^M h(m) x(n-m), m \in [1 \quad M], h(m) \neq 0$$

OA 的基本思想是将  $x(n)$  分段, 将分段后的每段与  $h(n)$  卷积:

$$x_k(n) \stackrel{\text{def}}{=} \begin{cases} x(n+kL), n=1, 2, \dots, L \\ 0, \text{otherwise} \end{cases}, \quad L \text{ 是任意的分段长度}$$

$$x(n) = \sum_k x_k(n-kL)$$

$$y(n) = \left( \sum_k x_k(n-kL) \right) * h(n) = \sum_k (x_k(n-kL) * h(n)) = \sum_k y_k(n-kL)$$

## 0.6 FIR 数字滤波器

FIR(Finite Impulse Response Digital Filter)的最大优点是可实现线性相位滤波。

线性相位 FIRDF

设 FIRDF 的单位脉冲响应  $h(n)$  的长度为  $N$ , 则其频响应函数为:

$$H(e^{j\omega}) = \sum_0^{N-1} h(n) e^{-j\omega n}$$

将  $H(e^{j\omega})$  表示成如下形式:

$$H(e^{j\omega}) = H_g(\omega) e^{j\theta(\omega)}$$

式中,  $H_g(\omega)$  是  $\omega$  的实函数, 如果满足  $\theta(\omega) = -\omega T$  则相位满足线性关系。

线性相位对时域和频域的约束

$$H(e^{j\omega}) = \sum_0^{N-1} h(n) e^{-j\omega n} = H_g(\omega) e^{-j\omega T} \quad \text{展开可得:}$$

$$\sum_{n=0}^{N-1} h(n) (\cos(\omega n) - j \sin(\omega n)) = H_g(\omega) (\cos(\omega T) - j \sin(\omega T))$$

系数偶对称。

## 0.7 窗函数设计 FIRDF

其设计思想是使用 FIRDF 逼近希望的滤波特性。

基本方法:

(1) 构造希望逼近的频响函数  $H_d(e^{j\omega})$ 。如:

$$H_d e^{j\omega} = \begin{cases} e^{-j\omega T}, & |\omega| \leq \omega_c \\ 0, & \omega_c < |\omega| \leq \pi \end{cases}$$

(2) 求出  $h_d(n)$ 。对  $H_d(e^{j\omega})$  进行 IFT 变换:

$$h_d(n) = \frac{1}{2\pi} \int_{-\pi}^{\pi} H_d(e^{j\omega}) e^{j\omega n} d\omega = \frac{1}{2\pi} \int_{-\pi}^{\pi} e^{-j\omega T} e^{j\omega n} d\omega = \frac{\sin[\omega_c(n-T)]}{\pi(n-T)}$$

(3) 加窗得到 FIRDF 的单位脉冲响应  $h(n)$ 。

$$h(n) = h_d(n) \omega(n)$$

如果要满足线性相位, 窗函数同样需要满足线性相位。

频域自适应滤波的主要原因是, 时域自适应滤波器需要很长的脉冲响应来处理同样长的回声持续时间。

# 第一章 重采样算法

## 1.1 信号重采样

假设信号  $x(t)$  是连续时间信号, 则以  $t = nT_s$  为间隔对信号  $x(t)$  采样, 采样后信号为



$x(nT_s)$ ， $n$  是整数， $T_s$  是采样周期，根据奈奎斯特采样定理，当  $x(t)$  是带限信号，且其频带范在  $\pm F_s/2$  之间，这时采样并不会导致频谱混叠，采样率  $F_s = 1/T_s$ ，设  $x(\omega)$  是  $x(t)$  的

傅里叶变换，则有  $X(\omega) = \int_{-\infty}^{+\infty} x(t)e^{-j\omega t} dt$ 。可以假定：

$$X(\omega) = 0, |\omega| \geq \pi F_s \quad (1.1)$$

根据香农定理，根据  $x(nT_s)$  重构  $x(t)$  的公式如下：

$$x(t) \triangleq \sum_{n=-\infty}^{\infty} x(nT_s) h_s(t - nT_s) \equiv x(t) \quad (1.2)$$

此处，

$$h_s(t) \triangleq \text{sinc}(F_s t) \triangleq \frac{\sin(\pi F_s t)}{\pi F_s t} \quad (1.3)$$

以下仅讨论  $x(t)$  在  $(F'_s = 1/T'_s)$ ，仅讨论方程 1.2 是整数倍情况下的  $T'_s$ 。

当新的采样率  $F'_s$  小于原始  $F_s$  时，低通滤波截止频率必须是新采样率的二分之一，即  $\pm F'_s/2$ 。对于理想的低通滤波，则有：

$$h(t) = \min\{1, F'_s / F_s\} \text{sinc}(\min\{F_s, F'_s\}t)$$

前面的增益因子保持通带内单位增益。

Sinc 函数  $\text{sinc}(t) \triangleq \sin(\pi t) / \pi t$  的波形如下：

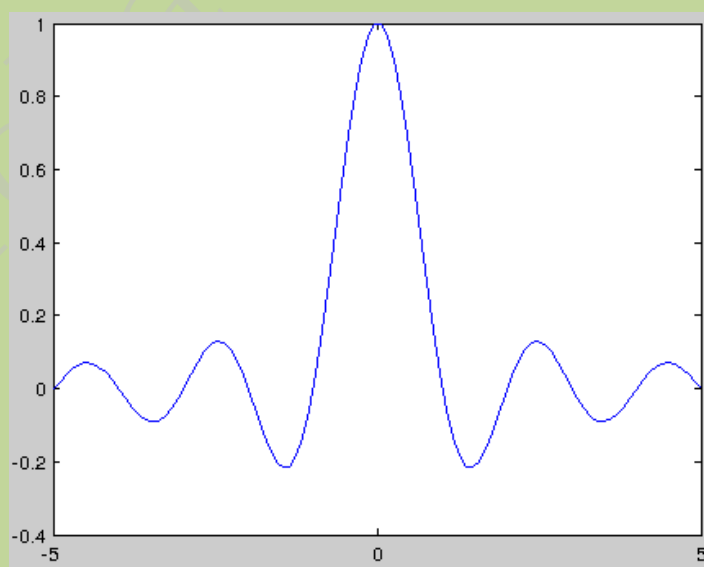


图 1.1 sinc 函数波形

公式 1.2 用卷积公式表示为： $(x * h_s(t))$ 。

该卷积过程可看成是 sinc 函数的平移叠加，每个 sinc 函数对应于一个采样点，并且幅

度被采样点调制过，所有 sinc 函数加在一起是原始信号。 $\text{sinc}(z)$  值为零的地方位于整数处（除  $z=0$ ）。这意味着  $t = nT_s$  时刻（对应于一个采样点），为和的唯一贡献是采样信号  $x(nT_s)$ 。

sinc 函数是理想低通滤波器，实际使用时会加窗处理。

## 1.2 音频重采样

在音频上，常常遇到重采样问题，如将 44.1k、48k、32k 以及 16k 音源之间的转换，通常这一过程被称为 SRC（sample rate converter）。采样率转换的基本思想是抽取和内插。从信号角度看，音频重采样是滤波。滤波函数的窗口大小以及插值函数一旦被确定，其重采样的性能也就确定了。

抽取可能引起频谱混叠，内插则产生镜频。通常在抽取前将抗混叠滤波，在内插后加抗镜像滤波。在语音识别里所需要的语音信号采样率实际上由其 ASR（automatic speech recognition）训练模型决定的。

### ➤ 重采样过程

1 设音频信号的原始采样率是  $L$ ，新的采样率是  $M$ ，原始信号长度是  $N$ ，则新采样率下信号的长度  $K$  满足如下关系：

$$K = (M / L) * N$$

2 对每个离散时间值： $k(1 \leq k \leq K)$ ，则实值  $n_k$  的值为：

$$n_k = (L / M) * k$$

$n_k$  为在原始采样间隔的情况下，要进行插值的位置。

3 确定两个加权系数的值，利用第二步计算得到的  $n_k$  值，求两个权值。选择恰当的权值，才能让线性差值所要插取的值更加接近插值点的理想幅度值，让  $w_1$  和  $w_2$  分别代表两个重采样权值。

$$w_1 = n_k - 1$$

$$w_2 = 1 - w_1$$

4 根据两个权值，估计插入点的具体幅度值。

$$y(k) = w_{1x}(n+1) + w_{2x}(n)$$

上采样会产生镜频，下采样会产生混叠。为了消除镜频和混叠，就需要将信号通过低通滤波器。上采样时，产生的镜频通过低通滤波器去除；下采样时，为了防止产生混叠，先通过低通滤波器进行处理。

### 1.3 sinc 重采样

这里之所以使用 sinc 重采样，一是因为开源语音处理算法 speex 使用的就是加窗的 sinc 方法，加窗函数是凯撒窗。

以整数因子抽取为例来说明算法实现过程，设  $x(n)$  是对模拟信号  $x_a(t)$  以奈奎斯特速率  $F_x$  采样得到的信号，其频谱为  $X(e^{j\omega})$ ，则在频率区间  $0 \leq |\omega| \leq \pi$ （模拟频率是  $|f| \leq F_x / 2 \text{Hz}$ ）， $X(e^{j\omega})$  是非零的。现在按整数倍  $D$  对  $x(n)$  进行抽取。为了避免频谱混叠，必须先对  $x(n)$  进行抗混叠低通滤波，将  $x(n)$  的有效频带限制在折叠频率（ $F_x / 2D \text{Hz}$ ）以内，等效的数字频率为  $\pi / D$  弧度以内，然后再按整数因子  $D$  对  $x(n)$  进行抽取，得到信号  $y(m)$ 。其原理框图如下：

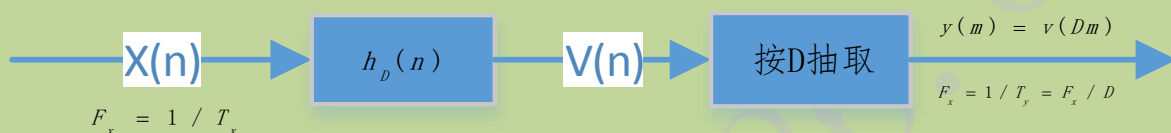


图 1.2 整数因子  $D$  抽取原理框图

理想情况下，抗混叠低通滤波器  $h_D(n)$  的频率响应：

$$H_D(e^{j\omega}) = \begin{cases} 1, & |\omega| < \pi / D \\ 0, & \pi / D \leq |\omega| \leq \pi \end{cases} \quad 1.4$$

经过抗混叠低通滤波器后输出为：

$$v(n) = h(n) * x(n) = \sum_{k=0}^{\infty} h_D(k)x(n-k) \quad 1.5$$

考虑到  $h_D(n)$  为因果稳定系统，所以式 1.5 中卷积求和从 0 开始。按整数因子  $D$  对  $v(n)$  抽取得：

$$y(m) = v(Dm) = \sum_{k=0}^{\infty} h_D(k)x(Dm-k) \quad 1.6$$

➤ 直接型 FIR 滤波器实现

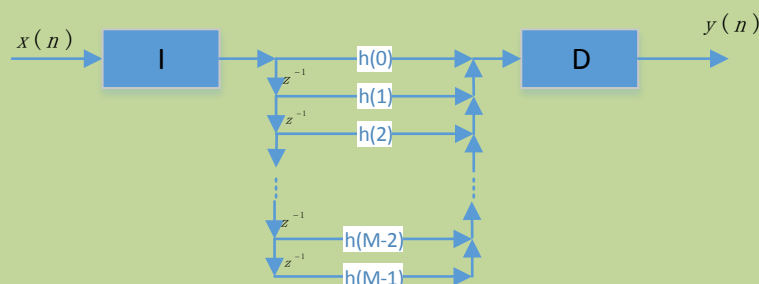


图 1.3 采样率转换系统的直接型 FIR 滤波器结构

图 1.3 中的结构很明白的表示了式 1.6 所表示的滤波过程，再经过  $D$  抽取器获得相应的

采样率，在  $I$  表示的插值过程中，会牵涉到在  $x(n)$  相邻样本值之间插入  $I-1$  个零值，如果  $I$  值比较大，则进入 FIR 滤波器的信号大部分为零，因此乘法运算的结果也大部分是零，即多数乘法是无效的，此种运算效率比较低。

#### ➤ 整数因子 $D$ 抽取直接型 FIR 滤波器结构

整数因子  $D$  抽取系统的直接型 FIR 滤波器结构如图 2.3 a 所示，该结构中，FIR 滤波器以最高采样率  $F_x$  运行，但其输出的每  $D$  个样值中抽取一个作为最终输出，丢弃  $D-1$  个样值，这样效率较低。

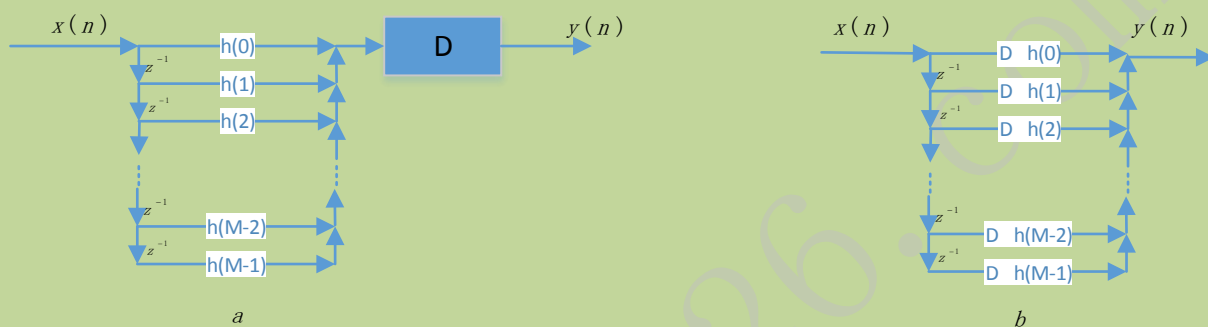


图 1.4 整数因子  $D$  抽取系统的直接型 FIR 滤波器

高效的 FIR 滤波器将抽取操作  $D$  嵌入 FIR 滤波器结构中，1.4 的 b 所示，a 的抽取器在  $n = Dm$  时刻开通，选通 FIR 滤波器的一个输出作为抽取系统输出序列的一个样值  $y(m)$ 。

$$y(m) = \sum_{k=0}^{M-1} h(k)x(Dm - k) \quad 2.4$$

b 的抽取器在  $n = Dm$  时刻同时开通，选通 FIR 滤波器输入信号  $x(n)$  的一组延迟： $x(Dm)$ ， $x(Dm-1)$ ， $x(Dm-2)$ ， $\dots$ ， $x(Dm-M+1)$ ，再进行乘法、加法运算，得到抽取系统输出序列的一个样值， $y(m) = \sum_{k=0}^{M-1} h(k)x(Dm - K)$ ，b 的运算量仅是 a 的  $1/D$  之一。

#### ➤ 整数因子 $I$ 内插系统直接型 FIR 滤波器结构

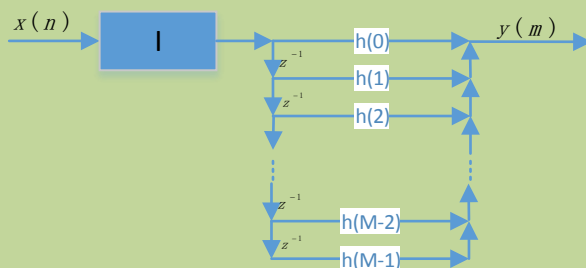


图 1.5 按整数因子  $I$  内插系统的直接型 FIR 滤波器

根据图 1.5 给出的结构，需要 FIR 滤波器以  $IF_x$  运行。该效率比较低。其高效的等效变换见图 1.6 b 所示。

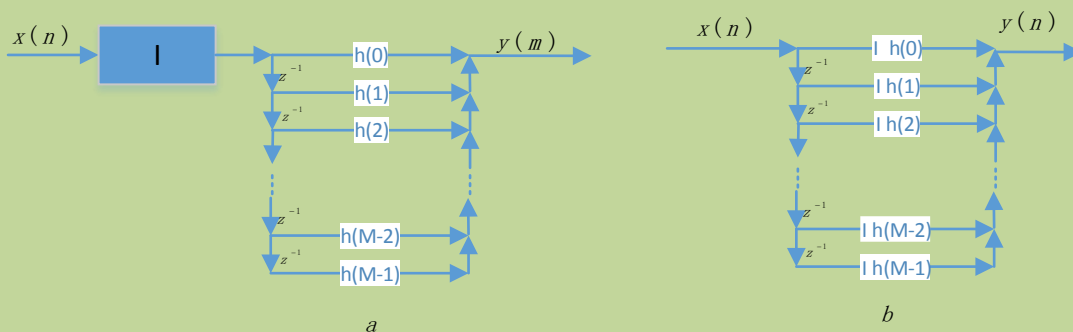


图 1.6 整数因子 I 内插系统的高效实现

### ➤ 多相滤波器结构

图 1.7 中 b 的实现结构可以采用较短的多项滤波器组来实现其内插功能，如果滤波器的总长度为  $M=NI$ ，则多相滤波器组有  $I$  个长度为  $N=M/I$  的短滤波器构成，且  $I$  个短滤波器轮流分时工作，所以称之为多相滤波器。整数因子  $I$  内插系统的直接型 FIR 滤波器的输出  $y(m) = h(m) * v(m)$ 。零值内插器的输出序列  $v(m)$  是在输入序列  $x(n)$  的两个相邻样值之间插入  $I-1$  个零样值得到，因此  $v(m)$  进入 FIR 滤波器的  $M$  个样值中只有  $N=M/I$  个非零值。即在任意时刻  $m$ ，计算  $y(m) = h(m) * v(m)$  时只有  $N$  个非零值与  $h(m)$  中的  $N$  个系数相乘。

$$v(m) = \begin{cases} x(m/I), & m=0, \pm I, \pm 2I, \pm 3I \\ 0, & \text{other} \end{cases} \quad 1.7$$

$$\text{所以 } y(m) = \sum_{n=0}^{M-1} h(n)v(m-n) = \sum_{n=0}^{N-1} h(nI)x(m-n), \quad \text{当 } m = jI + k, \quad k=0, 1, 2, \dots, I-1, j=0, 1, \dots$$

时，有：

$$y(m) = \sum_{n=0}^{M-1} h(n)v(m-n) = \sum_{n=0}^{N-1} h(k+nI)x(m-n) \quad 1.8$$

式 2.5 中的  $h(k+nI)$  看做长度  $N=M/I$  的子滤波器的单位脉冲响应，并用  $p_k(n)$  表示，则：

$$p_k(n) = h(k+nI), k=0, 1, \dots, I-1, n=0, 1, \dots, N-1$$

这样，从  $m=0$  开始，整数因子  $I$  内插系统的输出序列  $y(m)$  计算如下：

$$y(m) = \sum_{n=0}^{M-1} h(n)v(m-n) = p_k(n) * x(n) \quad 1.9$$

当  $m = jI + k$  从 0 开始增大时， $k$  从 0 开始以  $I$  为周期循环取值； $j$  表示循环周期数。所以式 2.6 对应的多相滤波器结构如图 2.6 所示。输出序列  $y(m)$  就是从  $k=0$  开始，依次循环选取  $I$  个子滤波器的输出所形成的序列。

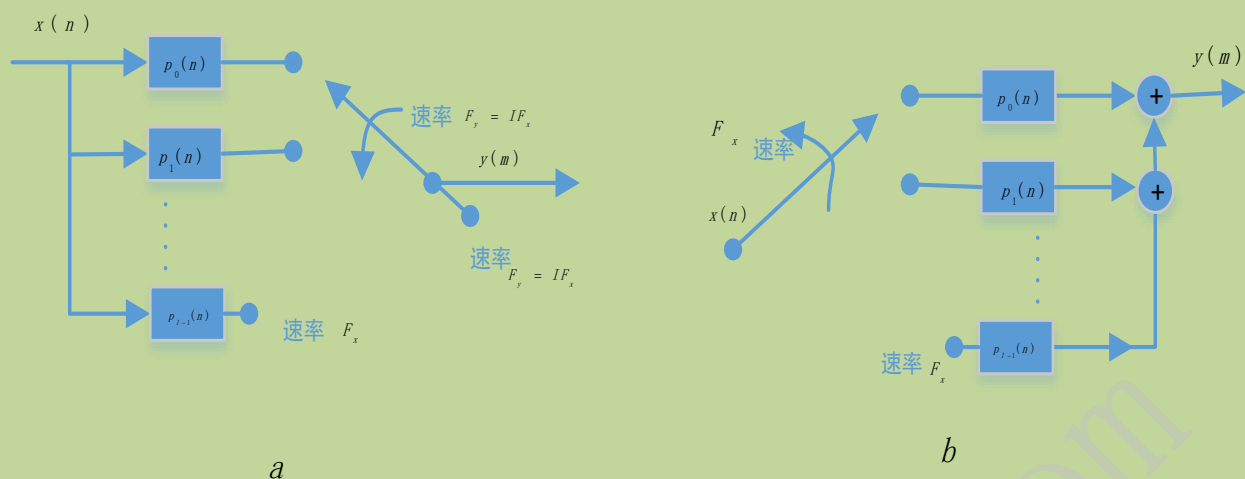


图 1.7 多相滤波结构，（a 整数因子 I 内插，b 整数因子 D 抽取）

多相滤波器的单位脉冲响应为：

$$p_k(n) = h(k + nD), k = 0, 1, 2, \dots, D-1; n = 0, 1, \dots, N-1 \quad 1.10$$

式中， $N$  为  $p_k(n)$  的长度。一般选择原抗混叠 FIR 滤波器总长度  $M=DN$ ， $N=M/D$ 。电子开关以速率  $F_x$  逆时针旋转，从子滤波器  $p_0(n)$  在  $m=0$  时刻开始，并输出  $y(0)$ ；然后以电子开关速率  $F_x$  逆时针旋转一周，即每次转到子滤波器  $p_0(n)$  时，输出端就以速率  $F_y = F_x / D$  送出一个  $y(m)$  样值。

从频域角度看，线性重采样适合应用与下采样。由于自己所在的工程所需的是将 48K 的音频信号重采样到 16k，这就是一个抽取过程。

## 1.4 MATLAB 实现

前述的内容属于理论方面，现在基于 MATLAB 来阐述重采样实现的一些工程实现细节。首先使用 MATLAB 的 decimate 函数将 48K 的音源抽取到 16k，并且绘制它们的频谱如下：

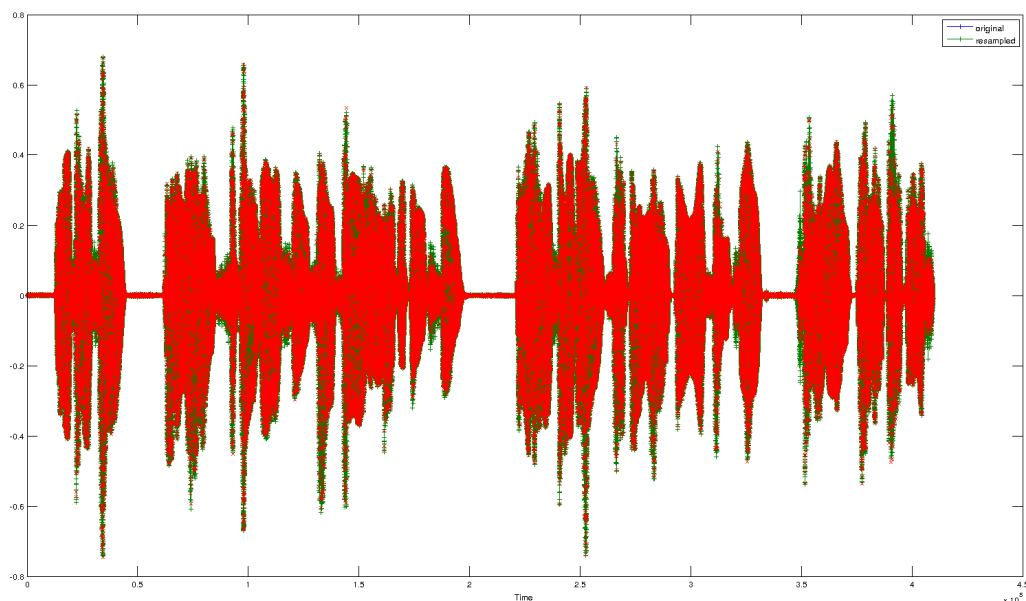


图 1.8 resample 抽取前后语音信号的时频域波形

上图深绿色的点数是  $4096 \times 100$ ，蓝色的点数是 136500； $4096 \times 100 / 136500 = 3.0007$ ，由时域可以看到信号确实数量降低到原来的  $1/3$  了。两种颜色叠加之后的波形如上，从时域波形上来看，两者的波形并未完全重叠，和其它地方相比，尖端处的差异大些。

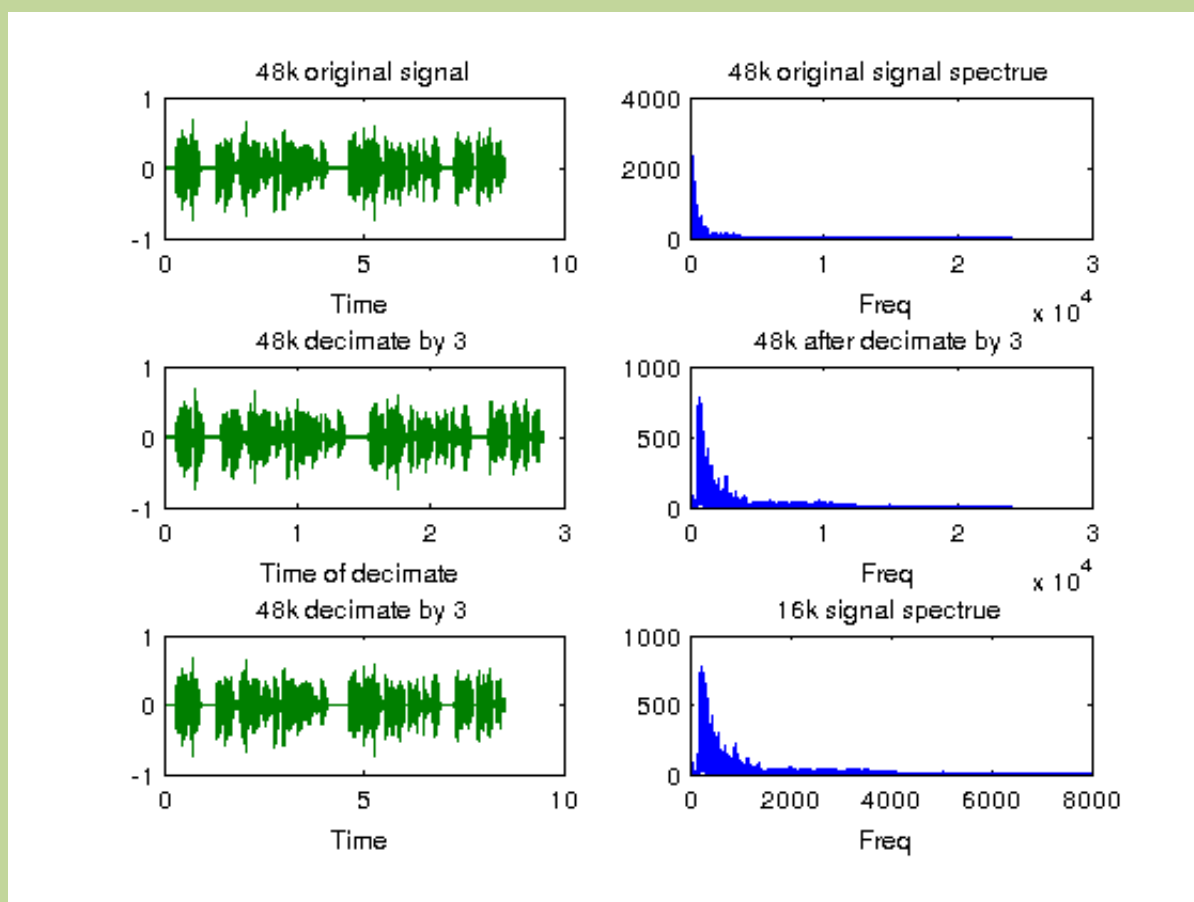


图 1.9 48k 和 16k 采样率下的信号时域和频域波形

图 1.9 的第一行是原始 48k 信号的时域和频域波形图，其频域波形在 24k 截止，第二行是 48K 采样率下抽取后信号的时域和频域波形图，对比时域和频域来看，它们的频谱区别还是有的，第三行是 16k 采样率下绘制重采样后信号的视频域波形，其频域截止频率在 8K，和第二行信号的频谱非常接近，但是它们还是存在差异的。

```
clear; close all; clc;
%read sound wave
[x,fs]=audioread('48k_sound.wav');
% sound(x,48000);
%sound(y,fs,nbits);
N=length(x);
tx=(0:N-1)/fs;%calc the time that correspond to original signal
xf=fft(x);
fx=(0:N/2)*fs/N;
figure(1);
subplot(321);plot(tx,x);xlabel('Time');title('48k original signal');
subplot(322);plot(fx,abs(xf(1:N/2+1)));xlabel('Freq ');title('48k original signal spectrue');
audiowrite('48k_sound.wav',x,48000);

%decimate signal
k=1/N/3;% Number to decimate 48k to 16k

%filter design kasier
fp=8000;fc=12000;as=100;ap=1;FS=48000;
wc=2*fc/FS; wp=2*fp/FS;% kaiser window
```



```

N=ceil((as-7.95)/(14.36*(wc-wp)/2))+1;
beta=0.1102*(as-8.7);
win=kaiser(N+1, beta);
wvtool(win);
B=fir1(N, wc, win);
freqz(B,1,512,fs);
y=filter(B,1,x);
y=downsample(y,3);
% 实际上这里是每隔三个点取一个点，即先低通滤波再取点，和紧接着的下面的三行结果是一样的
% y=x(3*k);% 48k-->16k,and data store in y
% y=decimate(x,1,3);
% y=resample(x,1,3);
M=length(y);% M is the length of y
% analysis y under the frequency of 48k
ty=(0:M-1)/fs; % fs equ 48k
subplot(323);plot(ty,y);xlabel('Time of decimate');title('48k decimate by 3');
yf=fft(y);% fft of 48k downsample by 3
fy=(0:M/2)*fs/M;
subplot(324);plot(fy,abs(yf(1:M/2+1)));xlabel('Freq');title('48k after decimate by 3');

%decimate data under sample 16k
tz=(0:M-1)/(fs/3);
subplot(325);plot(tz,y);xlabel('Time');title('48k decimate by 3');
fz=(0:M/2)*(fs/3)/M;
subplot(326);plot(fz,abs(yf(1:M/2+1)));xlabel('Freq ');title('16k signal spectrue');

sound(y,16000);

audiowrite('16k_sound.wav',y,16000);

% %ellipord
% fp=8000;fc=12000;as=100;ap=1;fs=48000;
% wc=2*fc/fs;wp=2*fp/fs;
% [n,wn]=ellipord(wp,wc,ap,as);%
% [b,a]=ellip(n,ap,as,wn);
% freqz(b,a,512,fs);

```

从上面的 MATLAB 运行的结果来看，从 48K 数据重采样到 16K，每隔三个点取一个点和先经过低通滤波再抽取实际上几乎没有差别。上面 MATLAB 代码经过滤波后的各种波形和图 1.9 几乎没有差别。

## 1.5 开源重采样 C 语言实现

## 1.6 基于 FPGA 的多相滤波

## 1.7 重采样性能评估

音频信号的好坏常用信噪比来评估，同样可以根据重采样后音频信号的性噪比来评估重采样的质量。信噪比的定义如下：

$$SNR = 10 \log_{10} \left\{ \frac{\sum_{n=0}^M s^2(n)}{\sum_{n=0}^M (s(n) - \hat{s}(n))^2} \right\} \quad 1.11$$



式中  $s(n)$  为原始信号,  $\hat{s}(n)$  是转换后的同频率音频信号,  $M$  为音频信号的长度。信噪比越大意味着总体上原始音频信号和重采样后的信号之间的差距较小, 接近程度较高, 在功放的测试中就有信噪比这么一个测试项, 现在信噪比在 80dB 以上的音响比较常见了。

由于语音信号是短时平稳随机信号, 可将语音信号分帧计算各帧的信噪比, 平均后可以得到信噪比。计算的公式如下:

$$SNR = \frac{1}{N} \sum_{m=0}^{M-1} \log_{10} \left\{ \frac{\sum_{n=0}^r x_n^2}{\sum_{n=0}^r (x_n - \hat{x}_n)^2} \right\} \quad 1.12$$

$M$  是总的帧数,  $r$  是每一帧的长度。

## 第 2 章 回声消除 (AEC) 原理及实现

### 2.1 回声消除原理

回声消除的基本原理是使用一个自适应滤波器对未知的回声信道  $\omega$  进行参数辨识, 根据扬声器信号与产生的多路回声的相关性为基础, 建立远端信号模型, 模拟回声路径, 通过自适应算法调整, 使其冲击响应和真实回声路径相逼近。然后将麦克风接收到的信号减去估计值, 即可实现回声消除功能。

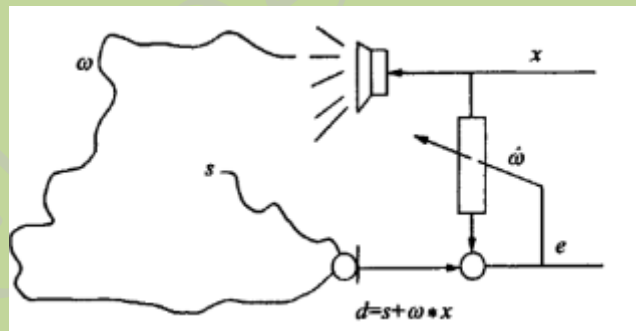


图 2.1 回声消除原理图

式中,  $\omega$  是回声通道的时域冲击响应函数;  $x$  是远端语音;  $echo$  是所得回声;  $s$  是近端说话人语音,  $d$  为麦克风采集到的信号;  $\hat{y}$  为对回声信号的估计值;  $e$  为误差。在电话、视频会议中这里的  $x$  通信另一端的语音信号, 而在机器语音识别中, 这里的  $x$  则指机器自身发出的声音。

$$echo = x * \omega \quad 2.1.1$$

$$d = s + echo \quad 2.1.2$$

$$\hat{y} = x * \hat{\omega} \quad 2.1.3$$

$$e = d - \hat{y} \quad 2.1.4$$

为了消除较长时间的回声，需要 FIR 滤波器的阶数尽量大。时域计算诸多不便，使用频域分块自适应滤波算法。

## 2.2 维纳滤波

均方误差（MSE， Mean Square Error）,对于离散时间系统，可定义期望响应  $d_k$  为一个希望自适应系统的输出  $y_k$  与之相接近的信号， $k$  为采样时刻。

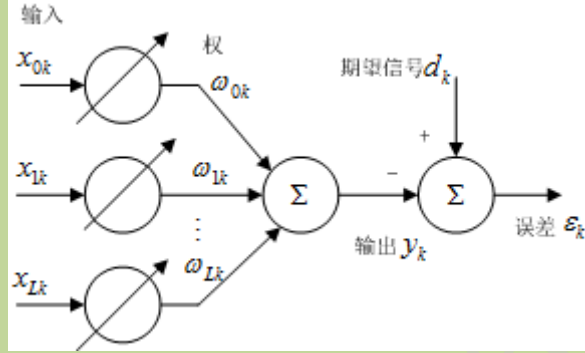


图 2.2 MSE 自适应系统

根据图 2.2，可以求得误差信号：

$$\varepsilon_k = d_k - y_k \quad 2.2.1$$

自适应线性组合器输出：

$$y_k = W_k^T X_k \quad 2.2.2$$

其中：

$$X_k = [x_{0k} x_{1k} \dots x_{Lk}], W_k = [w_{0k} w_{1k} \dots w_{Lk}]^T$$

分别为自适应系统在  $k$  时刻的输入信号向量和权向量，系统的均方误差为：

$$E(|\varepsilon_k|^2) = E[(d_k - y_k)^*(d_k - y_k)] = E(|d_k|^2) + W_k^H E[X_k^* X_k^T] W_k - 2 \operatorname{Re} \{W_k^T E[d_k^* X_k]\} \quad 2.2.3$$

$$R = E[X_k^* X_k^T] = \begin{bmatrix} E[x_{0k}^* x_{0k}] & E[x_{0k}^* x_{1k}] & \dots & E[x_{0k}^* x_{Lk}] \\ E[x_{1k}^* x_{0k}] & E[x_{1k}^* x_{1k}] & \dots & E[x_{1k}^* x_{Lk}] \\ \vdots & \vdots & \ddots & \vdots \\ E[x_{Lk}^* x_{0k}] & E[x_{Lk}^* x_{1k}] & \dots & E[x_{Lk}^* x_{Lk}] \end{bmatrix} \quad 2.2.4$$

定义期望响应和输入信号之间的互相关向量为：

$$P = E[d_k^* X_k] = \begin{bmatrix} d_k^* x_{0k} \\ \vdots \\ d_k^* x_{Lk} \end{bmatrix} \quad 2.2.5$$

将式 2.2.3 简化成下式：

$$\xi(\mathbf{w}) = E(|d_k|^2) + \mathbf{w}_k^H \mathbf{R} \mathbf{w}_k - 2 \operatorname{Re} \{ \mathbf{w}_k^T \mathbf{P} \} \quad 2.2.5$$

理想情况下  $E(|\varepsilon_k|^2)$  等于零，这时估计值等于观测值，如果不能达到理想，则应该是越小越好，这样估计值和观测值最接近。

对 2.2.5 求偏导数，得：

$$\nabla = \frac{\partial}{\partial \mathbf{W}} [\xi(\mathbf{W})] = 2\mathbf{R}\mathbf{W} - 2\mathbf{P}^* \quad 2.2.6$$

最佳权向量处的梯度值为零，于是：

$$\nabla = 2\mathbf{R}\mathbf{W}_{opt} - 2\mathbf{P}^* = 0 \quad 2.2.7$$

最小均方误差输出情况下的最佳权向量  $\mathbf{W}_{opt}$  满足维纳-霍夫方程：

$$\mathbf{W}_{opt} = \mathbf{R}^{-1} \mathbf{P}^* \quad 2.2.8$$

## 2.3 LMS 算法

$$\varepsilon_k = d_k - \mathbf{X}_k^T \mathbf{W}_k \quad 2.3.1$$

式中， $\mathbf{X}_k$  为输入样本向量，使用单次采样数据  $|\varepsilon_k|^2$  来代替均方误差  $\xi_k$ ，这样其梯度估计可表示为如下形式：

$$\begin{aligned} \hat{\nabla}_k &= \frac{\partial}{\partial \mathbf{W}_k} |\varepsilon_k|^2 = \frac{\partial}{\partial \mathbf{W}_k} [ |d_k|^2 + \mathbf{W}_k^H \mathbf{X}_k^* \mathbf{X}_k^T \mathbf{W}_k - 2 \operatorname{Re} (d_k^* \mathbf{X}_k^T \mathbf{W}_k) ] \\ &= 2\mathbf{X}_k^* \mathbf{X}_k^T \mathbf{W}_k - 2d \mathbf{X}_k^* = -2\varepsilon_k \mathbf{X}_k^* \end{aligned} \quad 2.3.2$$

基于最速下降法的权向量迭代如下：

$$\mathbf{W}_{k+1} = \mathbf{W}_k - \mu \hat{\nabla}_k = \mathbf{W}_k + 2\mu \varepsilon_k \mathbf{X}_k^* \quad 2.3.4$$

其中  $\mu$  是步长因子， $0 < \mu < \frac{1}{\lambda_{\max}}$ ， $\lambda_{\max}$  是  $\mathbf{R}_{xx}$  的最大特征值。 $\mathbf{W}(k)$  收敛于  $\mathbf{W}_{opt}$  由比值

$d = \frac{\lambda_{\max}}{\lambda_{\min}}$  决定，该比值叫做谱动态范围。大的  $d$  值意味着较长的时间才能收敛到最佳权值。

该算法用在语音增强的加性噪声消除功能上时，其工程实践并不完全按照式 2.3.1 意义来实现。

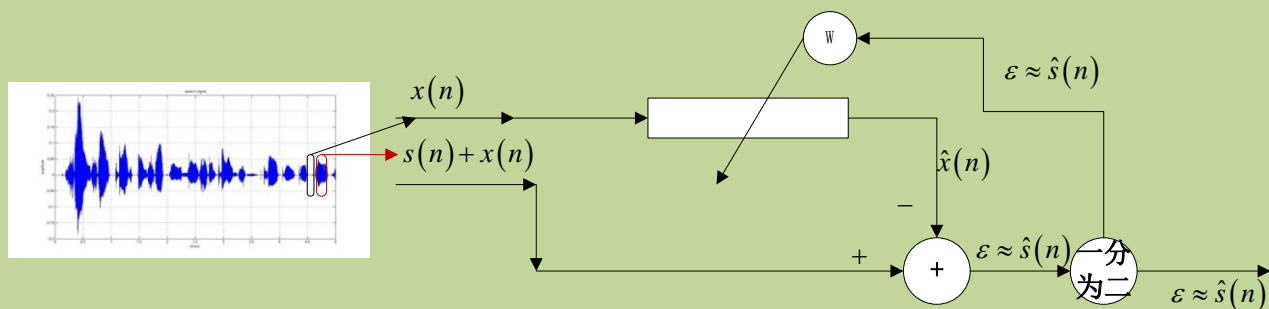


图 LMS 算法在语音增强中使用方法

在语音增强中，其目的是获得纯净的语音信号  $s(n)$ ，即上图中的最后输出信号，输入信号有两种，一种是带噪的语音信号  $s(n) + x(n)$ ，另一种是只有噪声的输入  $x(n)$ ，在没有人说话的情况下的输入信号，就仅仅是噪声输入。这里要使得噪声估计  $\hat{x}(n)$  非常接近  $x(n)$ ，这样  $\varepsilon = \sum_n s(n) + x(n) - \hat{x}(n)$ ，这时如果  $\varepsilon^2$  最小，则可以知道估计出的  $\varepsilon$  最接近  $s(n)$ 。

上述过程可以概述如下：

- 首先获取到噪声输入  $x(n)$ ，并存储下来，以 64 或者 128 点为总长度不断刷新存储噪声输入。
- 采集带噪声的语音信号  $s(n) + x(n)$ 。
- 用采集带噪语音信号减去估计到的噪声信号  $s(n) + x(n) - \hat{x}(n)$ 。
- 用 C 的输出作为误差，调节噪声权向量  $W$ 。

MATLAB 实现具体包括如下三个部分：

```
% Loop over input vector
for ii = 1:length(signal_with_noise)
    % Update buffer
    noise_buf = obj.update_buf(noise_buf, noise(ii)); //输入噪声估计
    % Filter this sample with current coefficient values
    filter_output = obj.data_filter(coefs, noise_buf); //通过权向量估计  $\hat{x}(n)$ 
    % Compute error
    err = signal_with_noise(ii) - filter_output; //相减得到  $\hat{s}(n)$ 
    % Update coefficients
    coefs = obj.update_coefs(coefs, noise_buf, obj.filter_params.step_size,
obj.filter_params.leakage, err); //用  $\hat{s}(n)$  调节权向量
    % Build output vector
```

dout(ii) = err; //存储输入信号的估计值  $\hat{s}(n)$

### 2.3.1 NLMS

输入信号较大时，会遇到梯度噪声放大的问题，使得能量低的信号算法收敛速度较慢。将输入信号按照自身的平均能量进行归一化处理，即得到归一化 NLMS 算法。设输入带噪信号可表示为： $x(n)$ ，其迭代算法的 NLMS 表示公式如下：

$$\mathbf{W}_{n+1} = \mathbf{W}_n - \mu \hat{\mathbf{V}}_n = \mathbf{W}_n + \frac{\mu}{N} \frac{e(n)\mathbf{x}(n)}{\hat{\sigma}_x^2(n)}$$

其中  $\hat{\sigma}_x^2(n) = \frac{1}{N} \sum_{k=0}^{N-1} x^2(n-k)$ ，其中 N 是噪声消除器和回波抵消器的长度，（常取 512

或者 1024）； $\mu$  是步长因子。当  $\hat{\sigma}_x^2(k)$  较小时， $\frac{\mu}{\hat{\sigma}_x^2(n)}$  的值可能较大，这时迭代算法变成如下形式：

$$\mathbf{W}_{n+1} = \mathbf{W}_n - \mu \hat{\mathbf{V}}_n = \mathbf{W}_n + \frac{\mu}{N} \frac{e(n)\mathbf{x}(n)}{\sigma + \hat{\sigma}_x^2(n)}$$

其计算过程如下：

参数：M=抽头系数（即 FIR 滤波器长度）

$\mu$  自适应常数

$$0 < \mu < 2 \frac{E[|\mathbf{x}(n)|^2] E[|\boldsymbol{\varepsilon}(n)|^2]}{E[|\mathbf{e}(n)|^2]}, \text{ 其中 } E[|\boldsymbol{\varepsilon}(n)|^2] = E[|\mathbf{W}_{opt} - \hat{\mathbf{W}}(n)|^2], \text{ 是权向量均方偏}$$

差， $\mathbf{W}_{opt}$  是最优维纳解， $\hat{\mathbf{W}}(n)$  是第 n 次迭代中得到的估值。 $E[|\mathbf{x}(n)|^2]$  是带噪输入信号的功率， $E[|\mathbf{e}(n)|^2]$  是误差信号功率。

初始化：

如果知道抽头权向量  $\hat{\mathbf{W}}(n)$  的先验知识，则用其来初始化  $\hat{\mathbf{W}}(0)$ ，否则令  $\hat{\mathbf{W}}(0) = \mathbf{0}$

数据：

A) 给定的  $\mathbf{X}(n)$  = 第 n 时刻  $M \times 1$  抽头输入向量，

$d(n)$  = 第 n 时间步的期望响应

B) 要计算的： $\hat{\mathbf{W}}(n+1)$  = 第 n+1 步抽头权向量的估计

计算：

对  $n=0,1,2, \dots$  计算

$$e(n) = d(n) - \hat{W}^H(n) X(n)$$

$$W_{n+1} = W_n - \mu \hat{\nabla}_n = W_n + \frac{\mu}{N} \frac{e(n) X(n)}{\hat{\sigma}_x^2(n)}$$

### 2.3.2 SE-LMS

Signed-error LMS (SE-LMS) 算法将误差  $e(n)$  用  $[-1 \ 0 \ 1]$  这三个量化值来代替, 如果误差大于 0, 则将  $e(n)$  赋值为 1, 其它类推。这时式 2.3.4 退化如下形式:

$$W_{k+1} = W_k - \mu \hat{\nabla}_k = W_k + 2\mu X_k^* \text{sign}(e[n])$$

该算法在加快运算速度的同时简化了电路结构,  $\mu$  设置成 2 的指数时, 通过移位就可以实现这里的乘法操作。降低了硬件实现的复杂度。

### 2.3.2 SD-LMS

Signal-dependent LMS 算法和 SE-LMS 很相似, 误差也是只取  $[-1 \ 0 \ 1]$  这三个值, 不同的是, 其选择是以采样到的误差信号为参考的, 如果  $x(n)$  大于 0, 则  $e(n)$  用 1 代替, 依次类推。

$$W_{k+1} = W_k - \mu \hat{\nabla}_k = W_k + 2\mu \text{sign}(x[n])$$

### 2.3.2 SS-LMS

Sign-sign LMS 算法既考虑输入信号的符号又考虑误差的信号。

$$W_{k+1} = W_k - \mu \hat{\nabla}_k = W_k + 2\mu \text{sign}(x[n]) \text{sign}(e[n])$$

### 2.3.2 LLMS

Leaky LMS 算法减轻了系数溢出问题。其不仅考虑了均方误差  $e^2(n)$  也考虑了滤波器系数。其权向量更新方程如下:

$$W_{k+1} = (1 - \mu\alpha)W_k - \mu \hat{\nabla}_k = W_k + 2\mu e(n) X(n)$$

### 2.3.2 LNLMS

Leaky NLMS 是归一化的 LLMS 算法。

$$W_{n+1} = (1 - \mu\alpha)W_n - \mu \hat{\nabla}_n = (1 - \mu\alpha)W_n + \frac{\mu}{N} \frac{e(n) X(n)}{\hat{\sigma}_x^2(n)}$$

## 2.4 块自适应滤波

### 2.4.1 块自适应滤波器

计算过程如下，对参考信号  $x$  分段并做 FFT 变换，分别对各段数据做频域滤波，累加后做 FFT 反变换，并只取后  $L$  ( $L$  是原始信号的分段后的长度) 点为有效的线性卷积结果，得到的是估计信号，将估计信号从回声信号中去除，得残差信号。计算子带步长，调整各段滤波器系数。这一过程表示如下图。

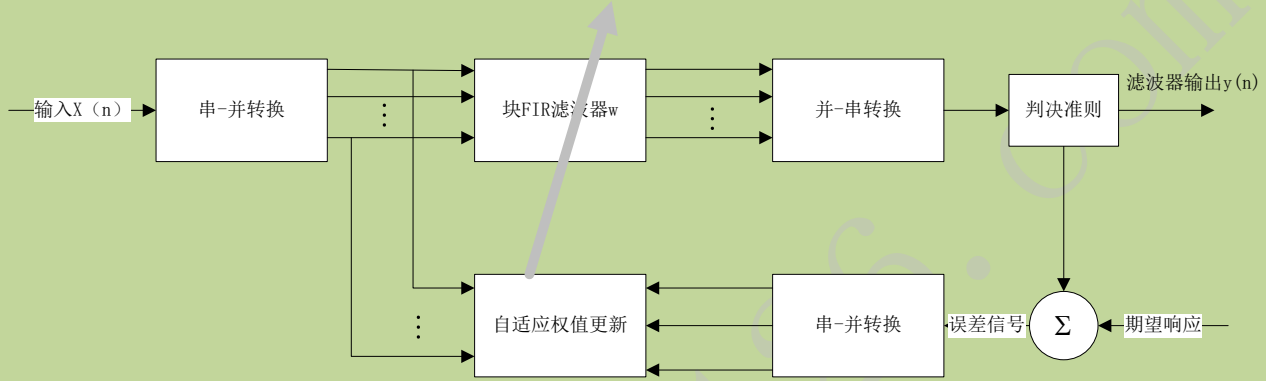


图 2.2 块自适应滤波器

设  $n$  时刻输入序列  $x(n)$  如下：

$$X(n) = [x(n), x(n-1), \dots, x(n-M+1)]^T \quad 2.4.1$$

对应于长度为  $M$  的 FIR 滤波器在  $n$  时刻的抽头权向量为：

$$\hat{W}(n) = [\hat{w}_0(n), \hat{w}_1(n-1), \dots, \hat{w}_{M-1}(n)]^T \quad 2.4.2$$

根据 FIR 滤波器原理：

$$y(n) = x(n) \times \hat{w}_0(n) + x(n-1) \times \hat{w}_1(n) + \dots + x(n-M+1) \times \hat{w}_{M-1}(n) \quad 2.4.3$$

用向量可以表示成如下：

$$y(n) = X(n)^T \hat{W}(n) \quad 2.4.4$$

下面对  $x(n)$  进行分块，设  $k$  表示块下标，它与原始样值时间  $n$  的关系为：

$$n = kL + i, i = 0, 1, \dots, L-1; k = 1, 2, \dots \quad 2.4.5$$

其中  $L$  是块的长度。第  $k$  块的数据为  $\{X(kL+i)\}_{i=0}^{L-1}$ ，其矩阵表示形式如下：

$$A^T(k) = [x(kL), x(kL+1), \dots, x(kL+L-1)] \quad 2.4.6$$

将滤波器对输入块  $A(k)$  的响应表示如下：

$$y(kL+i) = \hat{W}^T(k) A(k) = \sum_{j=0}^{M-1} \hat{w}_j(k) x(kL+i-j), i = 0, 1, \dots, L-1 \quad 2.4.7$$

设  $d(kL+i)$  表示期望信号，误差信号表示如下：

$$e(kL+i) = d(kL+i) - y(kL+i) \quad 2.4.8$$

考虑滤波器长度  $M=3$ ，块长度  $L=3$ ，其三个相邻的数据块是  $k-1$ ， $k$ ， $k+1$ ，则

$k-1$  滤波结果如下：

$$\begin{bmatrix} y(3k-3) \\ y(3k-2) \\ y(3k-1) \end{bmatrix} = \begin{bmatrix} x(3k-3) & x(3k-4) & x(3k-5) \\ x(3k-2) & x(3k-3) & x(3k-4) \\ x(3k-1) & x(3k-2) & x(3k-3) \end{bmatrix} \begin{bmatrix} w_0(k-1) \\ w_1(k-1) \\ w_2(k-1) \end{bmatrix} \quad 2.4.9$$

$k$  块滤波结果如下：

$$\begin{bmatrix} y(3k) \\ y(3k+1) \\ y(3k+2) \end{bmatrix} = \begin{bmatrix} x(3k) & x(3k-1) & x(3k-2) \\ x(3k+1) & x(3k) & x(3k-1) \\ x(3k+2) & x(3k+1) & x(3k) \end{bmatrix} \begin{bmatrix} w_0(k) \\ w_1(k) \\ w_2(k) \end{bmatrix} \quad 2.4.10$$

$k+1$  块滤波器结果如下：

$$\begin{bmatrix} y(3k+3) \\ y(3k+4) \\ y(3k+5) \end{bmatrix} = \begin{bmatrix} x(3k+3) & x(3k+2) & x(3k+1) \\ x(3k+4) & x(3k+3) & x(3k+2) \\ x(3k+5) & x(3k+4) & x(3k+3) \end{bmatrix} \begin{bmatrix} w_0(k+1) \\ w_1(k+1) \\ w_2(k+1) \end{bmatrix} \quad 2.4.11$$

上面的数据矩阵是托伯利兹矩阵，主对角线元素都相同。

## 2.4.2 块 LMS

权向量调整公式如下：

（权向量的调整）=（步长参数）\*（抽头输入向量）\*（误差信号）

因为在块 LMS 算法中误差信号随抽样速率而变，对于每一个数据块，我们有不同的用于自适应过程的误差信号值。因此，每一个块的抽头权向量更新公式如下：

$$\hat{w}(k+1) = \hat{w}(k) + \mu \sum_{i=0}^{L-1} x(kL+i)e(kL+i) \quad 2.4.12$$

其梯度向量的估计如下：

$$\hat{\nabla}(k) = -\frac{2}{L} \sum_{i=0}^{L-1} x(kL+i)e(kL+i) \quad 2.4.13$$

$\hat{\nabla}(k)$  的无偏估计如下：

$$\hat{w}(k+1) = \hat{w}(k) - \frac{1}{2} \mu_B \hat{\nabla}(k) \quad 2.4.14$$

## 2.4.3 块 LMS 算法收敛性

由于时间平均的缘故，它具有估计精度随快长度增加而大幅提高的特性。然而，长度的增加会导致其收敛速度进一步减慢。后文的快速 LMS 算法加速了这一过程。



### ➤ 平均时长数

$$\tau_{mse,av} = \frac{1}{2\mu_B \lambda_{av}} \quad 2.4.15$$

其中  $\lambda_{av}$  是输入自相关矩阵  $R = E[x(n)x(n)^T]$

上式中为了使零阶公式成立,  $\mu_B$  必须小于  $1/\lambda_{\max}$ , 其中  $\lambda_{\max}$  是相关矩阵的最大特征值。

### ➤ 失调

$$t \nu = \frac{\mu_B}{2L} tr[R] \quad 2.4.16$$

$tr[B]$  是相关矩阵的迹。

## 2.4.4 块长的选择

设滤波器长度  $M$  和块长度  $L$  的关系有三种可能:

1.  $L=M$ , 从计算的复杂性上看, 最佳
2.  $L<M$ , 有降低延迟的好处。
3.  $L>M$ , 将产生自适应过程冗余运算。

## 2.5 FLMS

FLMS (Fast LMS) 的基本思想是将时域块 LMS 放到频域来计算。利用 FFT 算法在频域上完成滤波器系数的自适应。快速卷积算法用重叠相加法和重叠存储法。重叠相加法是将长序列分成大小相等的短片段, 分别对各个端片段做 FFT 变换, 再将变换重叠的部分相加构成最终 FFT 结果, 重叠存储法在分段时, 各个短的段之间存在重叠, 对各个段进行 FFT 变换, 最后将 FFT 变换得结果直接相加即得最终变换结果。当块的大小和权值个数相等时, 运算效率达到最高。

根据重叠存储方法, 将滤波器  $M$  个抽头权值用等个数的零来填补, 并采用  $N$  点 FFT 进行计算, 其中  $N=2M$ , 因此,  $N \times 1$  的向量:

$$\hat{W}(k) = FFT \begin{bmatrix} \hat{w}(k) \\ 0 \end{bmatrix} \quad 2.5.1$$

表示 FFT 补零后的系数, 抽头权向量为  $\hat{w}(k)$ 。值得注意的是频域权向量  $\hat{W}(k)$  的长度是时域权向量  $\hat{w}(k)$  长度的两倍。相应的令:

$$X(k) = \text{diag} \left\{ \text{FFT} \begin{bmatrix} x(kM-M), \dots, x(kM-1), & x(kM), \dots, x(kM+M-1), \\ K-1, \text{block} & K, \text{block} \end{bmatrix} \right\} \quad 2.5.2$$

表示对输入数据的两个相继子块进行傅里叶变换得到一个  $N \times N$  对角阵。

将重叠存储法应用于 2.4.7 得。

$$\begin{aligned} y^T(k) &= [y(kM), y(kM+1), \dots, y(kM+M-1)] \\ &= \text{IFFT} [X(k) \hat{W}(k)], \text{last} M \end{aligned} \quad 2.5.3$$

是 2.5.3 只有最后  $M$  个元素被保留，因为最前面的  $M$  个元素是循环卷积的结果。

设第  $K$  块的  $M \times 1$  期望响应和误差信号分别如下：

$$d(k) = [d(kM), d(kM+1), \dots, d(kM+M-1)]^T \quad 2.5.4$$

$$e(k) = [e(kM), e(kM+1), \dots, e(kM+M-1)]^T = d(k) - y(k) \quad 2.5.5$$

根据式 2.5.3，可将  $e(k)$  变换到频域，即

$$E(k) = \text{FFT} \begin{bmatrix} \mathbf{0} \\ e(k) \end{bmatrix} \quad 2.5.6$$

则在更新权值的相关矩阵如下：

$$\Phi(k) = \sum_{i=0}^{L-1} x(kL+i)e(kL+i) = \text{IFFT} [X^T(k)E(k)] \quad , \text{ 的最前面 } M \text{ 个元素} \quad 2.5.7$$

则抽头的更新过程在频域中的表现如下：

$$\hat{W}(k+1) = \hat{W}(k) + \mu \text{FFT} \begin{bmatrix} \Phi(k) \\ \mathbf{0} \end{bmatrix} \quad 2.5.8$$

## 2.6 MDF 自适应权值调整

### 时域解

对于  $N$  阶 NLMS 算法，其误差调节向量如下式：

$$e(n) = d(n) - \hat{y}(n) = d(n) - \sum_{k=0}^{N-1} \hat{w}_k(n)x(n-k) \quad 2.6.1$$

权值更新如下：

$$\begin{aligned}\hat{w}_k(n+1) &= \hat{w}_k(n) + \mu \frac{e(n)x^*(n-k)}{\sum_{i=0}^{N-1}|x(n-i)|^2} \\ &= \hat{w}_k(n) + \mu \frac{(d(n) - \sum_i \hat{w}_i(n)x(n-i))x^*(n-k)}{\sum_{i=0}^{N-1}|x(n-i)|^2}\end{aligned}\quad 2.6.2$$

其中  $x(n)$  是参考信号， $\hat{w}_k(n)$  是  $n$  时刻和步长  $\mu$  的权值更新。假设滤波后的误差为  $\sigma_k(n) = \hat{w}_k(n) - w_k(n)$ ， $d(n) = v(n) + \sum_k \hat{w}_k(n)x(n-k)$ ，则误差的迭代关系如下：

$$\delta_k(n+1) = \delta_k(n) + \mu \frac{(v(n) - \sum_i \delta_i(n)x(n-i))x^*(n-k)}{\sum_{i=0}^{N-1}|x(n-i)|^2} \quad 2.6.3$$

在每一次调节中，滤波器的误差估计为  $\Lambda(n) = \sum_k \delta_k^*(n)\delta_k(n)$ ，展开后得如下形式：

$$\Lambda(n+1) = \sum_{k=0}^{N-1} \left| \delta_k(n) + \mu \frac{(v(n) - \sum_i \delta_i(n)x(n-i))x^*(n-k)}{\sum_{i=0}^{N-1}|x(n-i)|^2} \right|^2 \quad 2.6.4$$

如果  $x(n)$  和  $v(n)$  是不相关的白噪声信号，则下式：

$$E\{\Lambda(n+1) | \Lambda(n), x(n)\} = \Lambda(n) \left[ 1 - \frac{2\mu}{N} + \frac{\mu^2}{N} + \frac{\mu^2 \sigma_v^2}{\Lambda(n) \sum_{i=0}^{N-1}|x(n-i)|^2} \right] \quad 2.6.5$$

可以通过求解  $\partial E\{\Lambda(n+1)\} / \partial \mu = 0, \Lambda \neq 0$ ：

$$\frac{-2}{N} + \frac{2\mu}{N} + \frac{2\mu\sigma_v^2}{\Lambda(n) \sum_{i=0}^{N-1}|x(n-i)|^2} = 0 \quad 2.6.6$$

求解后得到最优步长：

$$\mu_{opt}(n) = \frac{1}{1 + \frac{\sigma_v^2}{\Lambda(n)(1/N) \sum_{i=0}^{N-1}|x(n-i)|^2}} \quad 2.6.7$$

期望  $\Lambda(n)(1/N) \sum_{i=0}^{N-1}|x(n-i)|^2$  等于剩余回声的方差  $\sigma_r^2(n)$ ，如果剩余回声的方差值等于 0，则步长因子等于 1， $r(n) = y(n) - \hat{y}(n)$ ，则有输出信号的方差是：

$$\sigma_e^2(n) = \sigma_v^2(n) + \sigma_r^2(n) \quad 2.6.8$$

这样可以求得这种情况下的最优步长因子为：

$$\mu_{opt}(n) \approx \frac{\sigma_r^2(n)}{\sigma_e^2(n)} \quad 2.6.9$$

则最优步长因子如下：

$$\hat{\mu}_{opt}(n) = \min\left(\frac{\hat{\sigma}_r^2(n)}{\hat{\sigma}_e^2(n)}, 1\right) \quad 2.6.10$$

当  $\Lambda(n) \approx \frac{\sigma_v^2}{\sigma_x^2 \left( \frac{2}{\mu} - 1 \right)}$  时，式 2.6.5 的迭代将停止（滤波器系数不在更新，

$E\{\Lambda(n+1)\} = \Lambda(n)$ ）。将 2.6.9 带入 2.6.10 得到在滤波器系数不更新情况下的剩余回声：

$$\sigma_r^2(n) \approx \min\left(\frac{1}{2}\hat{\sigma}_r^2(n), \sigma_v^2(n)\right) \quad 2.6.11$$

### 频域解

和时域相比，频域可以使步长因子  $\mu(k, l)$  按频域划分， $Y(k, l)$  和  $E(k, l)$  分别是频域中的记号，其和时域中的  $\hat{y}(n)$  和  $e(n)$  是对等的关系。 $k$  是频域索引， $l$  是帧索引。和 2.6.9 类似，可得频域步长因子如下：

$$\mu_{opt}(k, l) \approx \frac{\sigma_r^2(k, l)}{\sigma_e^2(k, l)} \quad 2.6.12$$

假设滤波器有一个和频谱无关的泄露（滤波器的误差）系数  $\eta(l)$ ，这将得到：

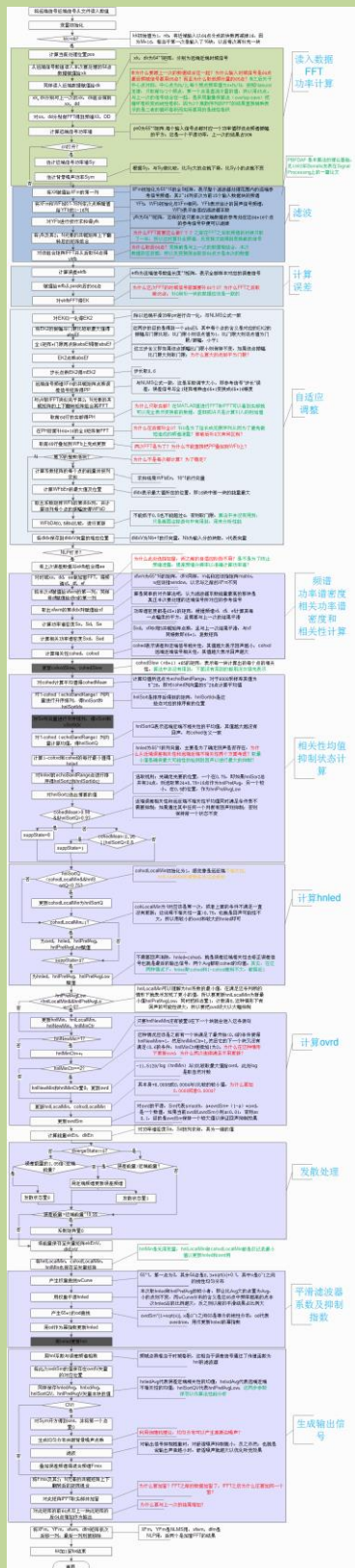
$$\hat{\sigma}_r^2(k, l) = \hat{\eta}(l) \hat{\sigma}_y^2(k, l) \quad 2.6.13$$

$\eta(l)$  实际上是滤波器的回声返回损失增强（ERLE）。

为了让步长因子调节的更快，使用瞬时估计， $\hat{\sigma}_y(k, l) = |Y(k, l)|^2$  和  $\hat{\sigma}_e(k, l) = |E(k, l)|^2$ ，这将使 2.6.10 步长因子调整变为下式：

$$\hat{\mu}_{opt}(k, l) = \min\left(\hat{\eta}(l) \frac{|\hat{Y}(k, l)|^2}{|E(k, l)|^2}, \mu_{\max}\right) \quad 2.6.14$$

$\mu_{\max}$  是小于等于 1 的数，以确保滤波器稳定。



## 第三章 语音阵列信号处理

和单麦克相比，麦克风阵列能够提供空域信息，这可以用来实现语音波束形成、盲信号分离以及声源定位。麦克风阵列主要分为两大类，它们是固定波束形成和自适应波束形成。固定波束形成包括：

- 延迟和波束形成
- 权重和波束形成
- 滤波和波束形成

自适应波束形成包括：

- LCMV 波束形成
- 广义互相关

这两类各有特点也各有用途。

麦克风阵列设计对语音信号的波束形成性能和噪声抑制性能影响较大。麦克风阵列的设计需要考虑的因数包括：麦克风个数，阵列的空间布局以及使用的麦克风类型。麦克风设计需要时刻考虑如下要数：

- 对噪声的抑制性能取决于信号的入射角和频率。
- 增加麦克风可以增加噪声的抑制能力。
- 在超过空域奈奎斯特定理的前提下，增加麦克风之间的间距能够增加噪声抑制能力，违反奈奎斯特准则时会导致空域频谱泄露，进而导致噪声抑制能力下降。
- 麦克风间的空间布局决定麦克风波束形成方向图。
- 和各向同性麦克相比，定向麦克风可以增加指向性。

麦克风阵列的语音波束方向图

- 波束方向图（指向图），通常针对远场语音模型而言。2D，3D 以及极坐标表示。
- 指向性指数，是频率的函数，常用 dB 为值来测量，指示的是麦克风阵列的指向性。
- 半功率波束宽度（波束宽度），指主瓣峰值的 1/2（3dB）的宽度。
- 第一空波束宽度，决定了空域分辨两个声源的能力，即两个声源接近到何程度，麦克风阵列无法分辨。
- 阵列增益。全方位信号和阵列输入信号之间的增益。

阵列设计的任务是在一个期望的方向图的前提下，设计在旁瓣幅度值允许的条件下，最小的波束宽度。阵列方向图设计的方法有如下几种：

- 延迟和技术
- 恒定束宽波束形成技术
- 最小方差无失真波束响应波束形成

- Frost 波束形成
- 广义旁瓣相消技术（GSC）

究竟如何选择合适的语音波束形成技术，需要理清如下问题：

- ✧ 设备使用场景，声源和麦克风的相对位置是运动还是静止的。
- ✧ 麦克风阵列的限制是什么？麦克风个数和它们的间距。
- ✧ 设备所在环境的噪声源通常是什么？
- ✧ 算法依赖平台的内存和处理能力。

## 3.1 阵列模型

### 3.1.1 线阵模型

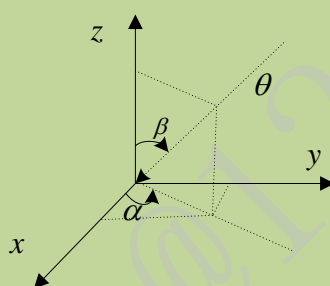


图 3.1 入射角三维坐标

入射方向的单位向量可表示为：

$$v(\theta) = [\sin(\beta)\cos(\alpha) \quad \sin(\beta)\sin(\alpha) \quad \cos(\beta)]^T \quad 3.1$$

其中， $\alpha$  和  $\beta$  分别水平方位角和垂直俯仰角。是对于该坐标系中的任意一点，其可以使用  $x$ ,  $y$ ,  $z$  三个坐标轴表示出来。

$$p_m = [x_m \quad y_m \quad z_m]^T \quad 3.1$$

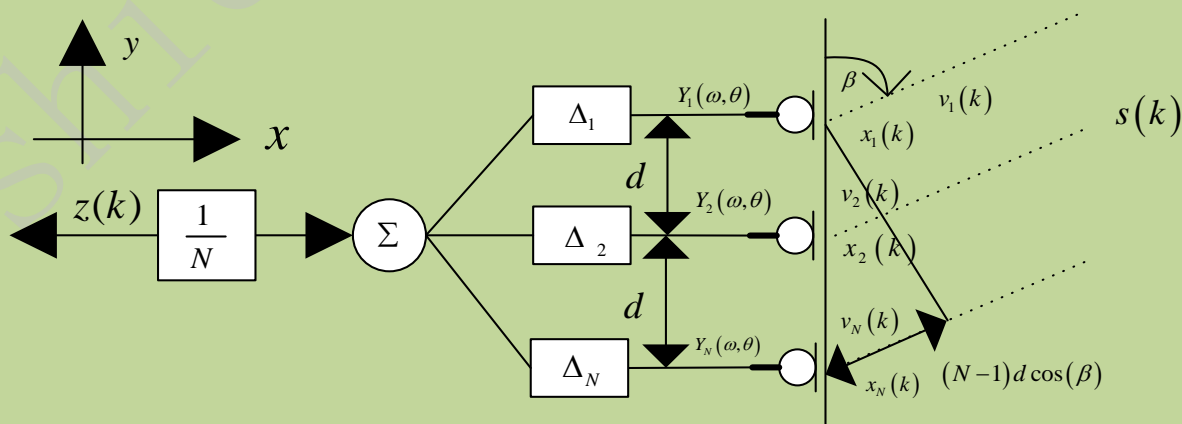


图 3.2 线阵模型

如上图所示线阵模型，设声波到麦克风阵列时，其波前可以被看成是平面波，入射角是  $\theta = (\alpha, \beta)$ ，即于平面法线方向的夹角。麦克风阵列 TDOA（time difference of array）的计算

公式如下：

$$TDOA = \tau_n(\theta) = \frac{v^T(\theta) p_n}{c}, n=1, \dots, N \quad 3.1$$

$c$  表示声速， $n$  是麦克风的索引， $\theta$  是声源的入射角（麦克风阵列波束形成方向）。为简化阵列模型，现将阵列降为二维处理，

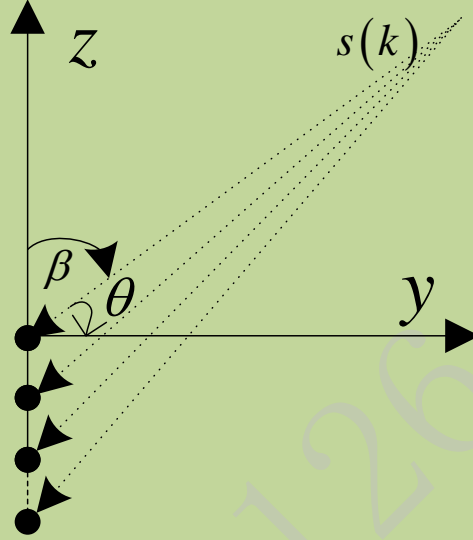


图 3.3 降维后入射角关系

虽然到达麦克风的声波并不平行，但在平面波模型下，可以认为同一声源发出的声音到达每个麦克风声波是平行的，即他们的入射角都是  $\theta$ ，这时  $\theta$  退化为  $yz$  平面的角度，在  $z$  轴方向上和  $y$  轴方向上声源到各个麦克的时间均有差异。设麦克风之间的距离为  $d$ ，则在  $z$  轴上的时间差（式 3.1 退化）为：

$$\tau_n(\theta) = \frac{(n-1)d \cos \theta}{c} \quad 3.2$$

则各个麦克风的输出  $y_n(k)$  为：

$$y_n(k) = s(k - \tau_n(\theta)) + v_n(k), n=1, 2, \dots, N \quad 3.3$$

其中， $v_n(k)$  为噪声， $s(k - \tau_n(\theta))$  是声源  $s(k)$  经过 TDOA 后的值。将阵列收到的信号使用向量表示如下：

$$x_s(t) = [s_1(t) \quad s_2(t) \quad \dots \quad s_N(t)] \quad 3.4$$

对式 3.3 进行傅里叶变换后，得如下公式：

$$Y_n(\omega, \theta) = \int_{-\infty}^{+\infty} y_n(t) e^{-j\omega t} dt = \int_{-\infty}^{+\infty} y(t - \tau_n(\theta)) e^{-j\omega t} dt = e^{-j\omega \tau_n(\theta)} S(\omega) \quad 3.5$$

再对  $y$  轴做一次类似的延迟变换求解：

$$z[k] = \frac{1}{N} \sum_{n=1}^N y_n[k + \Delta_n], \Delta_n = \frac{(n-1)d \cos(\beta)}{c} \quad 3.6$$

对式 3.6 做傅里叶变换得：

$$Z(\omega, \theta) = \frac{1}{N} \sum_{n=1}^N e^{j\omega \Delta_n} Y_n(\omega, \theta) = \frac{1}{N} \sum_{n=1}^N e^{j\omega \Delta_n} S(\omega) e^{-j\omega \tau_n(\theta)} \quad 3.7$$

根据  $Z(\omega, \theta)$  和  $S(\omega)$  可得到响应函数如下：



$$H(\omega, \theta) = \frac{Z(\omega, \theta)}{S(\omega)} = \frac{1}{N} \sum_{n=1}^N e^{-j\omega(n-1) \frac{d(\cos(\theta) - \cos(\beta))}{c}} \quad 3.8$$

由此可求得天线的方向图(使用到三角中的 2 倍角公式):

$$\begin{aligned} |H(\omega, \theta)| &= \left| \frac{1 - e^{-jN\omega \frac{d(\cos \theta - \cos \beta)}{c}}}{N(1 - e^{-j\omega \frac{d(\cos \theta - \cos \beta)}{c}})} \right| = \left| \frac{1 - \cos(N\omega \frac{d(\cos \theta - \cos \beta)}{c})}{N(1 - \cos(\omega \frac{d(\cos \theta - \cos \beta)}{c}))} \right| \\ &= \left| \frac{\sin(\frac{N\omega d(\cos \theta - \cos \beta)}{2c})}{N \sin(\frac{\omega d(\cos \theta - \cos \beta)}{2c})} \right| \end{aligned} \quad 3.9$$

根据麦克风个数，麦克风间距以及频率三个变量（同时设  $\beta$  成  $90^\circ$ ）可以画出如下三张图：

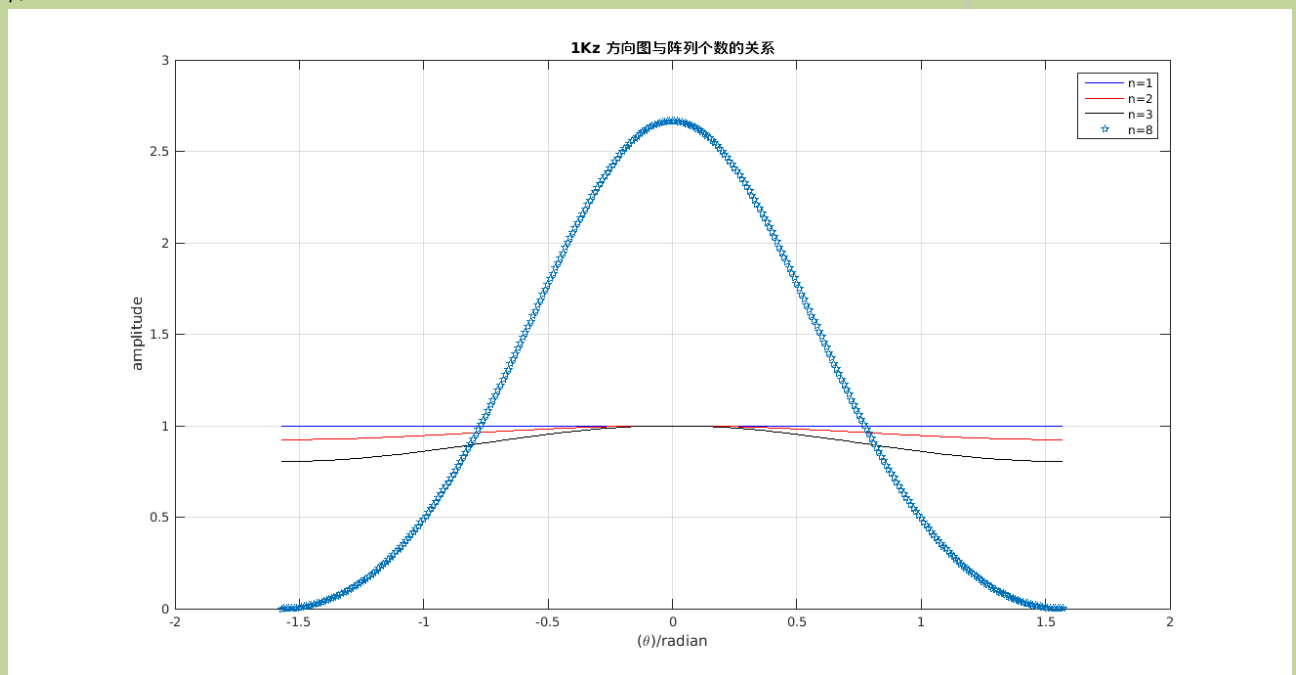


图 3.4 麦克风数量与波束方向图

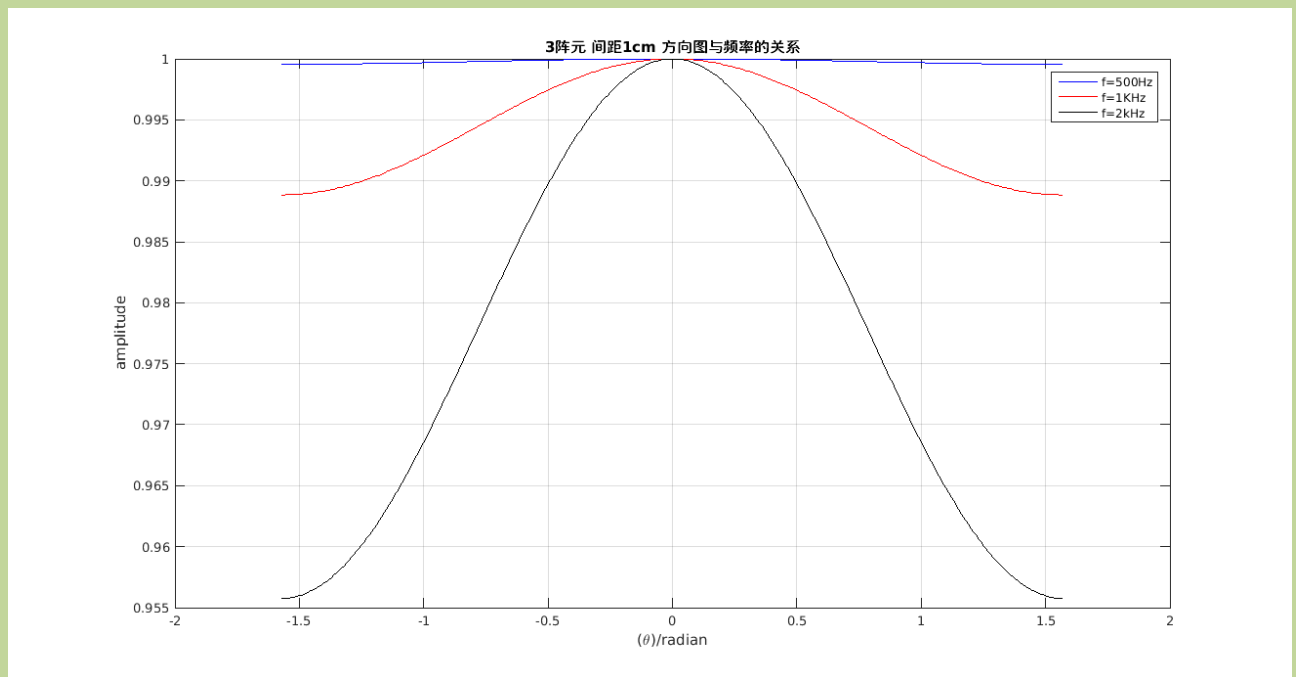


图 3.5.3 麦克风与波束方向图关系

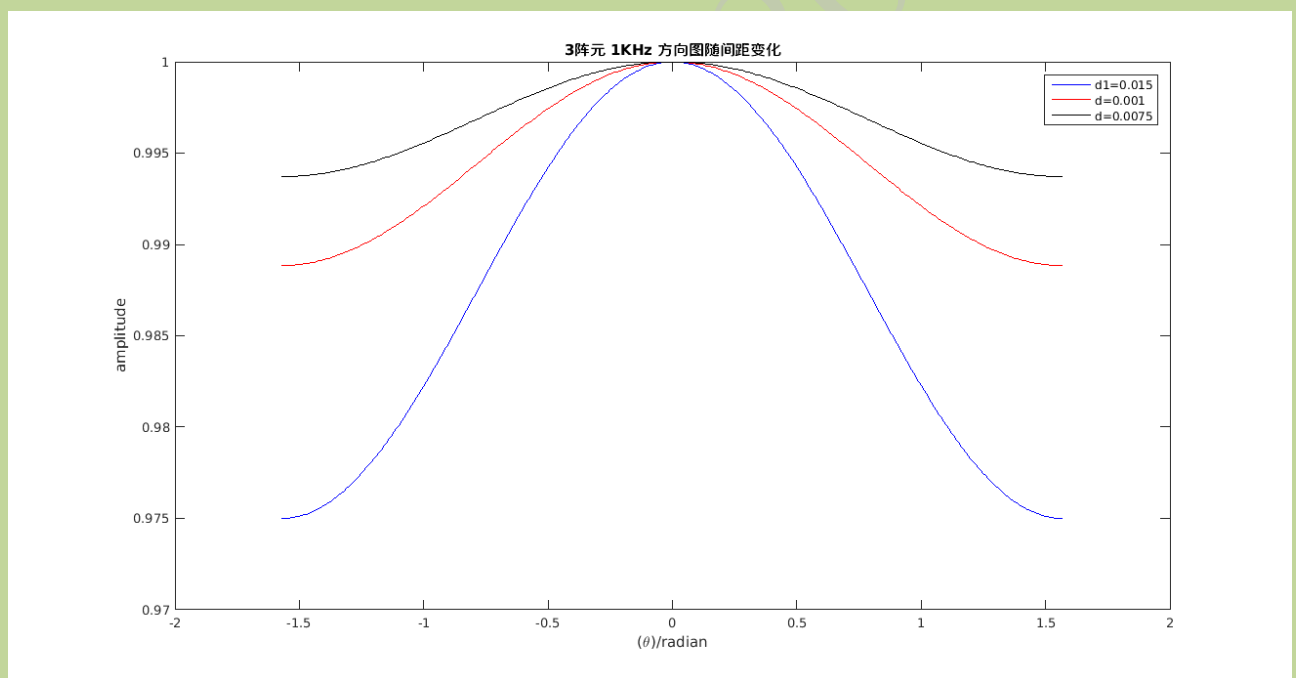


图 3.5.3 麦克风间距与波束方向图关系

概括上述三张图可知，麦克风数量越多，其波束指向性越好，在满足频域奈奎斯特采样定理前提下，麦克风间距越大其空间分辨率越好。当信号频率越高时，麦克风阵列的作用越明显。

图 3.5.3 麦克风与波束方向图关系

```
clear;
v=343;f=1000;
%Number of element
sita=-pi:0.01:pi;
lamda=v/f;
d=lamda/8;
nl=1;
beta=2*pi*d*sin(sita)/lamda;
```

```

z11=(n1/2)*beta;
z21=(1/2)*beta;
f1=sin(z11)./(n1*sin(z21));
F1=abs(f1);
figure(1);
% plot(sita,F1,'b');
polar(sita, F1,'b');
hold on;
n2=2;
beta=2*pi*d*sin(sita)/lamda;
z12=(n2/2)*beta;
z22=(1/2)*beta;
f2=sin(z12)./(n2*sin(z22));
F2=abs(f2);
% plot(sita,F2,'r');
polar(sita, F2,'r');
hold on;
n3=3;
beta=2*pi*d*sin(sita)/lamda;
z13=(n3/2)*beta;
z23=(1/2)*beta;
f3=sin(z13)./(n3*sin(z23));
F3=abs(f3);
polar(sita, F3,'k');
% plot(sita,F3,'k');
hold on;
n4=8;
beta=2*pi*d*sin(sita)/lamda;
z14=(n4/2)*beta;
z24=(1/2)*beta;
f4=sin(z14)./(n4*sin(z24));
F4=abs(f4);
% plot(sita, F4, 'p');
polar(sita, F4,'p');
hold off;
grid on;
xlabel('(\theta)/radian');
ylabel('amplitude');
title('1Kz 方向图与阵列个数的关系');
legend('n=1','n=2','n=3','n=8');

% Wave length
clear;figure;
v=343;f=1000;
sita=-pi/2:0.01:pi/2;
n=3;
d=0.01;
lamda1=v/(f/5);
beta=2*pi*d*sin(sita)/lamda1;
z11=(n/2)*beta;
z21=(1/2)*beta;

```

```

f1=sin(z11)./(n*sin(z21));
F1=abs(f1);
polar(sita, F1,'b');
hold on;
lamda2=(v/f);
beta=2*pi*d*sin(sita)/lamda2;
z12=(n/2)*beta;
z22=(1/2)*beta;
f2=sin(z12)./(n*sin(z22));
F2=abs(f2);
polar(sita, F2,'r');
hold on;
lamda3=v/(f*2);
beta=2*pi*d*sin(sita)/lamda3;
z13=(n/2)*beta;
z23=(1/2)*beta;
f3=sin(z13)./(n*sin(z23));
F3=abs(f3);
polar(sita, F3,'k');
% plot(sita,F1,'b',sita,F2,'r',sita,F3,'k');
hold off;
grid on;
xlabel('(\theta)/radian');
ylabel('amplitude');
title('3阵元 间距1cm 方向图与频率的关系');
legend('f=500Hz','f=1KHz','f=2kHz');

```

```

clear;figure;
v=343;f=1000;
sita=-pi:0.01:pi;
n=3;
lamda=v/f;
d1=0.015;
beta=2*pi*d1*sin(sita)/lamda;
z11=(n/2)*beta;
z21=(1/2)*beta;
f1=sin(z11)./(n*sin(z21));
F1=abs(f1);
polar(sita, 10*log(F1),'b');
% plot(sita,F1,'b');
hold on;
d2=0.01;
beta=2*pi*d2*sin(sita)/lamda;
z12=(n/2)*beta;
z22=(1/2)*beta;
f2=sin(z12)./(n*sin(z22));
F2=abs(f2);
% plot(sita,F2,'r');
polar(sita, 10*log(F2),'r');
hold on;

```

```

d3=0.0075;
beta=2*pi*d3*sin(sita)/lamda;
z13=(n/2)*beta;
z23=(1/2)*beta;
f3=sin(z13)./(n*sin(z23));
F3=abs(f3);
% plot(sita,F3,'k');
polar(sita, 10*log(F3),'k');
hold off;
hold on;
xlabel(' \theta /radian');
ylabel('amplitude dB');
title('3阵元 1KHz 方向图随间距变化');
legend('d1=0.015','d=0.001','d=0.0075');

```

### 3.1.2 面阵模型

### 3.1.3 圆阵模型

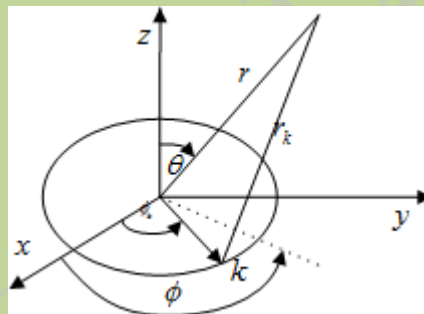


图 3.6  $N$  个均匀分布圆阵

在一个半径为  $R$  的圆周上，均匀分布  $N$  个相同的各向同性阵元。其中  $\phi$  是起始于  $x$  轴正向的方位角， $\theta$  是起始于  $z$  轴正向的俯仰角。

假设每个阵元的权重为  $W_n (n=0,1,\dots,N-1)$ ，则第  $N$  个阵元的方位角  $\phi_n = 2n\pi / N$ 。如果来自某一  $(\theta, \phi)$  方向的平面波照射在圆形阵列上，则第  $n$  个阵元相对于阵列中心的相对相位为：

$$\beta_n = -kR \cos(\phi - \phi_n) \sin(\theta) \quad 3.10$$

则  $N$  个均匀分布阵元所形成的圆形阵列因子的方向图函数为：

## 3.2 阵列波束形成技术

### 3.2.1 DS 模型

Delay and Sum（延迟和）

延迟调节的是接收天线的方向，和是滤波器系数，决定主瓣和旁瓣的大小和形状。

$N$  阶空域滤波器将有  $N-1$  个零点。可以设计何时的不同的权系数以抑制特定方向的干扰信号。

当麦克风阵列的间距增加时，阵列的波束宽度减小。所以如果当你需要一个尖锐的波束，可以简单的增加麦克风间距。但是这会增加麦克风阵列的孔径，这通常也意味着引入更多噪声。所以，在阵列设计时，我们期望麦克风的阵列的间距越大越好。但是当间距大于  $\frac{\lambda}{2} = \frac{c}{2f}$  时，将会产生空域混叠。 $\lambda$  是信号波长。

为了避免空域混叠，阵列间距必须满足  $d \leq \frac{\lambda}{2} = \frac{c}{2f}$ 。人耳朵的听觉范围是 20Hz~20KHz，人发声的范围是 85Hz~1.1KHz，对应的阵列间距是 13.2cm。

N 个麦克风阵列的波束形成方程如下：

$$y_n(k) = \alpha_n s[k - t - \mathcal{F}_n(\tau)] + v_n(k) = x_n(k) + v_n(k), n = 1, 2, \dots, N, \quad 3.11$$

式中  $\alpha_n$  取值介于 0 和 1 之间，是传输的衰减因子。 $s(k)$  是未知的语音信号（要识别的人说的话）， $t$  是从语音信源传递到第一个传感器的时间。 $v_n(k)$  是第  $n$  个传感器的加性噪声。 $\tau$  是第一个和第二个麦克风的相对延迟， $\mathcal{F}_n(\tau)$  是第一个麦克风相对第  $n$  个麦克风的传播延迟， $\mathcal{F}_1(\tau) = 0$ ， $\mathcal{F}_2(\tau) = \tau$ 。

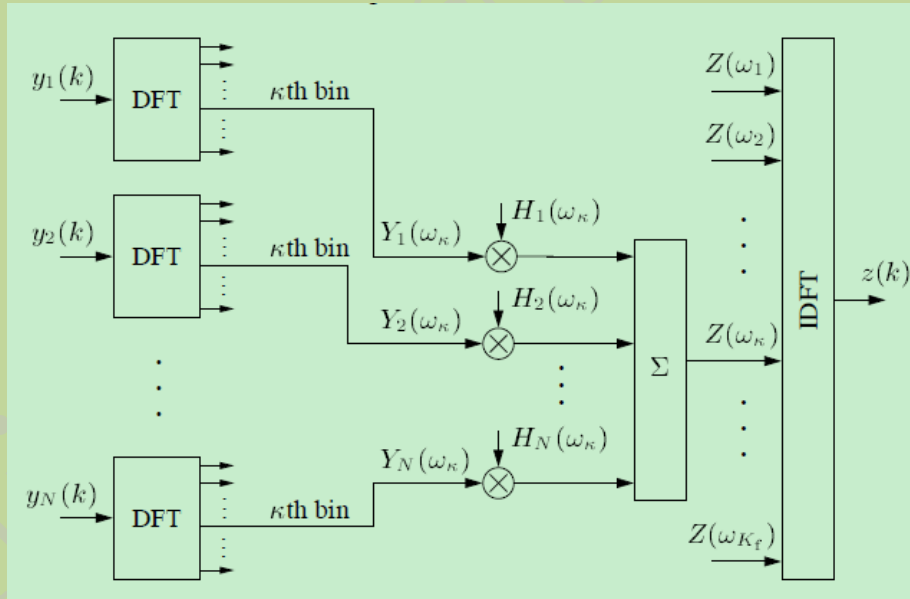


图 3.7 频域宽带波束形成

延迟和模型

将上式进行时移后得：

$$\begin{aligned} y_{a,n}(k) &= y_n[k + \mathcal{F}_n(\tau)] = \alpha_n s(k - t) + v_{a,n}(k) \\ &= x_{a,n}(k) + v_{a,n}(k), n = 1, 2, \dots, N \end{aligned} \quad 3.11$$

此时， $v_{a,n}(k) = v_n[k + \mathcal{F}_n(\tau)]$

接下来是对其的信号进行累加得：

$$z_{DS}(k) = \frac{1}{N} \sum_{n=1}^N y_{a,n}(k) = \alpha_s s(k-t) + \frac{1}{N} v_s(k) \quad 3.12$$

其中：

$$\alpha_s = \frac{1}{N} \sum_{n=1}^N \alpha_n, \quad v_s(k) = \sum_{n=1}^N v_{a,n}(k) = \sum_{n=1}^N v_n[k + F_n(\tau)]$$

则输入和输出的信噪比为：

$$SNR = \frac{\sigma_{x_1}^2}{\sigma_{v_1}^2} = \alpha_1^2 \frac{\sigma_s^2}{\sigma_{v_1}^2} \quad 3.13$$

其中  $\sigma_{x_1}^2 = E[x_1^2(k)]$ ,  $\sigma_{v_1}^2 = E[v_1^2(k)]$ ,  $\sigma_s^2 = E[s^2(k)]$ 。

则在 DS 后输出 SNR 可以表达如下：

$$oSNR = N^2 \alpha_s^2 \frac{E[s^2(k-t)]}{E[v_s^2(k)]} = N^2 \alpha_s^2 \frac{\sigma_s^2}{\sigma_{v_s}^2} = \left( \sum_{n=1}^N \alpha_n \right)^2 \frac{\sigma_s^2}{\sigma_{v_s}^2} \quad 3.14$$

$$\text{式 3.14 中 } \sigma_{v_s}^2 = E \left\{ \left[ \sum_{n=1}^N v_n[k + F_n(\tau)] \right]^2 \right\} = \sum_{n=1}^N \sigma_{v_n}^2 + 2 \sum_{i=1}^{N-1} \sum_{j=i+1}^N \rho_{v_i v_j}, \quad \rho_{v_i v_j}$$

其中， $\sigma_{v_n}^2 = E[v_n^2(k)]$  是噪声  $v_n(k)$  的方差， $\rho_{v_i v_j} = E\{v_i[k + F_i(\tau)]v_j[k + F_j(\tau)]\}$  是噪声  $v_i(k)$  和  $v_j(k)$  的协方差。

特殊情况 1：

$$\rho_{v_i v_j} = 0, \quad oSNR = N \cdot SNR \quad 3.15$$

特殊情况 2：

假设噪声的能量是一样的并且衰退因子等于 1。

$$oSNR = \frac{N}{1 + \rho_s} \cdot SNR \quad 3.16$$

$$\rho_s = \frac{2}{N} \sum_{i=1}^{N-1} \sum_{j=i+1}^N \rho_{v_i v_j}, \quad \rho_{v_i v_j} = \frac{E\{v_i[k + F_i(\tau)]v_j[k + F_j(\tau)]\}}{\sigma_{v_i v_j}} \quad 3.17$$

### 3.2.2 最大信噪比

将 3.11 重写为如下式：

$$\mathbf{y}_a(k) = s(k-t)\boldsymbol{\alpha} + \mathbf{V}_a(k) \quad 3.18$$

此处：

$$\mathbf{y}_a(k) = [y_{a,1}(k) \quad y_{a,2}(k) \quad \dots \quad y_{a,N}(k)]^T,$$

$$\mathbf{v}_a(k) = [v_{a,1}(k) \quad v_{a,2}(k) \quad \dots \quad v_{a,N}(k)]^T,$$

$$\boldsymbol{\alpha} = [\alpha_1 \quad \alpha_2 \quad \dots \quad \alpha_N]^T$$

假设噪声和信号是不相关的， $\mathbf{y}_a(k)$  的相关矩阵可以表示为如下形式：

$$\mathbf{R}_{y_a y_a} = \sigma_s^2 \boldsymbol{\alpha} \boldsymbol{\alpha}^T + \mathbf{R}_{v_a v_a} \quad 3.19$$

其中， $\mathbf{R}_{v_a v_a} = E[\mathbf{v}_a(k) \mathbf{v}_a^T(k)]$  是噪声的相关矩阵。则波束形成的输出可写为如下形式：

$$z(k) = \mathbf{h}^T \mathbf{y}_a(k) = \sum_{n=1}^N h_n y_{a,n}(k) = s(k-t) \mathbf{h}^T \boldsymbol{\alpha} + \mathbf{h}^T \mathbf{v}_a(k) \quad 3.20$$

此处：

$$\mathbf{h} = [h_1 \quad h_2 \quad \dots \quad h_N]^T \quad 3.21$$

可得到信噪比如下：

$$SNR(\mathbf{h}) = \frac{\sigma_s^2 (\mathbf{h}^T \boldsymbol{\alpha})^2}{\mathbf{h}^T \mathbf{R}_{v_a v_a} \mathbf{h}} \quad 3.22$$

上述信噪比最大解，等价于求解下述特征值方程。

$$\sigma_s^2 \boldsymbol{\alpha} \boldsymbol{\alpha}^T \mathbf{h} = \lambda \mathbf{R}_{v_a v_a} \mathbf{h} \quad 3.23$$

设  $\mathbf{R}_{v_a v_a}$  可逆，则最优解是特征向量  $\mathbf{h}_{\max}$  和特征值  $\lambda_{\max}$ ：

$$z_{\max}(k) = \mathbf{h}_{\max}^T \mathbf{y}_a(k) \quad 3.24$$

$$SNR(\mathbf{h}_{\max}) = \lambda_{\max} \quad 3.24$$

### 3.2.3 最小方差无失真响应滤波器

Minimum variance distortionless response (MVDR)，其基本思想是调节滤波系数  $\mathbf{h}$  最小化输出功率  $E[z^2(k)] = \mathbf{h}^T \mathbf{R}_{y_a y_a} \mathbf{h}$ ，则 MVDR 可以表述为如下方程：

$$\min_{\mathbf{h}} \mathbf{h}^T \mathbf{R}_{y_a y_a} \mathbf{h}, \mathbf{h}^T \boldsymbol{\alpha} = \alpha_1 \quad 3.25$$

拉格朗日乘法器可以解 3.25，得：

$$\mathbf{h}_{opt} = \alpha_1 \frac{\mathbf{R}_{y_a y_a}^{-1} \boldsymbol{\alpha}}{\boldsymbol{\alpha}^T \mathbf{R}_{y_a y_a}^{-1} \boldsymbol{\alpha}} \quad 3.26$$

### 3.2.4 线性约束最小方差

LCMV(linearly constrained minimum variance)，室内的麦克风阵列通常会接受到来自墙



壁、地板、天花板以及家具的回声，这一现象可以使用混响模型来表示：

$$y_n(k) = g_n * s(k) + v_n(k) = x_n(k) + v_n(k) \quad 3.27$$

$g_n$  是第  $n$  个麦克风对语音信号的脉冲响应函数。上述模型的矩阵表示方式是：

$y_n(k) = \mathbf{g}_n^T \mathbf{s}(k) + v_n(k)$ ，这里：

$$\mathbf{g}_n = [g_{n,0} g_{n,1} \cdots g_{n,L_g-1}]^T,$$

$$\mathbf{s}(k) = [s(k) s(k-1) \cdots s(k-L_h+1)]^T$$

麦克风阵列可以重写为如下形式：

$$\mathbf{y}_n(k) = \mathbf{G}_n \mathbf{s}_L(k) + \mathbf{v}_n(k), n=1,2,\dots,N \quad 3.28$$

这里：

$$\mathbf{y}_n(k) = [y_n(k) \ y_n(k-1) \ \cdots \ y_n(k-L_h+1)]^T,$$

$$\mathbf{v}_n(k) = [v_n(k) \ v_n(k-1) \ \cdots \ v_n(k-L_h+1)]^T,$$

$$\mathbf{S}_L(k) = [s(k) \ s(k-1) \ \cdots \ s(k-L+1)]^T,$$

$$\mathbf{G}_n = \begin{bmatrix} g_{n,0} & \cdots & g_{n,L_g-1} & 0 & 0 & \cdots & 0 \\ 0 & g_{n,0} & \cdots & g_{n,L_g-1} & 0 & \cdots & 0 \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ 0 & 0 & \cdots & 0 & g_{n,0} & \cdots & g_{n,L_g-1} \end{bmatrix}$$

$\mathbf{G}_n$  是  $L_h \times L, L = L_h + L_g - 1$

将  $N$  维观测的向量写作：

$$\mathbf{y}(k) = [\mathbf{y}_1^T \ \mathbf{y}_2^T \ \cdots \ \mathbf{y}_N^T]^T = \mathbf{G} \mathbf{S}_L(k) + \mathbf{v}(k) \quad 3.29$$

其中：

$$\mathbf{G} = \begin{bmatrix} \mathbf{G}_1 \\ \mathbf{G}_2 \\ \vdots \\ \mathbf{G}_N \end{bmatrix}_{NL_h \times L}$$

$$\mathbf{V}(k) = [\mathbf{v}_1^T(k) \ \mathbf{v}_2^T(k) \ \cdots \ \mathbf{v}_N^T(k)]^T$$

在信号和噪声不相关时， $NL_h \times NL_h$  的协方差系数  $\mathbf{y}(k)$  是：

$$\mathbf{R}_{yy} = E[\mathbf{y}(k) \mathbf{y}^T(k)] = \mathbf{G} \mathbf{R}_{ss} \mathbf{G}^T + \mathbf{R}_{vv} \quad 3.30$$

LCMV 滤波器的解可按下式获得：

在满足  $\mathbf{G}^T \mathbf{h} = \mathbf{u}$  条件下的  $\min_h \mathbf{h}^T \mathbf{R}_{yy} \mathbf{h}$

其中：

$\mathbf{h} = [h_1^T \ h_2^T \ \cdots \ h_N^T]^T$ ， $\mathbf{u} = [1 \ 0 \ \cdots \ 0]^T$ ，这里的意义是在去混响的前提下，使噪声最小。

上式的最优解是：

$$\mathbf{h}_R = \mathbf{R}_{yy}^{-1} \mathbf{G} (\mathbf{G}^T \mathbf{R}_{yy}^{-1} \mathbf{G})^{-1} \mathbf{u}^T, \quad (3.31)$$

其中  $\mathbf{R}$  表示的是混响模型，假设  $\mathbf{R}_{yy}$  和  $\mathbf{R}_{vv}$  是正定矩阵， $(\mathbf{G}^T \mathbf{R}_{yy}^{-1} \mathbf{G})$  非奇异的条件是  $N L_h \geq L$ ，这就意味着：

$$L_h \geq \frac{L_g - 1}{N - 1} \quad (3.32)$$

其解如下：

$$\mathbf{h}_R = \mathbf{R}_{vv}^{-1} \mathbf{G} (\mathbf{G}^T \mathbf{R}_{vv}^{-1} \mathbf{G})^{-1} \mathbf{u} \quad (3.33)$$

### 3.2.3 广义旁瓣相消

广义旁瓣相消技术将约束的 LCMV 最优问题变成了一个无约束问题。在广义旁瓣相消中，有：

$$\mathbf{h} = \mathbf{f} - \mathbf{B}\mathbf{w} \quad (7.23)$$

此处：

$$\mathbf{f} = \mathbf{G}_{\cdot 1} [\mathbf{G}_{\cdot 1}^T \ \mathbf{G}_{\cdot 1}^T]^{-1} \mathbf{u} \quad (7.24)$$

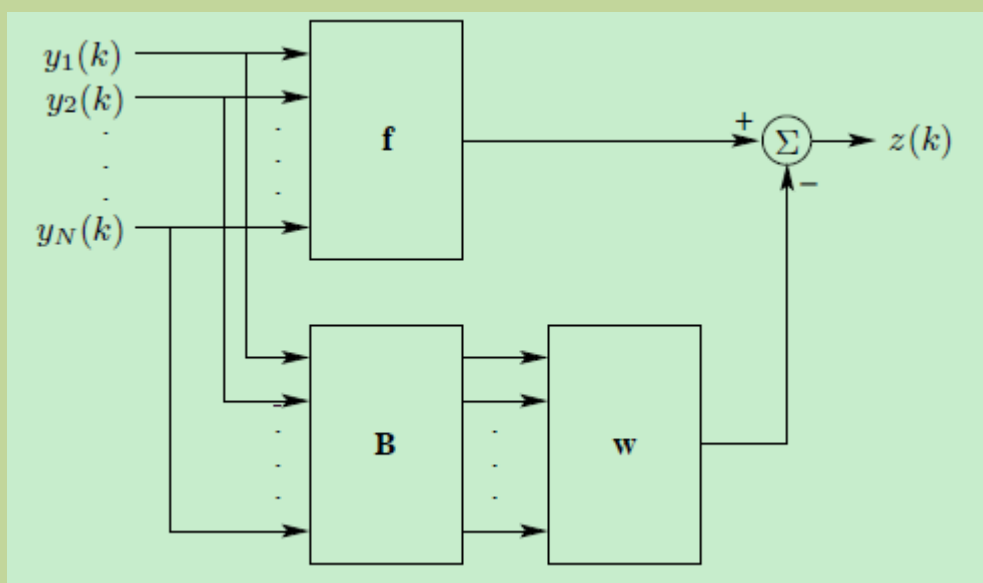
上式是  $\mathbf{G}_{\cdot 1}^T \mathbf{f} = \mathbf{u}$  的最小范数解。 $\mathbf{B}$  是阻塞矩阵， $\mathbf{W}$  是权向量。广义旁瓣是下式无限制的最优解：

$$\min_{\mathbf{w}} (\mathbf{f} - \mathbf{B}\mathbf{w})^T \mathbf{R}_{yy} (\mathbf{f} - \mathbf{B}\mathbf{w}) \quad (7.25)$$

其解是：

$$\mathbf{w}_{GSC} = [\mathbf{B}^T \mathbf{R}_{yy} \mathbf{B}]^{-1} \mathbf{B}^T \mathbf{R}_{yy} \mathbf{f} \quad (7.26)$$

其中  $\mathbf{f}$  代表 fixed Beamformer， $\mathbf{B}$  代表 Blocking matrix， $\mathbf{W}$  是自适应权向量。



广义旁瓣相消

$$h_{LCMV} = R_{yy}^{-1} G_{:1} \left[ G_{:1}^T R_{yy}^{-1} G_{:1} \right]^{-1} u = \left\{ I_{2Lh \times 2Lh} - B \left[ B^T R_{yy} B \right]^{-1} B^T R_{yy} \right\} f$$

$$= h_{GSC}$$

since every 4 - 5dB improvement of the Signal to Noise Ratio (SNR) may raise the speech intelligibility by 50% [1]

### 3.3 基于阵列定位和跟踪技术

麦克风阵列在空域实现了噪声抑制，此外，还能定位声源。即两个问题，一个是 DOA（direction of arrival），另一个是定位声源空间坐标定位。

DOA 算法的主要思想是找到各个麦克风收集到的信号的相位信息关系。主要有三类 DOA 算法：第一类是高分辨谱估计技术，这类技术在计算时需要用到相关矩阵，计算量较大，第二类是最大输出功率可控波速形成技术，利用加权函数求和得到，如 MVB（maximum variance beamformer），DSB（Delay and Sum Beamformer），多声源定位问题可以使用 MUSIC 和 ESPRIT 算法。第三类是基于时间差的定位技术，这种技术利用声源到各个麦克风阵列的时间差，然后求出角度，其计算量相对较小，对噪声处理也比较好，但只使用与单个声源的定位，对多个声源的定位效果比较差。

#### 3.3.1 互相关方法

考虑只有两个麦克风的场景，互相关函数如下：

$$\begin{aligned}
R_{y_1 y_2}(p) &= E[y_1(k)y_2(k+p)] \\
&= E[(\alpha_1 s(k-t) + v_1(k))(\alpha_2 s(k-t-\tau+p) + v_2(k+p))] \\
&= \alpha_1 \alpha_2 R_{ss}(p-\tau) + \alpha_1 R_{sv_2}(p+t) + \alpha_2 R_{sv_1}(p-t-\tau) + R_{v_1 v_2}(p)
\end{aligned}
\tag{3.3.1.1}$$

如果噪声之间不相关，并且和输入信号也不相关，则在  $p = \tau$  时，信号互相关函数达

$$\text{到极大值 } \hat{\tau} = \arg \max_p R_{y_1 y_2}(p), p \in [-\tau_{\max}, \tau_{\max}]
\tag{3.3.1.2}$$

这种方法估计 DOA 简单，但是也有显而易见的缺点，在有反射，混响的模型中，这一方法并不是很适用。在使用 TDOA 算法时，通常要注意空域混叠。

#### ➤ 互相关计算过程

$$r_{xy}(m) = \sum_{n=-\infty}^{\infty} x(n)y(n+m)
\tag{3.3.1.2}$$

其计算过程是保持  $x(n)$  不动， $y(n)$  左移  $m$  ( $m$  可正可负) 个抽样点后，两个序列相乘的结果。

```

x1 =
    1     2     3     4     5

>> x2=[1 0 2 3 4]
|
x2 =
    1     0     2     3     4

>> xcorr(x1,x2)

ans =
    4.0000    11.0000    20.0000    29.0000    39.0000    25.0000    13.0000     4.0000     5.0000

```

图 3.3.1.1 MATLAB 计算互相关结果

当  $m=-5$  时：

$$r_{xy}(4) = \sum_{n=-\infty}^{\infty} x(n)y(n+4) = \dots + x(-1)y(3) + x(0)y(4) + \dots = x(0)y(4) = 4$$

其它值依次类推可以计算出来。

对于截取长度为  $N$  的两路麦克风信号， $m$  可取值从  $-(N-1)$  到  $(N-1)$ ，对于  $N$  点信号，共可以计算出  $2N-1$  点互相关， $m$  值越大，对于  $N$  点有限长信号，其有意义计算的长度越短，实际中通常选择  $m \ll N$ 。其计算复杂度是  $O(N^2)$ ，当数据量较大时，计算较为耗时，可以采用 FFT 和 IFFT 来计算互相关函数以减少运算量，其计算复杂度是  $\frac{N}{2} \log_2 N$ 。

设  $y(n)$  是  $x(n)$  与  $h(n)$  的互相关函数，即

$$y(m) = \frac{1}{N} \sum_{n=0}^{N-1} x(n)h(n+m) \quad 3.3.1.3$$

则,

$$Y(e^{j\omega}) = X^*(e^{j\omega})H(e^{j\omega}) \quad 3.3.1.4$$

互相关可以表示成下式:

$$r_{xy}(m) = \sum_{n=-\infty}^{\infty} x(n)y(n+m) = x(-m)*y(m) \quad 3.3.1.5$$

相关函数在频域也不是完全乘积, 是一个信号的共轭乘以另一个信号的原始频谱。

### 3.3.2 广义互相关 (GCC Generalized Cross-Correlation)

GCC 算法用于远场语音 DOA 估计时, 其表达式和互相关法类似:

$$\hat{\tau} = \arg \max_p R_{y_1 y_2}(p), p \in [-\tau_{\max}, \tau_{\max}] \quad 3.3.2.1$$

$$R_{y_1 y_2}(p) = F^{-1}[\Psi_{y_1 y_2}(f)] = \int_{-\infty}^{\infty} \Psi_{y_1 y_2}(f) e^{j2\pi fp} df = \int_{-\infty}^{\infty} \mathcal{G}(f) \phi_{y_1 y_2}(f) e^{j2\pi fp} df \quad 3.3.2.2$$

$F^{-1}[\cdot]$  是离散傅里叶反变换 (IDTFT)。

$$\phi_{y_1 y_2}(f) = E[Y_1(f)Y_2^*(f)] \quad 3.3.2.3$$

是互相关谱

$$Y_n(f) = \sum_k y_n(k) e^{-j2\pi fk}, n=1,2, \quad 3.3.2.4$$

$\mathcal{G}(f)$  是频域权向量。选择不同的权向量可以产生不同的 GCC 方法。

#### 3.3.2.1 经典互相关

如果  $\mathcal{G}(f)=1$ , 则 GCC 退化为类似 3.3.2.1 节的互相关情况, 区别是互相关函数用了离散傅里叶变换 (DFT) 和离散傅里叶反变换 (IDFT)。对于远场模型:

$$Y_n(f) = \alpha_n S(f) e^{-j2\pi f[t-\tau(n)]} + V_n(f), n=1,2 \quad 3.3.2.5$$

则可得:

$$\psi_{y_1 y_2}(f) = \alpha_1 \alpha_2 e^{-j2\pi f\tau} E[|S(f)|^2] \quad 3.3.2.6$$

$\psi_{y_1 y_2}(f)$  取决于声源的这一性质可以用来估计 TDOA。

### 3.3.2.2 平滑相关变换

为了克服语音信号在估计 TDOA 时的波动性，一个有效的方法是在计算交叉谱之前先白化麦克风阵列的输出，这等价于下式：

$$g(f) = \frac{1}{\sqrt{E[|Y_1(f)|^2]E[|Y_2(f)|^2]}} \quad 3.3.2.7$$

由此可得：

$$\begin{aligned} \psi_{y_1 y_2}(f) &= \frac{\alpha_1 \alpha_2 e^{-j2\pi f \tau} E[|S(f)|^2]}{\sqrt{E[|Y_1(f)|^2]E[|Y_2(f)|^2]}} \\ &= \frac{\alpha_1 \alpha_2 e^{-j2\pi f \tau} E[|S(f)|^2]}{\sqrt{\alpha_1^2 E[|S(f)|^2] + \sigma_{v_1}^2(f)} \cdot \sqrt{\alpha_2^2 E[|S(f)|^2] + \sigma_{v_2}^2(f)}} \\ &= \frac{e^{-j2\pi f \tau}}{\sqrt{1 + \frac{1}{\text{SNR}_1(f)}} \cdot \sqrt{1 + \frac{1}{\text{SNR}_2(f)}}} \end{aligned} \quad 3.3.2.8$$

这就意味着 TDOA 估计的性能和 SNR 相关，当性噪比足够大时：

$$\psi_{y_1 y_2}(f) = e^{-j2\pi f \tau}$$

### 3.3.2.3 相位变换 PATH

TDOA 依赖于信号的相位信息而非信号的幅度，所以可以简单的丢掉幅度信息而保留相位信息。

$$g(f) = \frac{1}{|\phi_{y_1 y_2}(f)|} \quad 3.3.2.9$$

此时：

$$\psi_{y_1 y_2}(f) = e^{-j2\pi f \tau} \quad 3.3.2.10$$

这时可得广义旁瓣互相关函数为：

$$R_{y_1 y_2}(p) = \int_{-\infty}^{\infty} e^{j2\pi f(p-\tau)} df = \begin{cases} \infty, & p = \tau \\ 0, & \text{otherwise} \end{cases}$$

### 3.3.2.4 多通道互相关系数法

使用多通道自相关系数法，将新的信号向量表示成如下形式：

$$y_a(k, p) = [y_1(k) \quad y_2(k + \tau_2(p)) \quad \cdots \quad y_N(k + \tau_N(p))]^T \quad 3.3.2.11$$

由此可得空域互相关矩阵：

$$R_a(p) = E[y_a(k, p) y_a^T(k, p)] = \begin{bmatrix} \sigma_{y_1}^2 & r_{y_1 y_2}(p) & \cdots & r_{y_1 y_N}(p) \\ r_{y_1 y_2}(p) & \sigma_{y_2}^2 & \cdots & r_{y_1 y_2}(p) \\ \vdots & \vdots & \ddots & \vdots \\ r_{y_1 y_2}(p) & r_{y_1 y_2}(p) & \cdots & \sigma_{y_N}^2 \end{bmatrix} \quad 3.3.2.12$$

空域自相关矩阵可以

### 3.3.3 基于特征向量的方法

此方法最先用于雷达的 DOA 估计，该方法已经用于宽带麦克风阵列信号处理，考虑 N 元含高斯白噪声的麦克风阵列。

#### 3.3.3.1 窄带 MUSIC (Multiple Signal Classification) 算法

MUSIC 最早用在雷达上的超分辨技术上，且常用在多声源定位中。MUSIC 算法是正交法，其基本思想是信号子空间与和噪声的子空间域进行正交，这样就能构造一个空间函数，再利用计算机技术搜索信号的谱峰值，最后检查语音信号的 DOA。对于第 n 个麦克的输出，其频域表示如下：

$$Y_n(f) = X_n(f) + V_n(f) = S(f) e^{-j2\pi[t+\tau_n]f} + V_n(f) \quad 3.3.3.1$$

可以定义如下频域向量：

$$\bar{y} = [Y_1(f) \ Y_2(f) \ \cdots \ Y_N(f)]^T \quad 3.3.3.2$$

这样可以用向量表示接收到的信号  $\bar{y}$ ：

$$\bar{y} = \bar{x} + \bar{v} = \varsigma(\tau) S(f) e^{-j2\pi\tau f} + \bar{v} \quad 3.3.3.3$$

此处可得：

$$\varsigma(\tau) = [e^{-j2\pi\tau_1 f} \ e^{-j2\pi\tau_2 f} \ \cdots \ e^{-j2\pi\tau_N f}]^T \quad 3.3.3.4$$

则输出的协方差可以写为如下形式：

$$\mathbf{R}_Y = E(\bar{y}\bar{y}^H) = \mathbf{R}_X + \sigma_V^2 \mathbf{I} \quad 3.3.3.5$$

其中

$\mathbf{R}_X = \sigma_s^2 \varsigma(\tau) \varsigma^H(\tau)$ ， $\sigma_s^2 = E[|S(f)|^2]$ ， $\sigma_V^2 = E[|V_1(f)|^2] = \cdots = E[|V_N(f)|^2]$  使用特征向量分解矩阵：

$$\mathbf{R}_Y = \mathbf{B} \mathbf{\Lambda} \mathbf{B}^H \quad 3.3.3.6$$

这里：

$$\Lambda = \text{diag}[\lambda_{Y,1} \quad \lambda_{Y,2} \quad \cdots \quad \lambda_{Y,N}] = \text{diag}[\lambda_{X,1} \quad \sigma_V^2 \quad \cdots \quad \sigma_V^2] \quad 3.3.3.7$$

是矩阵  $R_Y$  的特征值矩阵， $\mathbf{B} = [\mathbf{b}_1 \quad \mathbf{b}_2 \quad \cdots \quad \mathbf{b}_N]$ ， $\mathbf{b}_n$  是特征值  $\lambda_{Y,n}$  对应的特征向量， $\lambda_{Y,1}$  是  $R_X$  的唯一非零正定特征值。即有：

$$\mathbf{R}_Y \mathbf{b}_n = \lambda_{Y,n} \mathbf{b}_n = \sigma_V^2 \mathbf{b}_n \quad 3.3.3.8$$

同时有：

$$\mathbf{R}_Y \mathbf{b}_n = \begin{bmatrix} \sigma_s^2 \zeta(\tau) \zeta(\tau)^H & \sigma_V^2 \mathbf{I} \end{bmatrix} \mathbf{b}_n \quad 3.3.3.9$$

由 3.3.3.8 和 3.3.3.9 可得：

$$\sigma_s^2 \zeta(\tau) \zeta(\tau)^H \mathbf{b}_n = \mathbf{0} \quad 3.3.3.10$$

这等价于：

$$\zeta(\tau)^H \mathbf{b}_n = \mathbf{0} \quad 3.3.3.11$$

这意味着  $R_Y$  的  $N-1$  个特征向量和实际的 TDOA 向量正交。假设代价函数如下：

$$J_{MUSIC}(p) = \frac{1}{\sum_{n=2}^N |\mathbf{b}_n^H \zeta(p)|^2} \quad 3.3.3.12$$

则有 TDOA 的估计值如下：

$$\hat{\tau} = \arg \max_p J_{MUSIC}(p) \quad 3.3.3.13$$

### 3.3.3.2 宽带 MUSIC 算法

窄带 MUSIC 算法并在语音中效果并不好，这是因为语音是非平稳信号，语音信号的中的频谱成分随着语音变换而变化，一个抵消非平稳的方法是将式 3.3.3.12 的代价函数在整个频谱上去求代价函数。

$$\mathbf{y}_{1:N}(k, p) = [y_1[k + \tau_1(p)] \quad y_2[k + \tau_2(p)] \quad \cdots \quad y_N[k + \tau_N(p)]] \quad 3.3.3.14$$

空域相关矩阵如下：

$$\mathbf{R}_a(p) = E[\mathbf{y}_{1:N}(k, p) \mathbf{y}_{1:N}^T(k, p)] = \mathbf{R}_s(p) + \sigma_v^2 \mathbf{I} \quad 3.3.3.15$$

原始信号的协方差如下：

$$\mathbf{R}_s(p) = \begin{bmatrix} \sigma_s^2 & r_{ss,12}(p, \tau) & \cdots & r_{ss,1N}(p, \tau) \\ r_{ss,21}(p, \tau) & \sigma_s^2 & \cdots & r_{ss,2N}(p, \tau) \\ \vdots & \vdots & \ddots & \vdots \\ r_{ss,N2}(p, \tau) & r_{ss,N2}(p, \tau) & \cdots & \sigma_s^2 \end{bmatrix} \quad 3.3.3.16$$



$$r_{ss,ij}(p, \tau) = E \left\{ s[k-t-\tau_i + p_i] s[k-t-\tau_j + p_j] \right\} \quad 3.3.3.17$$

如果  $\tau = p$  则有：

$$\mathbf{R}_s(p) = \sigma_s^2 \begin{bmatrix} 1 & 1 & \cdots & 1 \\ 1 & 1 & \cdots & 1 \\ \vdots & \vdots & \ddots & \vdots \\ 1 & 1 & \cdots & 1 \end{bmatrix} \quad 3.3.3.18$$

该矩阵的秩是 1，如果  $\tau \neq p$ ，则矩阵的秩取决于输入信号的自身的特性。如果输入信号是白化过程的信号，则输入信号的自相关矩阵  $\mathbf{R}_s(p) = \text{diag}[\sigma_s^2 \quad \sigma_s^2 \quad \cdots \quad \sigma_s^2]$ 。在这一特殊情况下， $\mathbf{R}_s(p)$  是满秩的。通常来说，当  $\tau \neq p$ ， $\mathbf{R}_s(p)$  是半正定的。对  $\mathbf{R}(p)$  和  $\mathbf{R}_s(p)$  进行特征值分解，假设  $\lambda_{s,1}(p) \geq \lambda_{s,2}(p) \geq \cdots \geq \lambda_{s,N}(p)$  是  $\mathbf{R}_s(p)$  的  $N$  个特征值。则  $\mathbf{R}(p)$  的  $N$  个特征向量可表示为如下：

$$\lambda_{y,n}(p) = \lambda_{s,n}(p) + \sigma_v^2 \quad 3.3.3.19$$

设  $\mathbf{b}_1(p), \mathbf{b}_2(p), \dots, \mathbf{b}_N(p)$  是  $\mathbf{R}(p)$  特征向量， $\mathbf{R}(p)$  是对称的拓普利兹矩阵，其所有特征向量都是实数组成的，则有：

$$\mathbf{R}(p)\mathbf{B}(p) = \mathbf{B}(p)\Lambda(p) \quad 3.3.3.20$$

此处：

$$\mathbf{B}(p) = [\mathbf{b}_1(p) \quad \mathbf{b}_2(p) \quad \cdots \quad \mathbf{b}_N(p)]$$

$$\Lambda(p) = \text{diag}[\lambda_{y,1}(p) \quad \lambda_{y,2}(p) \quad \cdots \quad \lambda_{y,N}(p)]$$

此处  $p = \tau$ ，由于  $\mathbf{R}_s(p)$  的秩等于 1，所以对于  $n \geq 2$  情况，有下式成立：

$$\mathbf{R}(\tau)\mathbf{b}_n(\tau) = [\mathbf{R}_s(\tau) + \sigma_v^2 \mathbf{I}] \mathbf{b}_n(\tau) = \sigma_v^2 \mathbf{b}_n(\tau) \quad 3.3.3.21$$

这意味着：

$$\mathbf{b}_n^T(p) \mathbf{R}(p) \mathbf{b}_n(p) = \begin{cases} \sigma_v^2, p = \tau \\ \lambda_{y,n(p)} \geq \sigma_v^2, p \neq \tau \end{cases} \quad 3.3.3.22$$

所以可以使用如下的代价函数：

$$J_{BMUSIC}(p) = \frac{1}{\sum_{n=2}^N \mathbf{b}_n^T(p) \mathbf{R}(p) \mathbf{b}_n(p)} \quad 3.3.3.23$$

该代价函数的峰值对于 TDOA 值  $\tau$ 。

### 3.3.4 最小熵法

基于高阶统计的 TDOA 估计法，熵是一个对随机或者不确定变量的一个统计度量，由香农在信息论中最先提出来，对于一个随机变量  $y$ ，其概率密度函数（PDF，probability density function）是  $p(y)$ ，其熵定义为：

$$H(y) = -\int p(y) \ln p(y) dy = -E[\ln p(y)] \quad 3.3.4.1$$

对于语音接收到的信号  $y(k, p)$  而言，其熵如下：

$$H[y(k, p)] = -\int p y(k, p) \ln p y(k, p) dy(k, p) \quad 3.3.4.2$$

其最小值（熵），对应的解  $p$  就是 TDOA：

$$\hat{\tau}^{ME} = \arg \min_p H[y(k, p)] \quad 3.3.4.3$$

ME 是 minimum entropy 的缩写。

#### 3.3.4.1 高斯信号源

如果信源是高斯类型的且麦克风输出没有噪声。假设麦克信号是零均值联合随机高斯信号，则它们的联合概率密度函数是：

$$p[y(k, p)] = \frac{e^{-\eta(k, p) / 2}}{\sqrt{(2\pi)^N \det[R(p)]}} \quad 3.3.4.4$$

此处：

$$\eta(k, p) = y^T(k, p) R^{-1}(p) y(k, p) \quad 3.3.4.5$$

将式 3.3.4.5 代入 3.3.4.2 得：

$$H[y(k, p)] = \frac{1}{2} \ln \left\{ (2\pi e)^N \det[R(p)] \right\} \quad 3.3.4.6$$

则 3.3.4.3 将变为：

$$\hat{\tau}^{ME} = \arg \min_p \det[R(p)] \quad 3.3.4.7$$

#### 3.3.4.2 语音信号源

虽然语音信号是复杂的随机信号，并没有明确的数学表达式表示其熵，但是语音信号可以用拉普拉斯模型来近似描述。单变量拉普拉斯零均值，方差为  $\sigma_y^2$  的分布如下：

$$p(y) = \frac{\sqrt{2}}{2\sigma_y} e^{-\sqrt{2}|y|/\sigma_y} \quad 3.3.4.8$$

其熵可表示如下：

$$H[y] = 1 + \ln(\sqrt{2}\sigma_y) \quad 3.3.4.9$$

对于零均值多变量协方差为  $\mathbf{R}(p)$  的拉普拉斯累计分布如下：

$$p[y(k, p)] = \frac{2 \left[ e^{-\eta(k, p)} / 2 \right]^{Q/2} K_Q \left[ \sqrt{2\eta(k, p)} \right]}{\sqrt{(2\pi)^N \det[\mathbf{R}(p)]}} \quad 3.3.4.10$$

这里：

$Q = (2 - N) / 2$ ， $K_Q(\cdot)$  是第三类修正贝塞尔函数：

$$K_Q(a) = \frac{1}{2} \left( \frac{a}{2} \right)^Q \int_0^\infty z^{-Q-1} e^{-z - \frac{a^2}{4z}} dz, a > 0 \quad 3.3.4.11$$

其联合熵函数如下：

$$H[y(k, p)] = \frac{1}{2} \ln \left\{ \frac{(2\pi)^N \det[\mathbf{R}(p)]}{4} \right\} - \frac{Q}{2} E \left\{ \ln \left[ \frac{\eta(k, p)}{2} \right] \right\} - E \left\{ \ln K_Q \left[ \sqrt{2\eta(k, p)} \right] \right\}$$

### 3.3.5 自适应特征向量分解法

该方法对于混响模型下，具有较好的健壮性，在加性噪声情况下：

$$y_1(k) * g_2 = x_1(k) * g_1 * g_2 = x_2(k) * g_1 = y_2(k) * g_1 \quad 3.3.4.12$$

则  $k$  时刻可得下式：

$$\mathbf{y}^T(k) \mathbf{w} = \mathbf{y}_1^T(k) \mathbf{g}_2 - \mathbf{y}_2^T(k) \mathbf{g}_1 = 0 \quad 3.3.4.13$$

将上式左乘以  $\mathbf{y}(k)$  可得：

$$\mathbf{R}_{yy} \mathbf{W} = \mathbf{0}_{2L \times 1} \quad 3.3.4.14$$

$\mathbf{W}$  是  $\mathbf{R}_{yy}$  对应的特征值 0 的特征向量，实际上噪声总是存在的， $\mathbf{R}_{yy}$  实际上是正定矩阵，所以， $\mathbf{W}$  是相对  $\mathbf{R}_{yy}$  于最小的特征值的特征向量。即有：

$$\hat{\mathbf{W}} = \arg \min_{\mathbf{w}} \mathbf{W}^T \mathbf{R}_{yy} \mathbf{W}, \|\mathbf{w}\| = 1 \quad 3.3.4.15$$

式 3.3.4.15 可以使用约束性 LMS 算法来求解：

初始化

$$\hat{\mathbf{g}}_n(0) = \begin{bmatrix} \frac{\sqrt{2}}{2} & 0 & \cdots & 0 \end{bmatrix}, n = 1, 2$$

$$\hat{\mathbf{w}}(0) = [\hat{\mathbf{g}}_2^T(0) - \hat{\mathbf{g}}_1^T(0)]^T$$

按下式计算：

$$e(k) = \hat{\mathbf{w}}^T(k) \mathbf{y}(k)$$

$$\hat{\mathbf{w}}(k+1) = \frac{\hat{\mathbf{w}}(k) - \mu e(k) \mathbf{y}(k)}{\|\hat{\mathbf{w}}(k) - \mu e(k) \mathbf{y}(k)\|}$$

在上述方法收敛后，

$$\hat{\tau} = \arg \max_l |\hat{g}_{1,l}| - \arg \max_l |\hat{g}_{2,l}|$$

### 3.3.6 自适应盲信号分离（BSS, Blind Source Separation）

独立分量分析法（ICA, Independent Component Analysis）是盲源分离采用最多的算法。

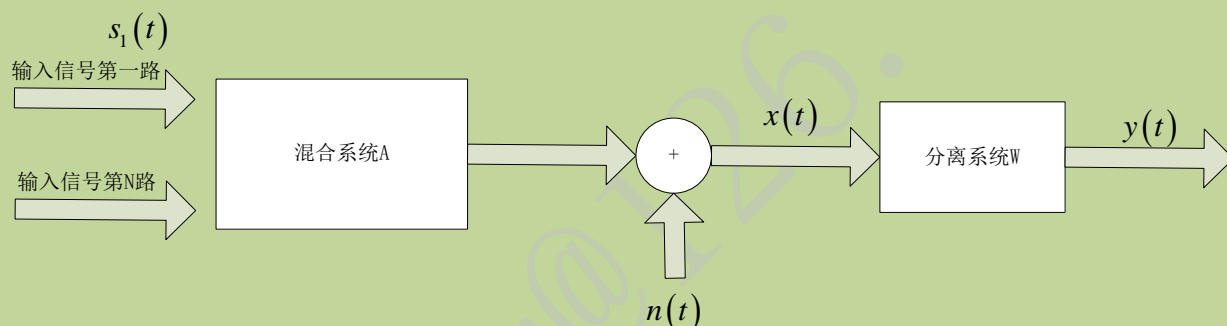


图 3.3.6.1 盲源分离原理

BSS 分离的数学模型如上图，

$\mathbf{s}(t) = [s_1(t) \cdots s_N(t)]$  是输入 N 维未知信源，A 为未知混合系统， $\mathbf{x}(t) = [x_1(t) \cdots x_M(t)]$  是观测到的信号向量， $\mathbf{n}(t) = [n_1(t) \cdots n_M(t)]^T$ ，盲源分离的目的就是在源信号 s 和混合系统 A 均未知前提下，调节分离系统 W，使得输出信号 y 是 s 的估计，即：

$$\mathbf{y}(t) = \mathbf{x} \cdot \mathbf{W} = \hat{\mathbf{s}} \quad 3.3.6.1$$

FastICA 算法

1. 对观测的数据 X 去均值，使其均值为 0，其目的是简化 ICA 算法。
2. 对数据进行白化  $\mathbf{X} \rightarrow \mathbf{Z}$ ，去除观测信号之间的相关性，简化后续独立分量的提取过程。
3. 更新准则，采用随机梯度法，定义自然梯度：

$$\mathbf{W}(k+1) = \mathbf{W}(k) - \theta(k) \frac{\partial J(\mathbf{W}(k))}{\partial \mathbf{W}(k)} \mathbf{W}^T(k) \mathbf{W}(k) \quad 3.3.6.2$$

根据信息最大或者互信息最大可以得出如下代价函数：

$$J(\mathbf{w}, y) = -\log(\det(\mathbf{w})) - \sum_{i=1}^n \log(p_i(y_i)) \quad 3.3.6.3$$

$p_i(y_i)$  是  $y_i$  信号的概率密度函数， $\det(\mathbf{W})$  是矩阵  $\mathbf{W}$  的行列式的值。将代价函数代入 3.3.6.2 得：

$$\Delta \mathbf{W}(k) = \theta(k) [I - f(\mathbf{y}(k) \mathbf{y}^T(k))] \mathbf{W}(k) \quad 3.3.6.4$$

FastICA:

1. 初始化非零权向量  $\mathbf{W}$
2. 向量叠加， $\mathbf{w}_p^+ = E\{xg(\mathbf{w}_p^T x)\} - E\{g'(\mathbf{w}_p^T x)\} \mathbf{w}_p$ ， $g$  是非线性函数， $g'$  是其一阶时域导数。
3. 归一化成酉空间向量

$$\mathbf{w}_p = \frac{\mathbf{w}_p^+}{\|\mathbf{w}_p^+\|}$$

4. 根据施密特正交，解其和之前向量的相关性。

$$\mathbf{w}_p = \mathbf{w}_p - \sum_{j=1}^{p-1} \mathbf{w}_p^T \mathbf{w}_j \mathbf{w}_j; \mathbf{w}_p = \frac{\mathbf{w}_p}{\|\mathbf{w}_p\|} \quad 3.3.6.5$$

5. 如果  $\mathbf{W}$  还是没收敛，则返回步骤一。

### 3.3.1 TDOA

一般实时系统使用 TDOA (time difference of arrival) 技术，该技术主要是估计声源的入射角，该方法的基本步骤是：首先根据各麦克风之间的相对延迟，然后利用相对之间的相位延迟估计声源的位置。在 3.1.1 节中，在入射角是  $\theta$  的情况下，信号到达麦克风之间的延迟  $\tau = \frac{d \sin \theta}{c}$ ，则将对  $\theta$  的估计转换为对  $\tau$  的估计。由于噪声和混响的存在，精确估计延迟比较困难，常采用的时延估计法主要有广义互相关 (GCC, Generalized Cross Correlation) 和 LMS 自适应滤波器法。考虑两麦克风接收信号为：

$$y_1(n) = s_1(n) + n_1(n)$$

$$y_2(n) = s_2(n + \tau) + n_2(n)$$

其中  $\tau$  是时延,  $s_1(n)$  和  $s_2(n + \tau)$  是语音信号,  $n_1(n)$  和  $n_2(n)$  是背景噪声, 它们的互相关函数为:

$$R_{12}(p) = F^{-1}[\psi_{y_1 y_2}(f)] = \int_{-\infty}^{+\infty} \psi_{y_1 y_2}(f) e^{j2\pi fp} df = \int_{-\infty}^{+\infty} \phi_{y_1 y_2}(f) e^{j2\pi fp} df \quad 3.3.1.1$$

其中  $\phi_{y_1 y_2}(f) = E[Y_1(f)Y_2^*(f)]$ ,  $Y_n(f) = \sum_k y_n(k) e^{-j2\pi fk}$ ,  $n=1,2$ ,  $\tau$  为延迟,  $f$  为频率。

该麦克风对的时延就是函数  $R_{12}(\tau)$  的峰值对应的  $\tau$  值。通常为了使时延  $\tau$  估计的更准确, 通常会采用加权的方式进行计算:

$$R_{12}(p) = \int_{-\infty}^{+\infty} W_{12} Y_1(f) Y_2^*(f) e^{j2\pi fp} df \quad 3.3.1.2$$

$W_{12}$  是广义互相关加权函数。

#### ● 平滑相关变换

为了克服语音信号在估计 TDOA 时的波动影响, 一个有效的方法是在计算麦克风之间的互相关前白化麦克风输出, 即:

$$W_{12} = \frac{1}{\sqrt{E[|Y_1|^2] E[|Y_2|^2]}} \quad 3.3.1.3$$

则可得:

$$\begin{aligned} \psi_{y_1 y_2}(f) &= \frac{\alpha_1 \alpha_2 e^{-j2\pi fp} E[|S(f)|^2]}{\sqrt{E[|Y_1(f)|^2] E[|Y_2(f)|^2]}} \\ &= \frac{\alpha_1 \alpha_2 e^{-j2\pi fp} E[|S(f)|^2]}{\sqrt{\alpha_1^2 E[|s(f)|^2] + \sigma_{v_1}^2(f)} \cdot \sqrt{\alpha_2^2 E[|s(f)|^2] + \sigma_{v_2}^2(f)}} \\ &= \frac{e^{-j2\pi fp}}{\sqrt{1 + \frac{1}{SNR_1(f)}} \cdot \sqrt{1 + \frac{1}{SNR_2(f)}}} \end{aligned} \quad 3.1.3.4$$

可得到信噪比如下:

$$SNR_n(f) = \frac{\alpha_n^2 E[|S(f)|^2]}{E[|V_n(f)|^2]}, n=1,2 \quad 3.1.3.5$$

#### ● 相位变化法

由式 3.3.1.2 可知, TDOA 实际上是相位的估计, 所以可以将  $W_{12}$  设置为如下形式:

$$W_{12} = \frac{1}{|\phi_{y_1 y_2}(f)|} \quad 3.1.3.6$$

则有：

$$\psi_{y_1 y_2}(f) = e^{-j2\pi f \tau} \quad 3.1.3.7$$

这样可以求得：

$$R_{12}(p) = \int_{-\infty}^{+\infty} e^{j2\pi f(p-\tau)} df = \begin{cases} \infty, p = \tau \\ 0, otherwise \end{cases} \quad 3.1.3.8$$

### 3.3.3 空域线性预测法

其基本思想是使用冗余的麦克风来增加 TDOA 估计的准确度。假设有三个麦克风，它们之间的延迟分别是  $\tau_{12}$   $\tau_{13}$  和  $\tau_{23}$ ，它们之间的关系是  $\tau_{13} = \tau_{12} + \tau_{23}$ ，考虑 N 元线阵远场语音模型：

$$y_n(k) = \alpha_n s(k - t - \tau_{n1}) + v_n(k), 1 \leq n \leq N \quad 3.1.3.10$$

忽略噪声的影响，则有以下式成立：

$\alpha_n s(k - t) = y_n(k + \tau_{n1}), 1 \leq n \leq N$ ，如果取  $y_1(k)$  和  $y_n(k + \tau_{n1})$  对齐，则可得前向预测误差：

$$e_1(k, p) = y_1(k) - \mathbf{y}_{a,2:N}^T(k, p) \boldsymbol{\alpha}_{2:N}(p) \quad 3.1.3.11$$

其中  $p$  代表了时延。

$$\mathbf{y}_{a,2:N}(k, p) = [y_2(k + p_{21}) \cdots y_N(k + p_{N1})]^T \quad 3.1.3.12$$

根据最小均方误差准则，可以求得延迟  $p$ 。

$$J_1(p) = E[e_1^2(k, p)] \quad 3.1.3.13$$

$$R_{2:N}(p) \boldsymbol{\alpha}_{2:N(p)} = r_{2:N(p)} \quad 3.1.3.14$$

此处：

$$R_{2:N}(p) = E[\mathbf{y}_{2:N}(k, p) \mathbf{y}_{2:N}^T(k, p)] = \begin{bmatrix} \sigma_{y_2}^2 & r_{y_2 y_3}(p) & \cdots & r_{y_2 y_N}(p) \\ r_{y_3 y_2}(p) & \sigma_{y_3}^2 & \cdots & r_{y_3 y_N}(p) \\ \vdots & \vdots & \ddots & \vdots \\ r_{y_N y_2}(p) & r_{y_N y_3}(p) & \cdots & \sigma_{y_N}^2 \end{bmatrix} \quad 3.1.3.15$$

是空域自相关矩阵。

$$\sigma_{y_n}^2 = E[y_n^2(k)], 1 \leq n \leq N \quad 3.1.3.16$$

$$r_{y_i y_j}(p) = E \left\{ y_i \left[ k + \tau_i(p) \right] y_j \left[ k + \tau_j(p) \right] \right\}, i, j = 1, 2 \cdots N \quad 3.1.3.17$$

并且：

$$\mathbf{r}_{2:N}(p) = E \left\{ r_{y_1 y_2}(p) \quad r_{y_1 y_3}(p) \quad \cdots \quad r_{y_1 y_N}(p) \right\}, i, j = 1, 2 \cdots N \quad 3.1.3.18$$

将上式代入 3.1.3.14 得：

$$e_{1,\min}(k, p) = y_1(k) - \mathbf{r}_{2:N}^T(p) \mathbf{R}_{2:N}^T(p) \mathbf{r}_{2:N}(p)$$

可以得：

$$J_{1,\min}(p) = E \left\{ e_{1,\min}^2(k, p) \right\} = \sigma_{y_1}^2 - \mathbf{r}_{2:N}^T(p) \mathbf{R}_{2:N}^{-1}(p) \mathbf{r}_{2:N}(p)$$

然后可得：

$\hat{\tau} = \arg \min_p J_{1,\min}(p)$ ，如果只有两个麦克风，则退化为互相关算法，当麦克风数量增加时，该算法将利用冗余信息提高 TDOA 的估计值。

当信噪比达到 10dB 时，两个麦克风可以达到准确的 TDOA 估计，当信噪比到 -5dB 时，基于两个麦克风的 TDOA 的估计将不准确。但是当使用 4 或者更多的麦克风时，系统的结果将是正确的。

### 3.3.

#### 3.3.2 SRP-PHAT

SRP-PHAT (steering response power-phase transform) 可控功率响应和相位变换技术，基于粒子滤波需要的运算量太大，使用 SRC (stochastic region contraction) 随机区域收缩方法简化计算量。SRP-PHAT 方法的波束输出功率 (帕斯瓦尔定律) 为：

$$P(\bar{q}) = \sum_{m=1}^M \sum_{l=1}^M \int_{-\infty}^{+\infty} \psi_{ml}(k) X_m(k) X_l^*(k) e^{j2\pi k(\tau_1 - \tau_m)} dk$$

其中  $\tau_1$  和  $\tau_m$  为阵列的可控时延， $\bar{q}$  为声源的空间位置矢量， $X_m(k)$  为第  $m$  个麦克风信号  $x_m(t)$  的加窗傅里叶变换。上式即第  $m$  和第 1 个麦克风的交叉功率谱，在多信道下 PHAT 加权系数为：

$$\psi_{ml}(k) = \frac{1}{\left[ X_m(k) X_l^*(k) \right]}$$

设  $R_{ml}(\tau)$  为第  $m$  和第 1 个麦克风接收到信号的 PHAT 加权广义互相关函数，则 SRP-PHAT 算法的时域表示为：

$$P(\bar{q}) = \sum_{m=1}^M \sum_{l=1}^M \int_{-\infty}^{+\infty} R_{ml}(\tau_1(\bar{q}) - \tau_m(\bar{q}))$$



该定位方法就是在所有可能覆盖声源的位置中， $P(\bar{q})$  最大值所对应的  $\bar{q}$  值。也就是

$$\hat{\bar{q}} = \arg \max_{\bar{q}} P(\bar{q})$$

随机区域收敛技术 SRC (stochastic region contraction)

其基本思想是：给定一个起始的长方体作为搜索区域，并且假定全局最优解（说话人的位置）就在这个长方体里，在求解最优解的过程中通过迭代的方式缩小长方体的体积直到体积小到的范围。

假设  $V_{peak}$  是要缩小的范围， $V_{room}$  是起始考察的范围。则为了使其 miss 概率（左栏）小于该栏时计算的次数列在了对应的表中。人脸的面积约 330~400 平方厘米，实际按 800 平方厘米作为搜索区域，则 8000 平方分米（8 平方米，3\*3 米）。

$\frac{v_{peak}}{v_{room}}   P_{(miss v_{peak})}$	0.1	0.01	0.001	0.0001
1%	44	459	4063	46050
0.1%	66	688	6905	69,075
0.001%	88	917	9206	92099

表：计算  $\hat{P}(\bar{q})$  的次数

定义  $J_i$  是第  $i$  次迭代的点数； $N_i$  是定义新空间的点数； $V_{i+1}$  有一个长方体边界向量  $\bar{B}_{i+1} \equiv [x_{\max}(i+1)x_{\min}(i+1)y_{\max}(i+1)y_{\min}(i+1)z_{\max}(i+1)z_{\min}(i+1)]$ ， $I$  迭代的次数， $FE_i$  是第  $i$  次迭代的所有 fe (functional evaluation, 任一点的  $\hat{P}(\bar{q})$ ) 的估计次数之和， $\Phi$  是 fe 的最大个数，则 SRC 算法如下：

1.  $i = 0$

2. 初始化  $J_0$ 、 $N_0$  以及  $V_0 = V_{room}$

3. 对  $J_i$  点计算  $\hat{P}(\bar{q})$

4. 排序，找出最优的  $N_i \leq J_i$  点

5 将原始区域缩小为  $V_{i+1}$ ， $\bar{B}_{i+1}$  包括这些  $N_i$  个点。

6 IF:  $V_{i+1} < V_u$ ，或者  $FE_i > \Phi$  并且  $V_{i+1} < T_1 V_u$ ， $T_1$  是一个大约等于 10 的参数， $V_u$  是目标区域大小；这时求  $\hat{x}_s^n(i^*)$ ， $I = i$ ，Return x

7 ELSE IF:  $FE_i > \Phi$ ，stop，return NULL

8 ELSE: 在  $N_i$  点中选取大于等于均值  $\mu_i$  ( $N_i$  的均值) 的子集  $G_i$

9 在  $V_{i+1}$  中评估  $J_{i+1}$

10 将  $N_{i+1}$  点作为  $G_i$  和刚刚计算的  $J_{i+1}$  (源于  $N_{i+1} - G_i$  空间) 的集合,

11  $i = i + 1$  , 跳至第 5 步。

计算开销:

对于  $M$  个麦克风,  $\hat{P}(\bar{q})$  的计算次数  $Q = M(M-1)/2$  相位变换, 对于 DFT 的长度是  $L$

则有: DFT:  $M * \frac{5}{2} L \log_2^L$  , 频谱计算: 约  $10QL$ , IDFT:

语音阵列处理的其它问题

说话人数量估计

背景噪声下信号分离和盲信号分离

生理语音学

## 下篇-语音识别

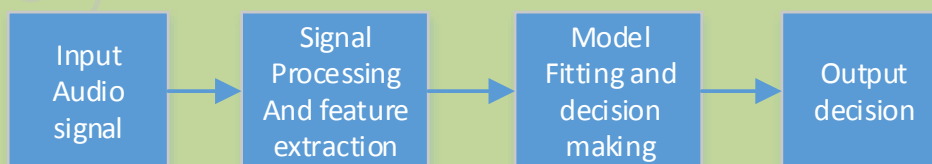
### 语音信号预加重

语音信号随着频谱增加幅度值越小, 对于频率相差两边的语音信号, 其功率谱差 6dB; 需要对高频信号进行预加重, 一般是将语音信号通过一个高通滤波器  $1-aZ^{-1}$ , 即预加重滤波器,

其目的是滤除 50Hz, 60Hz 的工频干扰。

### VAD 算法

当前 VAD 采用的算法依赖于以下几种特征: 基于能量变换, 周期性, 频谱差异。其可以用在语音增强, 语音编码和语音识别。VAD 算法模型有基于神经元模型和基于统计模型。比较流行的统计模型是高斯分布模型和拉普拉斯分布模型。



VAD 处理框图

## 第五章 ASR (automatic speech recognition)

语音识别问题就是模式分类问题。将音素以及音素组用离散的类来模拟。语音识别的目标是预测正确的类序列。如果  $\mathbf{z}$  表示从声波提取的特征向量序列，则语音识别系统根据最优分类方程来工作：

$$\hat{\omega} = \arg \max_{\omega \in W} P(\omega | \mathbf{z}) \tag{5.1}$$

$\hat{\omega}$  是识别系统判断的说话的词序列， $W$  是所有可能的说话词的序列集合。实际上用贝叶斯准则来计算该值。

$$\hat{\omega} = \arg \max_{\omega \in W} \frac{P(\mathbf{Z} | \omega)P(\omega)}{P(\mathbf{Z})} \tag{5.2}$$

$P(\mathbf{Z} | \omega)$  是声学似然（声学打分），代表了在词  $\omega$  被说了的情况下，语音序列  $\mathbf{Z}$  出现的概率。 $P(\omega)$  是语言打分，是语音序列出现的先验概率，其计算依赖于语言模型。在忽略语音序列出现概率的情况下，式 5.2 可以简化为：

$$\hat{\omega} = \arg \max_{\omega \in W} P(\mathbf{Z} | \omega)P(\omega) \tag{5.3}$$

这样语音识别可以分为两个主要步骤，特征提取和解码。

ASR 主要包括四个部分：信号处理和特征提取，声学模型（AM acoustic model），语言模型（LM language model），假设搜索（hypothesis search）。

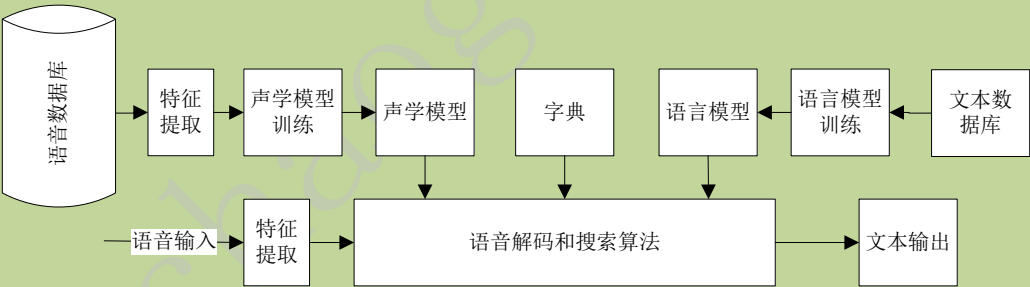


图 5.1 连续语音识别框图

在任意一个 ASR（automatic speech recognition）系统中第一步是特征提取，声道的形状（舌头，牙齿）在语音短时功率谱的包络中显示出来。而 MFCC（Mel Frequency Cepstral Coefficient，梅尔频率倒谱系数）就是一种准确描述该包络的一种特征向量。目前业界广泛使用的是基于 deep Learning 的特征学习，确切的说是深度神经网络（Deep Neural Network（DNN））。

## 5.1 ASR 模型

### 5.1.1 高斯混合模型 GMM(Gaussian mixture model)

对于说话人识别，语言特性降噪以及语音识别适用。

若随机变量  $X$  服从一个位置参数为  $\mu$ 、尺度参数为  $\sigma$  的概率分布，其概率密度函数

是：

$$f(x) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-\mu)^2}{2\sigma^2}} \quad 5.4$$

则称该变量是高斯分布(正态分布)的。记作  $X \sim N(\mu, \sigma^2)$ ；正态随机向量  $\mathbf{x} = (x_1, x_2, \dots, x_D)^T$  的高斯分布是：

$$f(\mathbf{x}) = \frac{1}{\sqrt{2\pi}^D \sqrt{|\Sigma|}} e^{-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \Sigma^{-1}(\mathbf{x} - \boldsymbol{\mu})} \quad 5.5$$

也常记作： $X \sim N(\boldsymbol{\mu} \in R^D, \Sigma \in R^{D \times D})$ ，一个连续标量  $x$  的混合高斯分布的概率密度函数是：

$$f(x) = \sum_{m=1}^M \frac{c_m}{\sqrt{2\pi}\sigma_m} e^{-\frac{1}{2}\left(\frac{x-\mu_m}{\sigma_m}\right)^2} = \sum_{m=1}^M c_m N(x; \mu_m, \sigma_m^2), (-\infty < x < +\infty; \sigma_m > 0; c_m > 0) \quad 5.6$$

混合权重的和等于一，即  $\sum_{m=1}^M c_m = 1$ 。和单峰高斯分布相比，5.6 是一个具有多个峰值的分布（混合高斯分布），体现在  $M > 1$ 。

混合高斯分布（5.6）随机变量  $x$  的期望是  $E(x) = \sum_{m=1}^M c_m \mu_m$ 。

多元混合高斯分布的联合概率密度函数是：

$$f(\mathbf{x}) = \sum_{m=1}^M \frac{c_m}{\sqrt{2\pi}^D \sqrt{|\Sigma_m|}} e^{-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_m)^T \Sigma_m^{-1}(\mathbf{x} - \boldsymbol{\mu}_m)} = \sum_{m=1}^M c_m N(\mathbf{x}; \boldsymbol{\mu}_m, \Sigma_m), (c_m > 0) \quad 5.7$$

### 5.1.2 参数估计

$$\Theta = \{c_m, \mu_m, \Sigma_m\}$$

使用最大似然估计法估计混合高斯分布的参数，

$$c_m^{(j+1)} = \frac{1}{N} \sum_{t=1}^N h_m^{(j)}(t) \quad 5.8$$

$$\mu_m^{(j+1)} = \frac{\sum_{t=1}^N h_m^{(j)}(t) x^{(t)}}{\sum_{t=1}^N h_m^{(j)}(t)} \quad 5.9$$

$$\Sigma_m^{(j+1)} = \frac{\sum_{t=1}^N h_m^{(j)}(t) [x^{(t)} - \mu_m^{(j)}][x^{(t)} - \mu_m^{(j)}]^T}{\sum_{t=1}^N h_m^{(j)}(t)} \quad 5.10$$

后验概率的计算如下：

$$h_m^j(t) = \frac{c_m^{(j)} N(x^{(t)}; \mu_m^{(j)}, \Sigma_m^{(j)})}{\sum_{i=1}^n c_i^{(j)} N(x^{(t)}; \mu_i^{(j)}, \Sigma_i^{(j)})} \quad 5.11$$

基于当前（第  $j$  次）的参数估计， $x^{(t)}$  的条件概率取决于每一个的采样。

### 5.1.2 隐马尔可夫模型 HMM (hidden Markov model)

#### 马尔科夫链

$X_1, X_2, X_3, \dots$  马尔科夫链 (Markov Chain)，描述一种状态序列，其每个状态值取决于前面有限个状态。马尔科夫链是具有马尔科夫性质的随机变量的一个数列。这些变量的范围，就是它们所有可能的集合（又称状态空间）， $X_n$  是在状态  $n$  的状态。如果  $X_{n+1}$  对于过去状态的条件概率分布仅是  $X_n$  的一个函数，则：

$$P(X_{n+1} = x \mid X_1 = x_1, X_2 = x_2, \dots, X_n = x_n) = P(X_{n+1} = x \mid X_n = x_n)$$

这里  $x$  为过程中的某个状态。

隐马尔可夫模型是统计模型，其被用来描述一个含有隐含位置参数的马尔科夫过程。其难点是从可观察的参数中确定该过程的隐含参数。然后利用这些参数来作进一步的分析，例如模式识别。

## 5.2 MFCC

MFCC 将语音被分为很多帧，每帧语音对应于一个频谱（通过短时 FFT 计算），频谱表示频率与能量的关系。在实际使用中，频谱图有三种，即线性振幅谱、对数振幅谱、自功率谱（对数振幅谱中各谱线的振幅都作了对数计算，所以其纵坐标的单位是 dB（分贝））。这个变换的目的是使那些振幅较低的成分相对高振幅成分得以拉高，以便观察掩盖在低幅噪声中的周期信号）。

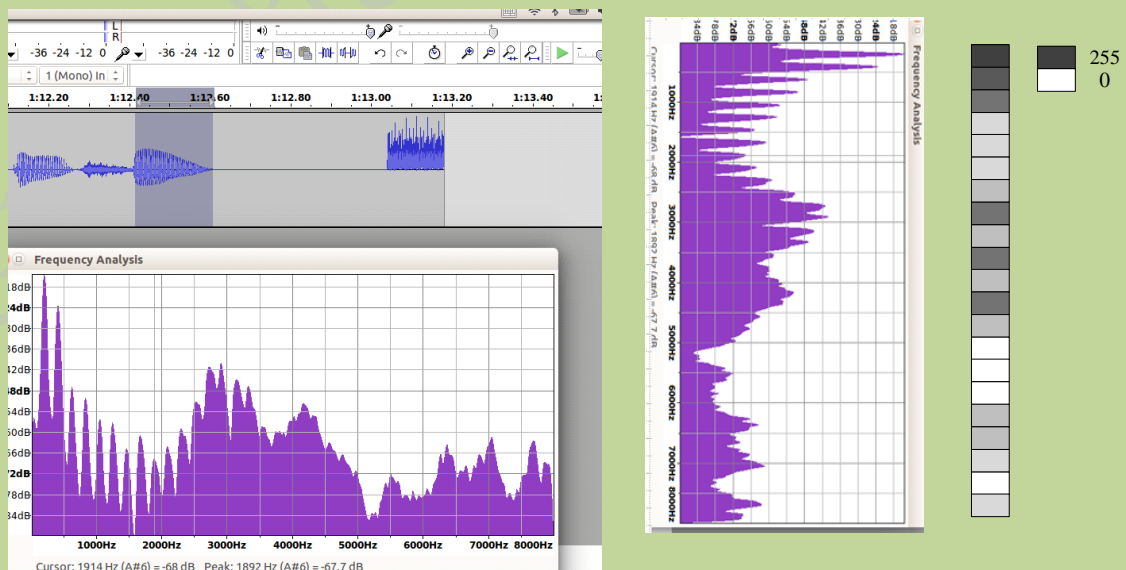


图 5.2 一帧信号的频谱变换

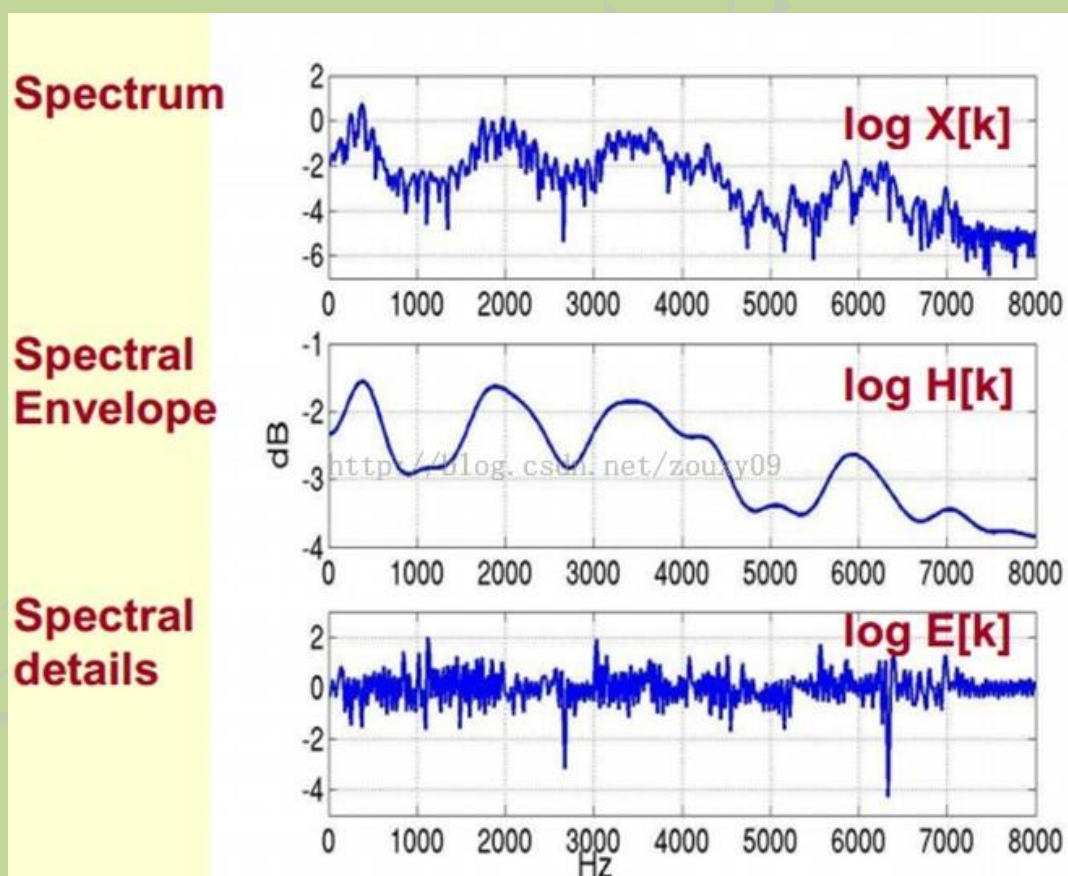
先将其中一帧语音（蓝色）的频谱（红色）通过坐标表示出来，如上图。现在将频谱旋

转 90 度。然后把这些幅度映射到一个灰度级表示，255 表示黑，0 表示白色。幅度值越大，相应的区域越黑。这样我们会得到一个随着时间变化的频谱图，这个就是描述语音信号的 spectrogram 声谱图。

音素（Phones）的属性可以更好的在声谱图里观察出来。另外，通过观察共振峰和它们的转变可以更好的识别声音。隐马尔科夫模型（Hidden Markov Models）就是隐含地对声谱图进行建模以达到好的识别性能。还有一个作用就是它可以直观的评估 TTS 系统（text to speech）的好坏，直接对比合成的语音和自然的语音声谱图的匹配度即可。图 5.2 是中文“区”这个字的时域和频域语音的频谱图。峰值就表示语音的主要频率成分，这些峰值称为共振峰（formants），而共振峰就是携带了声音的辨识属性。用它就可以识别不同的声音。

要提取的不仅仅是共振峰的位置，还得提取它们转变的过程。所以提取的是频谱的包络（Spectral Envelope）。这包络就是一条连接这些共振峰点的平滑曲线。

可以这么理解，原始的频谱由两部分组成：包络和频谱的细节。这里用到的是对数频谱，所以单位是 dB。那现在我们需要把这两部分分开，这样我们就可以得到包络了。

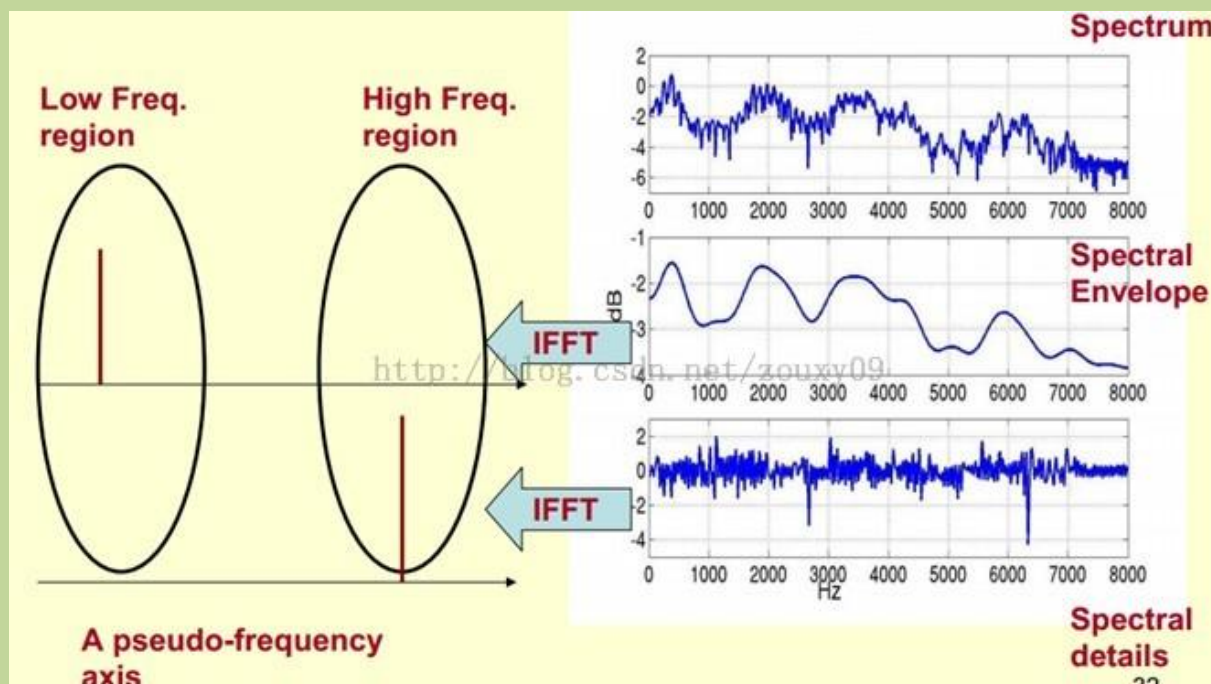


那怎么把他们分离开呢？也就是，怎么在给定  $\log X[k]$  的基础上，求得  $\log H[k]$  和  $\log E[k]$  以满足  $\log X[k] = \log H[k] + \log E[k]$ 。

为了达到这个目标，需要 Play a Mathematical Trick。这个 Trick 是什么呢？就是对频谱做 FFT。在频谱上做傅里叶变换就相当于逆傅里叶变换 Inverse FFT (IFFT)。需要注意的一点

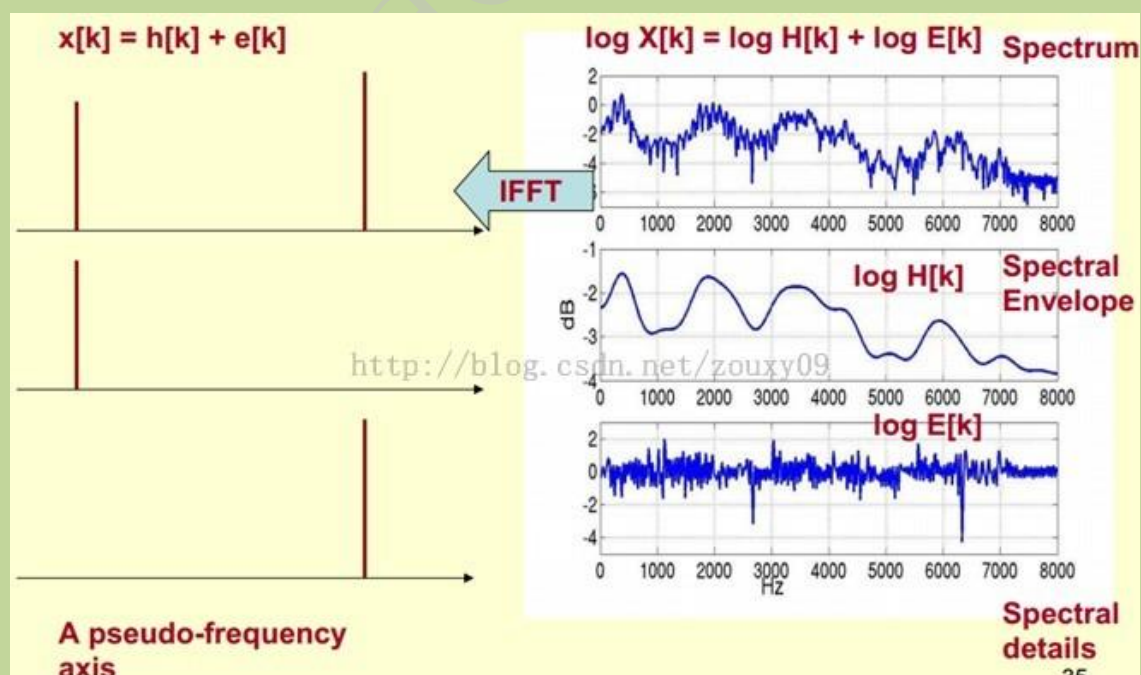


是，是在频谱的对数域上面处理的，这也属于 Trick 的一部分。这时候，在对数频谱上面做 IFFT 就相当于在一个伪频率（pseudo-frequency）坐标轴上面描述信号。



由上面这个图可以看到，包络是主要是低频成分（这时候需要转变思维，这时候的横轴就不要看成是频率了，可以看成时间），把它看成是一个每秒 4 个周期的正弦信号。这样在伪坐标轴上面的 4Hz 的地方给它一个峰值。而频谱的细节部分主要是高频。把它看成是一个每秒 100 个周期的正弦信号。这样在伪坐标轴上面的 100Hz 的地方给它一个峰值。

把它俩叠加起来就是原来的频谱信号了。



在实际中已经知道  $\log X[k]$ ，所以可以得到了  $x[k]$ 。那么由图可以知道， $h[k]$  是  $x[k]$  的低频部分，那么将  $x[k]$  通过一个低通滤波器就可以得到  $h[k]$  了。

$x[k]$ 实际上就是倒谱 Cepstrum（这个是一个新造出来的词，把频谱的单词 spectrum 的前面四个字母顺序倒过来就是倒谱的单词了）。 $h[k]$ 就是倒谱的低频部分。 $h[k]$ 描述了频谱的包络，它在语音识别中被广泛用于描述特征。

现在总结下倒谱分析，它实际上是这样一个过程：

1) 将原语音信号经过傅里叶变换得到频谱： $X[k]=H[k]E[k]$ ；

只考虑幅度就是： $|X[k]|=|H[k]||E[k]|$ ；

2) 我们在两边取对数： $\log|X[k]|=\log|H[k]|+\log|E[k]|$ 。

3) 再在两边取逆傅里叶变换得到： $x[k]=h[k]+e[k]$ 。

这实际上有个专业的名字叫做同态信号处理。它的目的是将非线性问题转化为线性问题的处理方法。对应上面，原来的语音信号实际上是一个卷性信号（声道相当于一个线性时不变系统，声音的产生可以理解为一个激励通过这个系统），第一步通过卷积将其变成了乘性信号（时域的卷积相当于频域的乘积）。第二步通过取对数将乘性信号转化为加性信号，第三步进行逆变换，使其恢复为卷性信号。这时候，虽然前后均是时域序列，但它们所处的离散时域显然不同，所以后者称为倒谱频域。

总结下，倒谱（cepstrum）就是一种信号的傅里叶变换经对数运算后再进行傅里叶反变换得到的谱。它的计算过程如下：

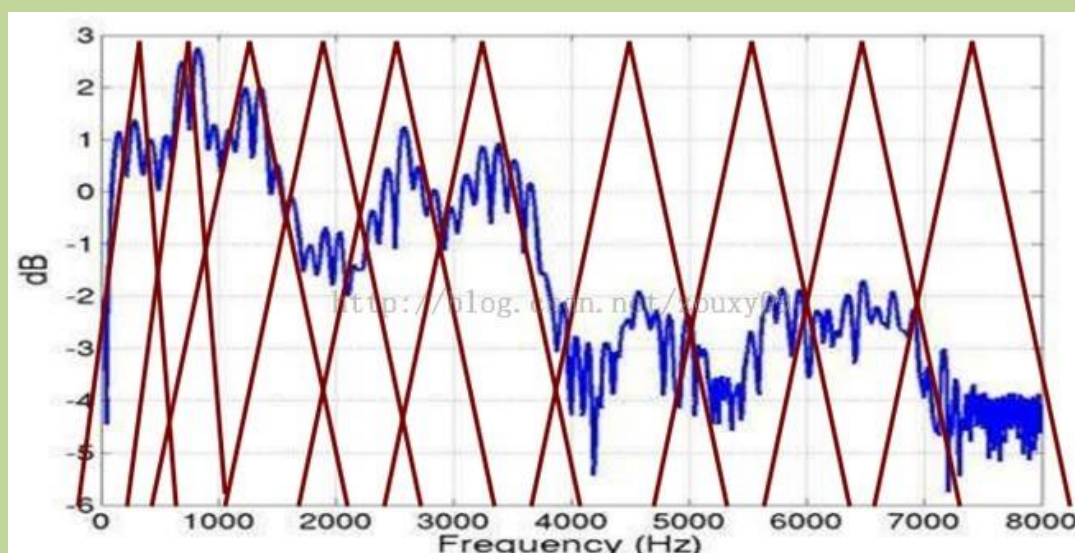


### 5.2.1 Mel 频率分析 (Mel-Frequency Analysis)

给一段语音，可以得到了它的频谱包络（连接所有共振峰值点的平滑曲线）了。但是，对于人类听觉感知的实验表明，人类听觉的感知只聚焦在某些特定的区域，而不是整个频谱包络。

而 Mel 频率分析就是基于人类听觉感知实验的。实验观测发现人耳就像一个滤波器组一样，它只关注某些特定的频率分量（人的听觉对频率是有选择性的）。也就是说，它只让某些频率的信号通过，而压根就直接无视它不想感知的某些频率信号。但是这些滤波器在频率坐标轴上却不是统一分布的，在低频区域有很多的滤波器，他们分布比较密集，但在高频区域，滤波器的数目就变得比较少，分布很稀疏。





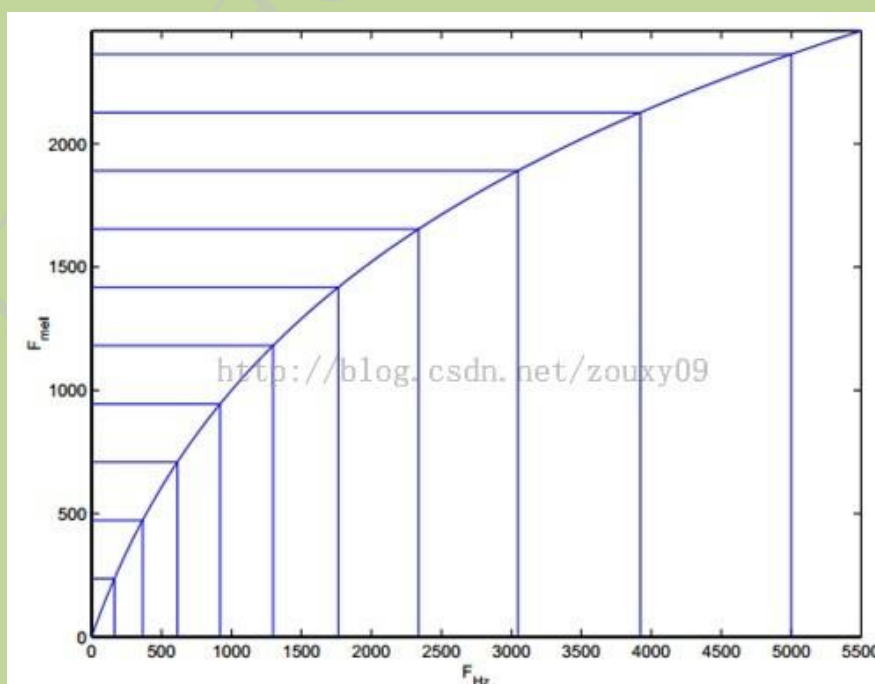
人的听觉系统是一个特殊的非线性系统，它响应不同频率信号的灵敏度是不同的。在语音特征的提取上，人类听觉系统做得非常好，它不仅能提取出语义信息，而且能提取出说话人的个人特征，这些都是现有的语音识别系统所望尘莫及的。如果在语音识别系统中能模拟人类听觉感知处理特点，就有可能提高语音的识别率。

梅尔频率倒谱系数（Mel Frequency Cepstrum Coefficient, MFCC）考虑到了人类的听觉特征，先将线性频谱映射到基于听觉感知的 Mel 非线性频谱中，然后转换到倒谱上。

将普通频率转化到 Mel 频率的公式是：

$$mel(f) = 2595 * \log_{10}(1 + f / 700)$$

由下图可以看到，它可以将不统一的频率转化为统一的频率，也就是统一的滤波器组。



在 Mel 频域内，人对音调的感知度为线性关系。举例来说，如果两段语音的 Mel 频率相差两倍，则人耳听起来两者的音调也相差两倍。

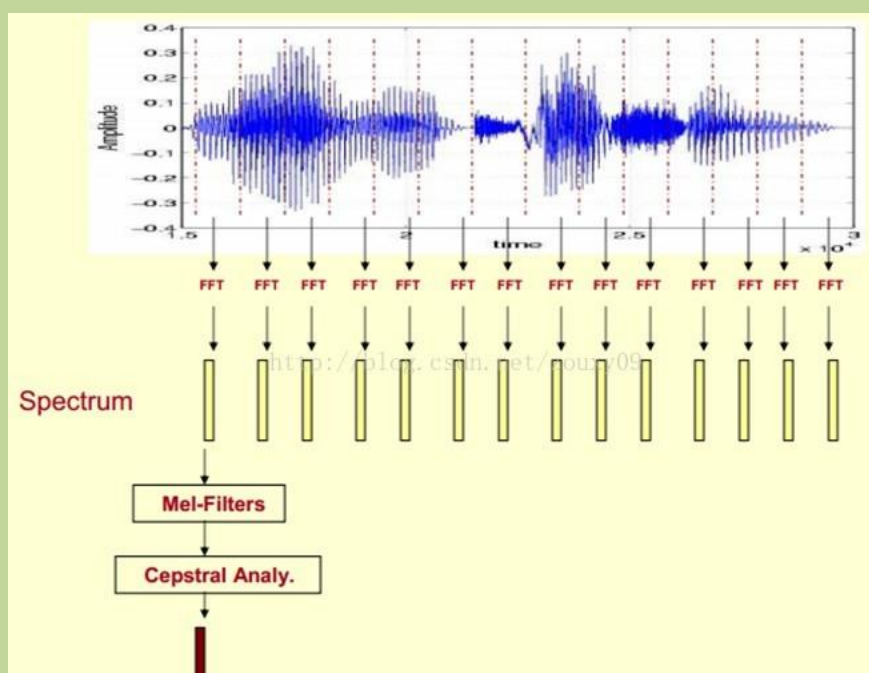
### 5.2.2 Mel 频率倒谱系数 (Mel-Frequency Cepstral Coefficients)

将频谱通过一组 Mel 滤波器就得到 Mel 频谱。公式表述就是： $\log X[k] = \log (\text{Mel-Spectrum})$ 。这时候我们在  $\log X[k]$  上进行倒谱分析：

1) 取对数： $\log X[k] = \log H[k] + \log E[k]$ 。

2) 进行逆变换： $x[k] = h[k] + e[k]$ 。

在 Mel 频谱上面获得的倒谱系数  $h[k]$  就称为 Mel 频率倒谱系数，简称 MFCC。



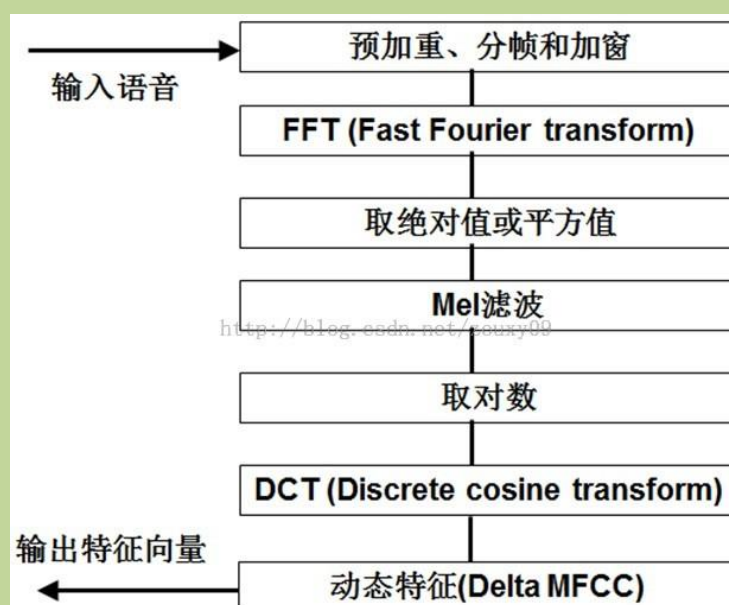
总结下提取 MFCC 特征的过程：（具体的数学过程网上太多了，这里就不想贴了）

1) 先对语音进行预加重、分帧和加窗；

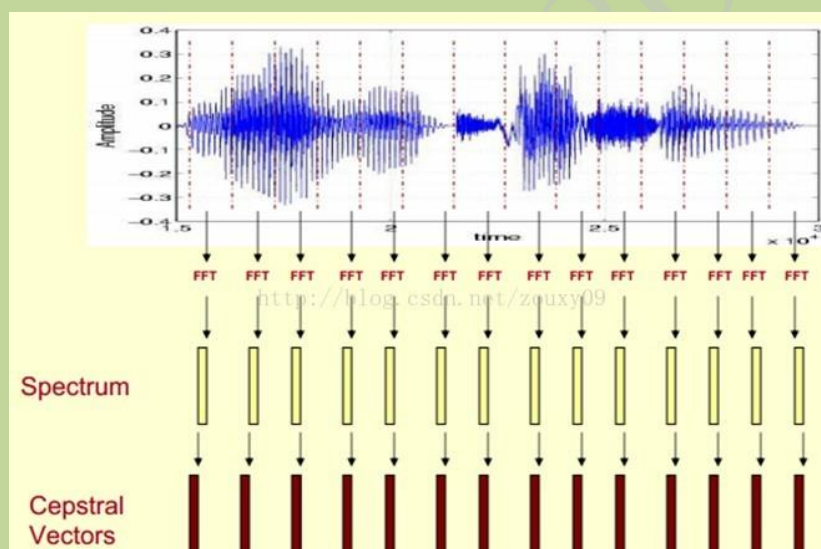
2) 对每一个短时分析窗，通过 FFT 得到对应的频谱；

3) 将上面的频谱通过 Mel 滤波器组得到 Mel 频谱；

4) 在 Mel 频谱上面进行倒谱分析（取对数，做逆变换，实际逆变换一般是通过 DCT 离散余弦变换来实现，取 DCT 后的第 2 个到第 13 个系数作为 MFCC 系数），获得 Mel 频率倒谱系数 MFCC，这个 MFCC 就是这帧语音的特征；



这时候，语音就可以通过一系列的倒谱向量来描述了，每个向量就是每帧的 MFCC 特征向量。



这样就可以通过这些倒谱向量对语音分类器进行训练和识别了。

## DNN（Deep neutral network）深度神经网络

### 15.1 深度神经元网络架构

深度神经元网络是传统的多层感知系统（MLP，multiplayer perceptron）。式 15.1 是一个五层神经元网络，一个输入层，三个隐藏层和一个输出层。

$$\mathbf{V}^l = f(\mathbf{z}^l) = f(\mathbf{W}^l \mathbf{v}^{l-1} + \mathbf{b}^l), 0 < l < L \quad 15.1$$

此处， $\mathbf{z}^l = \mathbf{W}^l \mathbf{v}^{l-1} + \mathbf{b}^l \in \mathbb{R}^{N_l \times 1}$ ， $\mathbf{v}^{l-1} \in \mathbb{R}^{N_{l-1} \times 1}$ ， $\mathbf{W}^l \in \mathbb{R}^{N_l \times N_{l-1}}$ ， $\mathbf{b}^l \in \mathbb{R}^{N_l \times 1}$ ， $N_l \in \mathbb{R}$ ；它们

分别是激励向量， activation 向量， 权重矩阵， 偏移向量以及  $l$  层神经元数。 $\mathbf{V}^0 = \mathbf{0} \in \Re^{N_0 \times 1}$  是观测（特征）向量。 $N_0 = D$  是特征维度。 $f(\bullet): \Re^{N_l \times 1} \rightarrow \Re^{N_l \times 1}$  是激励函数。在大多数应用中，激励函数是 s 型函数是：

$$\sigma(z) = \frac{1}{1 + e^{-z}}, \quad 15.2$$

也可能是双曲正切函数：

$$\tanh(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}} \quad 15.3$$

对于回归型问题，线性层常采用下式模型：

$$\mathbf{v}^L = \mathbf{z}^L = \mathbf{W}^L \mathbf{v}^{L-1} + \mathbf{b}^L \quad 15.4$$

生成输出向量  $\mathbf{v}^L \in \Re^{N_L}$ ， $N_L$  是 输出维数。

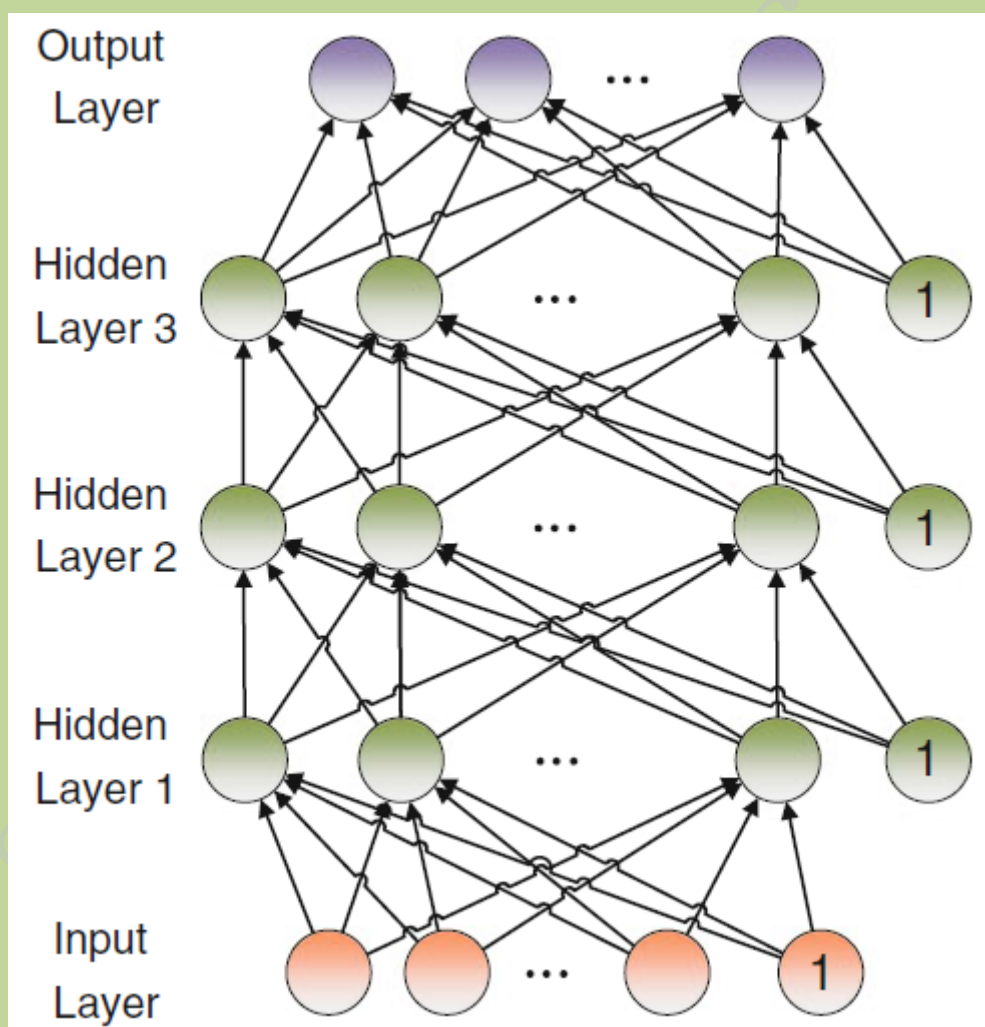


图 15.1 深度神经元网络

对于分类问题，神经元的每一个输出代表一个类  $i \in \{1, \dots, C\}$ ， $C = N_L$  是类的总数。第  $i$  个输出  $v_i^L$  代表了观测向量  $\mathbf{o}$  的观测概率  $P_{dm}(i|\mathbf{o})$ 。

DNN 前向计算过程:

```

 $V^0 \leftarrow o$ 
for  $l \leftarrow 1; l < L; l \leftarrow l+1$  do
     $Z^l \leftarrow W^l V^{l-1} + B^l$ 
     $V^l \leftarrow f(Z^l)$ 
end for
 $Z^L \leftarrow W^L V^{L-1} + B^L$ 
if regression then
     $V^l \leftarrow Z^l$ 
else
     $V^l \leftarrow \text{soft max}(Z^l)$ 
return  $V^l$ 
end procedure

```

## 15.2 误差反向传播参数估计

DNN 的模型参数  $\{W, b\}$  是未知的, 可以从训练样本  $S = \{(o^m, y^m) | 0 \leq m < M\}$  中估计,  $M$  是训练的样本数,  $o^m$  是  $m$  阶观测向量,  $y^m$  是期望的输出向量。这一过程通常称为训练过程或者参数估计过程。

### 15.2.1 训练评价准则

最小期望损准则:

$$J_{EL} = E(J(W, b; o, y)) = \int_o J(W, b; o, y) p(o) d(o) \quad 15.7$$

$\{W, b\}$  是参数模型,  $o$  是观测向量,  $y$  是输出向量,  $p(o)$  是观测向量  $o$  概率密度函数,  $J(W, b; o, y)$  是损耗函数。

回归问题常采用均方误差准则:

$$J_{MSE}(W, b; S) = \frac{1}{M} \sum_{m=1}^M J_{MSE}(W, b; o^m, y^m) \quad 15.8$$

$$J_{MSE}(W, b; o, y) = \frac{1}{2} \|v^L - y\|^2 = \frac{1}{2} (v^L - y)^T (v^L - y) \quad 15.9$$

对于分类问题,  $y$  是概率分布, 则交叉熵准则是:

$$J_{CE}(W, b; S) = \frac{1}{M} \sum_{m=1}^M J_{CE}(W, b; o^m, y^m) \quad 15.10$$

$$J_{CE}(W, b; o, y) = \sum_{i=1}^C y_i \log v_i^L \quad 15.11$$

$y_i = P_{emp}(i|o)$  是观测向量  $o$  所属分类  $i$  的经验概率,  $v_i^L = P_{dnn}(i|o)$  是和上式一样的 DNN 的概率。

### 15.2.2 训练算法

误差反向传递算法:

$$\mathbf{W}_{t+1}^l \leftarrow \mathbf{W}_t^l - \varepsilon \Delta \mathbf{W}_t^l \quad 15.14$$

$$\mathbf{b}_{t+1}^l \leftarrow \mathbf{b}_t^l - \varepsilon \Delta \mathbf{b}_t^l \quad 15.15$$

$\mathbf{W}_t^l$  和  $\mathbf{b}_t^l$  分别是  $l$  层的第  $t$  次更新的权向量和偏移向量。

$$\Delta \mathbf{W}_t^l = \frac{1}{M_b} \sum_{m=1}^{M_b} \nabla_{\mathbf{W}_t^l} J(\mathbf{W}, \mathbf{b}; \mathbf{o}^m, \mathbf{y}^m) \quad 15.16$$

$$\Delta \mathbf{b}_t^l = \frac{1}{M_b} \sum_{m=1}^{M_b} \nabla_{\mathbf{b}_t^l} J(\mathbf{W}, \mathbf{b}; \mathbf{o}^m, \mathbf{y}^m) \quad 15.17$$

分别是  $t$  时刻的平均权重梯度和平均偏移向量。

最上层的权重矩阵和偏移向量的梯度取决于选择的准则, 对于回归问题, MSE 训练准则 (15.9) 且其线性输出层是式 15.5, 则输出层的权重矩阵的梯度是:

$$\nabla_{\mathbf{W}_t^L} J_{MSE}(\mathbf{W}, \mathbf{b}; \mathbf{o}, \mathbf{y}) = \nabla_{z_t^L} J_{MSE}(\mathbf{W}, \mathbf{b}; \mathbf{o}, \mathbf{y}) \frac{\partial z_t^L}{\partial \mathbf{W}_t^L} = e_t^L (v_t^{L-1})^T$$

## 15.3 实际处理中的问题

### 15.3.1 数据预处理

在语音识别中, CMN(Cepstral Normalization)倒谱系数归一化, 即减去每个音节 MFCC 特征的均值能减小通道引入的畸变。CMN 的归一化可以通过首先估算样本的均值, 然后用样本值减去估算均值而得到。

### 15.3.2 模型初始化

### 15.3.3 权重衰减

控制过拟合的方法, 过拟合是最最大限度减少预期损失带来的, 因为实际上是减少的训练集的损失。控制过拟合的方法是规范训练准则使参数模型并不是十分适合样本。常用的规范标准包括: 基于  $L_1$  规范

$$R_1(\mathbf{W}) = \|\text{vec}(\mathbf{W})\|_1 = \sum_{l=1}^L \|\text{vec}(\mathbf{W}^l)\|_1 = \sum_{l=1}^L \sum_{i=1}^{N_l} \sum_{j=1}^{N_{l-1}} |\mathbf{W}_{ij}^l| \quad 15.37$$

基于  $L_2$  规范,

$$R_2(\mathbf{W}) = \|\text{vec}(\mathbf{W})\|_2^2 = \sum_{l=1}^L \|\text{vec}(\mathbf{W}^l)\|_2^2 = \sum_{l=1}^L \sum_{i=1}^{N_l} \sum_{j=1}^{N_{l-1}} (\mathbf{W}_{ij}^l)^2 \quad 15.38$$



$\mathbf{W}_{ij}$  是矩阵  $\mathbf{W}$  的  $(i, j)$  的第  $n$  个值,  $\text{vec}(\mathbf{W}) \in \mathbb{R}^{[N_t \times N_{t-1}] \times 1}$  是矩阵  $\mathbf{W}^l$  所有列的组合。

$\|\text{vec}(\mathbf{W}^l)\|_2$  等于  $\|\mathbf{W}^l\|_F$ , 即  $\mathbf{W}^l$  的弗罗贝尼乌斯范数。

当引入规范准则, 训练模型准则变为:

$$\ddot{J}(\mathbf{W}, \mathbf{b}; S) = J(\mathbf{W}, \mathbf{b}; S) + \lambda R(\mathbf{W}) \quad 4.39$$

$J(\mathbf{W}, \mathbf{b}; S)$  是  $J_{MSE}(\mathbf{W}, \mathbf{b}; S)$  或者  $J_{CE}(\mathbf{W}, \mathbf{b}; S)$  在训练集  $S$  上对应的经验损失最优。 $\lambda$  是插入权重。

$$\nabla_{\mathbf{W}_t^l} \ddot{J}(\mathbf{W}, \mathbf{b}; \mathbf{o}, \mathbf{y}) = \nabla_{\mathbf{W}_t^l} J(\mathbf{W}, \mathbf{b}; \mathbf{o}, \mathbf{y}) + \lambda \nabla_{\mathbf{W}_t^l} R(\mathbf{W}) \quad 4.40$$

$$\nabla_{\mathbf{b}_t^l} \ddot{J}(\mathbf{W}, \mathbf{b}; \mathbf{o}, \mathbf{y}) = \nabla_{\mathbf{b}_t^l} J(\mathbf{W}, \mathbf{b}; \mathbf{o}, \mathbf{y}) \quad 4.40$$

此处:

$$\nabla_{\mathbf{W}_t^l} R_1(\mathbf{W}) = \text{sgn}(\mathbf{W}_t^l), \quad \nabla_{\mathbf{W}_t^l} R_2(\mathbf{W}) = 2\mathbf{W}_t^l$$

权重衰减在样本数较少时很有用。在语音识别类型的问题中权重矩阵通常有一百万个参数, 当样本集很大时, 插入权重  $\lambda$  应当很小甚至是零。

#### 15.3.4 丢弃

过拟合问题也会采用丢弃的方式解决问题, 丢弃的基本思想是对每一个隐藏层随机忽略一部分神经元。这就要求检测模式的神经元之间的依赖性较弱。

### 15.4 高级模型初始化方法

#### 15.4.1 受限玻尔兹曼机

RBM (Restricted Boltzmann machine) 是一个随机神经生成网络。本质上是一个无向图, 该图包括一层随机可见的神经和一层随机不可见神经。对于二值 RBM,  $\mathbf{v} \in \{0, 1\}^{N_v \times 1}$ ,  $\mathbf{h} \in \{0, 1\}^{N_h \times 1}$ , 对每一个可见向量  $\mathbf{v}$  和不可见向量  $\mathbf{h}$  RBM 注入能量如下:

$$E(\mathbf{v}, \mathbf{h}) = -\mathbf{a}^T \mathbf{v} - \mathbf{b}^T \mathbf{h} - \mathbf{h}^T \mathbf{W} \mathbf{v} \quad 15.4.1$$

$N_v$  和  $N_h$  分别是可见和隐藏神经元数。 $\mathbf{W} \in \mathbb{R}^{N_h \times N_v}$  是权重矩阵,  $\mathbf{a}$  和  $\mathbf{b}$  是偏移矩阵。如果可见矩阵是实数, 则 RBM 变成高斯-伯努利分布 (Gaussian-Bernoulli), 注入能量变成:

$$E(\mathbf{v}, \mathbf{h}) = \frac{1}{2} (\mathbf{v} - \mathbf{a})^T (\mathbf{v} - \mathbf{a}) - \mathbf{b}^T \mathbf{h} - \mathbf{h}^T \mathbf{W} \mathbf{v} \quad 15.4.2$$

根据能量, 每一种情况分配了一个概率:

$$P(\mathbf{v}, \mathbf{h}) = \frac{e^{-E(\mathbf{v}, \mathbf{h})}}{Z} \quad 15.4.3$$

其中  $Z = \sum_{\mathbf{v}, \mathbf{h}} e^{-E(\mathbf{v}, \mathbf{h})}$  是已知分区函数的归一化因子。

在 RBM 算法中， $P(\mathbf{v} | \mathbf{h})$  和  $P(\mathbf{h} | \mathbf{v})$  的后验概率是很容易计算的。例如在伯努利-伯努利模型中：

$$\begin{aligned} P(\mathbf{v} | \mathbf{h}) &= \frac{e^{-E(\mathbf{v}, \mathbf{h})}}{\sum_{\mathbf{v}, \mathbf{h}} e^{-E(\mathbf{v}, \mathbf{h})}} = \frac{e^{a^T \mathbf{v} + b^T \mathbf{h} + \mathbf{h}^T \mathbf{W} \mathbf{v}}}{\sum_{\tilde{\mathbf{h}}} e^{a^T \mathbf{v} + b^T \tilde{\mathbf{h}} + \tilde{\mathbf{h}}^T \mathbf{W} \mathbf{v}}} = \frac{\prod_i e^{b_i h_i + h_i \mathbf{W}_{i,*} \mathbf{v}}}{\sum_{\tilde{\mathbf{h}}} \dots \sum_{\tilde{h}_N} \prod_i e^{b_i \tilde{h}_i + \tilde{h}_i \mathbf{W}_{i,*} \mathbf{v}}} \\ &= \prod_i \frac{e^{b_i h_i + h_i \mathbf{W}_{i,*} \mathbf{v}}}{\sum_{\tilde{h}_N} e^{b_i \tilde{h}_i + \tilde{h}_i \mathbf{W}_{i,*} \mathbf{v}}} = \prod_i P(h_i | \mathbf{v}) \end{aligned} \quad 15.4.4$$

$\mathbf{w}_{i,*}$  是矩阵  $\mathbf{w}$  的第  $i$  行。15.4.4 表明相对于可见向量隐藏神经元是有条件独立的。

$$P(h_i | \mathbf{v}) = \frac{e^{b_i 1 + \mathbf{W}_{i,*} \mathbf{v}}}{e^{b_i 1 + \mathbf{W}_{i,*} \mathbf{v}} + e^{b_i 0 + 0 \mathbf{W}_{i,*} \mathbf{v}}} = \sigma(b_i + \mathbf{W}_{i,*} \mathbf{v}) \quad 15.4.5$$

#### 15.4.2 RBM 性能

定义如下自由能量：

$$F(\mathbf{v}) = -\log \left( \sum_{\mathbf{h}} e^{-E(\mathbf{v}, \mathbf{h})} \right) \quad 15.4.6$$

使用自由能量，可以将边际概率  $P(\mathbf{v})$  写成如下：

$$P(\mathbf{v}) = \sum_{\mathbf{h}} P(\mathbf{v}, \mathbf{h}) = \sum_{\mathbf{h}} \frac{e^{-E(\mathbf{v}, \mathbf{h})}}{Z} = \frac{\sum_{\mathbf{h}} e^{-E(\mathbf{v}, \mathbf{h})}}{Z} = \frac{e^{-F(\mathbf{v})}}{\sum_{\mathbf{v}} e^{-F(\mathbf{v})}} \quad 15.4.7$$

#### RBM 参数学习

训练 RBM 时，采用了随机梯度下降法（SGD）来最小化负对数似然（NLL，negative log likelihood）的方法。

$$J_{NLL}(\mathbf{W}, \mathbf{a}, \mathbf{b}; \mathbf{v}) = -\log P(\mathbf{v}) = F(\mathbf{v}) + \log \sum_{\mathbf{v}} e^{-F(\mathbf{v})} \quad 15.4.8$$

参数更新方法如下：

$$\mathbf{W}_{t+1} \leftarrow \mathbf{W}_t - \varepsilon \Delta \mathbf{W}_t$$

$$\mathbf{a}_{t+1} \leftarrow \mathbf{a}_t - \varepsilon \Delta \mathbf{a}_t$$

$$\mathbf{b}_{t+1} \leftarrow \mathbf{b}_t - \varepsilon \Delta \mathbf{b}_t$$

这里  $\varepsilon$  是学习速率：

$$\Delta \mathbf{W}_t = \rho \Delta \mathbf{W}_{t-1} + (1 - \rho) \frac{1}{M_b} \sum_{m=1}^{M_b} \nabla_{\mathbf{W}_t} J_{NLL}(\mathbf{W}, \mathbf{a}, \mathbf{b}; \mathbf{v}^m)$$

$$\Delta \mathbf{a}_t = \rho \Delta \mathbf{a}_{t-1} + (1 - \rho) \frac{1}{M_b} \sum_{m=1}^{M_b} \nabla_{\mathbf{a}_t} J_{NLL}(\mathbf{W}, \mathbf{a}, \mathbf{b}; \mathbf{v}^m)$$



$$\Delta \mathbf{b}_t = \rho \Delta \mathbf{b}_{t-1} + (1-\rho) \frac{1}{M_b} \sum_{m=1}^{M_b} \nabla_{\mathbf{b}_t} J_{NLL}(\mathbf{W}, \mathbf{a}, \mathbf{b}; \mathbf{v}^m)$$

$\rho$  是动量参数， $M_b$  是分批处理量的大小， $\nabla_{\mathbf{W}_t} J_{NLL}(\mathbf{W}, \mathbf{a}, \mathbf{b}; \mathbf{v}^m)$ ， $\nabla_{\mathbf{a}_t} J_{NLL}(\mathbf{W}, \mathbf{a}, \mathbf{b}; \mathbf{v}^m)$ ， $\nabla_{\mathbf{b}_t} J_{NLL}(\mathbf{W}, \mathbf{a}, \mathbf{b}; \mathbf{v}^m)$  分别是 NLL 准则模型参数  $\mathbf{W}$ ， $\mathbf{a}$  和  $\mathbf{b}$  的梯度。

NLL 的倒数如下：

$$\nabla_{\theta} J_{NLL}(\mathbf{W}, \mathbf{a}, \mathbf{b}; \mathbf{v}) = - \left[ \left\langle \frac{\partial E(\mathbf{v}, \mathbf{h})}{\partial \theta} \right\rangle_{data} - \left\langle \frac{\partial E(\mathbf{v}, \mathbf{h})}{\partial \theta} \right\rangle_{model} \right] \quad 15.4.9$$

$\theta$  是一些模型参数， $\langle x \rangle_{data}$  和  $\langle x \rangle_{model}$  分别是数据模型估计的  $x$  期望。对于可见和隐藏权重的更新有下式：

$$\nabla_{w_{ij}} J_{NLL}(\mathbf{W}, \mathbf{a}, \mathbf{b}; \mathbf{v}) = - \left[ \langle v_i h_j \rangle_{data} - \langle v_i h_j \rangle_{model} \right] \quad 15.4.10$$

$\langle v_i h_j \rangle_{data}$  是样本集中可见神经元和隐藏神经元激发的频率， $\langle v_i h_j \rangle_{model}$  是对应的模型激发频率。当隐藏值未知时，需要计算指数次才能得到  $\langle \bullet \rangle_{model}$ 。所以常常会采用接近的算法。

常采用对比分歧（CD，contrastive divergence）快速算法来接近，可见隐藏权重梯度的一步对比分歧接近估计如下：

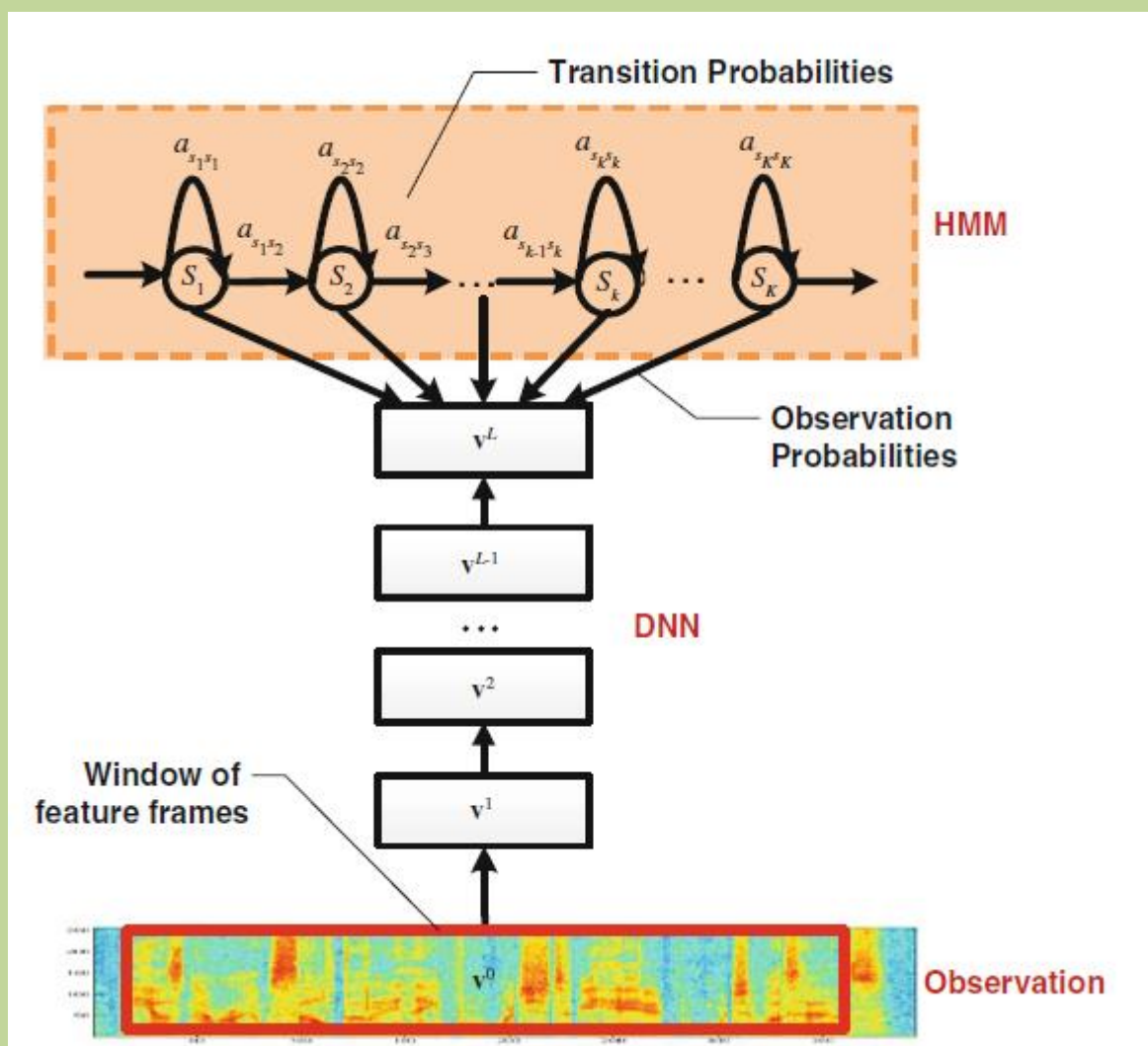
$$\nabla_{w_{ji}} J_{NLL}(\mathbf{W}, \mathbf{a}, \mathbf{b}; \mathbf{v}) = - \left[ \langle v_i h_j \rangle_{data} - \langle v_i h_j \rangle_{\infty} \right] \approx - \left[ \langle v_i h_j \rangle_{data} - \langle v_i h_j \rangle_1 \right] \quad 15.4.11$$

$\langle \bullet \rangle_{\infty}$  和  $\langle \bullet \rangle_1$  分别是运行无限步和一步吉布斯采样器的期望，

15.4.4

## 15.5 深度神经元和马尔科夫混合模型（DNN-HMM）

在 DNN-HMM 算法中，语音的力度变化由 HMM 模型来匹配，观测概率使用 DNN 方法来计算。这一过程如下图。



DNN-HMM 混合系统架构

DNN 的每个神经元输出

Deep Learning 理论

Deep learning 概念源于人造神经网络研究。分为三类：

- 1 无监督或者生成学习，用于模式的分析和综合，找到观测数据的高阶相关性。
- 2 有监督学习，用于模式分类。
- 3 混合学习，

有监督的识别方法是后项传递神经网络（BPNN）。

NLP (native language processing)

TTS (text to speech)

AEC 算法:

AEC 算法主要包括如下几个重要的模块:

1 回声延迟估计

2 NLMS (normalize least mean square) 最小均方误差

3 NLP ( ) 非线性滤波

4 CNG 舒适噪声产生, 一般经典 aec 算法还包括双端检测 (DT)

回声延迟估计

回声延迟长短对回声抵消器的性能有较大影响, 过长的滤波器抽头会带来较大的延迟, 并且语音信号是短时平稳信号, 过长的滤波器抽头也不适合短时平稳信号特征。

基于相关的延迟算法。该算法的主要思想是:

设 1 表示有说话音, 0 表示无说话声 (或者很弱的说话声), 参考端 (远端) 信号  $x(t)$  和接收端信号可能的组合方式如下:

(0, 0); (0, 1); (1, 0); (1, 1)

webrtc 默认 (1, 0) 和 (0, 1) 是不可能发生的。设在时间间隔  $p$  上, 即  $p = 1, 2, \dots, P$ , 频带  $q, q = 1, 2, \dots, Q$ , 输入信号  $x$  加窗后的功率谱用  $X_w(p, q)$  表示, 其角标表示其加了窗函数。对每个频带的功率谱设定一个门限  $X_w(p, q)$ ,

如果  $X_w(p, q) \geq X_w(p, q)_{threshold}$ , 则  $X_w(p, q) = 1$ ;

如果  $X_w(p, q) < X_w(p, q)_{threshold}$ , 则  $X_w(p, q) = 0$ ;

同理, 对于信号  $y(t)$ , 加窗信号功率谱  $Y_w(p, q)$  和门限  $Y_w(p, q)_{threshold}$ ,

如果  $Y_w(p, q) \geq Y_w(p, q)_{threshold}$ , 则  $Y_w(p, q) = 1$ ;

如果  $Y_w(p, q) < Y_w(p, q)_{threshold}$ , 则  $Y_w(p, q) = 0$ ; 考虑到实际处理的方便, 在 webrtc 的 c 代码中, 将经过 fft 变换后的频域功率谱分为 32 个子带, 这样每个特定子带  $X_w(p, q)$  的值可以用 1 个比特来表示, 总共需要 32 个比特, 只用一个 32 位数据类型就可以表示。

2) NLMS 归一化最小均方自适应算法

LMS/NLMS/RLS 都是经典自适应滤波算法, 设远端信号为  $x(n)$ , 近端信号为  $d(n)$ , 则误差信号  $e(n) = d(n) - w(n)x(n)$ , 此处'表示转置, 因为信号一般使用列向量表示的。NLMS 对滤波器的系数更新使用变步长方法, 即步长  $u = u_0 / (\gamma + x'(n) * x(n))$ ; 其中  $u_0$  为更新

步长因子,  $\gamma$  是稳定因子, 则滤波器系数更新方程为  $W(n+1)=W(n)+\gamma e(n)x(n)$ ;

### 3)NLP(非线性滤波)

webrtc 采用了维纳滤波器, 此处只给出传递函数的表达式, 设估计的语音信号的功率谱为  $P_s(\omega)$ , 噪声的功率谱为  $P_n(\omega)$ , 则滤波器的传递函数为  $H(\omega)=P_s(\omega)/(P_s(\omega)+P_n(\omega))$ 。

### 4) CNG (舒适噪声产生)

首先生成在  $[0, 1]$  上均匀分布随机噪声矩阵, 再用噪声的功率谱开方后去调制噪声的幅度。

fullaec 算法需要注意两点:

1) 延迟要小, 因为算法默认滤波器长度是分为 12 块, 每块 64 点, 按照 8000 采样率, 也就是  $12 \times 8\text{ms} = 96\text{ms}$  的数据, 超过这个长度就处理不了了。

2) 延迟抖动要小, 因为算法是默认 10 块也计算一次参考数据的位置 (即滤波器能量最大的那一块), 所以如果抖动很大的话, 找参考数据不准确的, 这样回声就消除不掉了。

PBFDFAF (基于分段块频域自适应滤波算法的长延时回声消除, partitioned Multiple Block Frequency Domain Adaptive Filter)