

Mel Log Spectrum Approximation (MLSA) Filter for Speech Synthesis

Satoshi Imai, Kazuo Sumita* and Chieko Furuichi, Regular Members

Research Laboratory of Precision Machinery and Electronics,
Tokyo Institute of Technology, Yokohama, Japan 227

SUMMARY

The spectral envelope of speech can be represented efficiently by the log magnitude spectrum on the nonlinear frequency scale, which is close to mel scale (called mel-log spectrum). The mel cepstrum defined by its Fourier coefficients is also considered to have a suitable property as the parameter to represent the spectral envelope. So far, no satisfactory filter has been reported for the synthesis approximating the mel-log spectrum. This paper presents a method of constructing the mel-log spectrum approximation (MLSA) filter, which has a relatively simple structure and a low coefficient sensitivity, together with a design example of MLSA filter for speech synthesis. The transfer function of MLSA filter is represented by Padé approximation, which approximates the exponential of the transfer function of the filter (so-called basic filter). Since the transfer function of the basic filter is represented by a polynomial with the transfer function of the first-order all-pass filter as the variable, it is necessary in the realization of the filter to delete from the feedback loop the path without a delay. By the construction method of MLSA filter shown in this paper, the path without delay can easily be deleted from the feedback loop in the MLSA filter. The obtained MLSA filter is of relatively simple structure and has low coefficient sensitivity. The quantization characteristics of the coefficient are also satisfactory.

1. Introduction

Since the spectral envelope of the speech is generally of pole-zero type and

the sensation to the sound level is logarithmic, the spectral envelope can well be represented by the log spectrum. The cepstrum is defined by the Fourier coefficient of the log spectrum, and is matched to the characteristics of the hearing sensation while retaining the correspondence to the original spectrum. The sensitivity of the log spectral envelope to the cepstrum is low and the cepstrum has a good quantization property.

In the log spectral envelope, the cepstrum has low distortion in the interpolation along the time axis [1]. The extraction of the cepstrum is relatively easy, and can be used as the parameter in the pole-zero model. In the speech synthesis, the log magnitude approximation filter (LMA filter) can be used, which applies the cepstrum directly to the filter coefficients. However, the order of the cepstrum required in speech synthesis with satisfactory quality, generally, is slightly larger than that of the parameter in the linear prediction method.

By considering the hearing characteristics in regard to the sound level, it seems more desirable to represent the speech spectral envelope by the log magnitude spectrum on a nonlinear frequency scale such as mel scale (called mel-log spectrum). This will lead to a more efficient representation than the usual log spectrum, and by using the mel cepstrum defined by its Fourier coefficients, the spectral envelope can be represented by the order which is lower than that in usual cepstrum representation.

Until now, however, no satisfactory filter has been reported which can approximate the mel-log spectrum for speech synthesis. It is considered that the log magnitude spectrum can better be approximated by recursive filter than by nonrecursive filter, and it is a natural consequence that the mel-log spectrum is approximated by a recursive

*Presently with Central Res. Lab., Toshiba Co.

filter. As the recursive filter realizing a frequency response on the nonlinear frequency scale such as mel scale, a construction method has been proposed of a nonlinear frequency linear prediction synthesis filter by Strube [3]. His method, however, is difficult to apply to the realization of mel-log spectrum approximation (MLSA) filter.

The transfer function of MLSA filter can be represented by a rational function with the transfer function of a first-order, all-pass filter as the variable. Consequently, it may contain a path without a delay in the feedback loop of the filter, which must be deleted. In Strube's method, to delete the path without a delay from the feedback loop, the denominator polynomial of the transfer function, with the transfer function of the first-order all-pass filter as the variable, is a polynomial with the transfer function of the first-order low-pass filter with the same parameter as the all-pass filter as the variable. In this method, the sensitivity of the log magnitude response to the filter coefficient becomes very large, and at the same time, the coefficient values exhibit a wide distribution. For those reasons, it is almost impossible to apply the method to MLSA filter.

This paper presents a construction method for MLSA filter, resulting in a relatively simple structure and low coefficient sensitivity. A design example is shown for the MLSA filter for speech synthesis. In the proposed construction method, in contrast to Strube's method, the path without a delay is deleted from the feedback loop by essentially modifying the denominator polynomial of the transfer function, which uses the transfer function of the first-order, all-pass filter as the variable, into a polynomial in which the terms not containing the unit delay variable as a factor and the terms containing the variable as a first-order factor are separated.

2. Mel-Log Spectrum and Mel Cepstrum

2.1 Approximate representation of mel-scale

One of the frequency scales that can well approximate the mel scale, leading to easy modelling of the spectrum, is as follows. Consider an all-pass filter for which the transfer function $H^{(a)}(z)$ and the frequency response $H^{(a)}(e^{j\omega})$ ($\omega = \omega \Delta t$, ω : radian frequency, Δt : unit delay in the digital filter) are given by

$$H^{(a)}(z) = (z^{-1} - \alpha) / (1 - \alpha z^{-1}) \quad (1)$$

$$H^{(a)}(e^{j\omega}) = \exp(-j\beta_a(\omega)) \quad (2)$$

The phase response is given by

$$\beta_a(\omega) = \tan^{-1} \frac{(1 - \alpha^2) \sin \omega}{(1 + \alpha^2) \cos \omega - 2\alpha} \quad (3)$$

Then the nonlinear frequency scale is defined as

$$\tilde{\omega} = \beta_a(\omega) \quad (4)$$

In the phase response $\beta_a(\omega)$ of this all-pass filter, the unit delay of the digital filter Δt is set as 100 μs , 125 μs ($= (8 \text{ kHz})^{-1}$), or 156.25 μs ($= (6.4 \text{ kHz})^{-1}$), and α in the transfer function $H^{(a)}(z)$ of Eq. (1) is set as 0.35, 0.31, or 0.28. Then the error to the mel scale is less than 2.5% of the full-scale value which is satisfactory. When $\alpha = 0$, it is the same as the linear scale. The parameter α of the all-pass filter is called the frequency compression parameter. The scale defined by Eq. (4) includes the linear scale as a special case and can well represent the nonlinear frequency scale such as mel scale. It is called simply mel frequency scale in this paper.

2.2 Mel-log spectrum

Letting the usual log spectrum on the linear frequency scale be $G_o(\omega)$, the log magnitude spectrum $G_a(\tilde{\omega})$ on the mel-frequency scale $\tilde{\omega}$ is given by

$$G_a(\tilde{\omega}) = G_o(\beta_a^{-1}(\tilde{\omega})) \quad (5)$$

$$\beta_a^{-1}(\tilde{\omega}) = \tan^{-1} \frac{(1 - \alpha^2) \sin \tilde{\omega}}{(1 + \alpha^2) \cos \tilde{\omega} + 2\alpha} \quad (6)$$

where $\beta_a^{-1}(\tilde{\omega})$ is the inverse function of the phase response $\beta_a(\omega)$ of the all-pass filter given by Eq. (3).

The log magnitude spectrum $G_a(\tilde{\omega})$ on the mel-frequency scale $\tilde{\omega}$ is simply called mel-log spectrum. The mel-log spectrum representing the spectral envelope of the speech waveform is called mel-log spectral envelope.

2.3 Mel cepstrum

Mel cepstrum $c_a(m)$ is defined by the Fourier coefficients of the mel-log spectrum $G_a(\tilde{\omega})$ as

$$g_a(m) = \frac{1}{2\pi} \int_{-\pi}^{\pi} G_a(\tilde{\omega}) e^{jm\tilde{\omega}} d\tilde{\omega} = g_a(-m) \quad (7)$$

as

$$c_a(m) = \begin{cases} 2g_a(m) & (m > 0) \\ g_a(m) & (m = 0) \\ 0 & (m < 0) \end{cases} \quad (8)$$

Consequently, the mel-log spectrum $G_a(\tilde{\omega})$ can be represented by the mel spectrum $c_a(m)$ as

$$G_a(\tilde{\omega}) = \sum_{m=0}^{\infty} c_a(m) \cos(m\tilde{\omega}) \quad (9)$$

3. Direct Approximation of Mel-Log Spectrum

3.1 Frequency response on mel-frequency scale

Let the transfer function $H^{(a)}(z)$ of Eq. (1) be

$$\tilde{z}^{-1} = H^{(a)}(z) \quad (10)$$

Then if there exists the relation

$$H_a(\tilde{z}) = H(z) \quad (11)$$

for any transfer function $H(z)$, it follows from Eqs. (1), (2) and (3) that

$$H_a(e^{j\tilde{\omega}}) = H(e^{j\omega}) \quad (12)$$

In other words, $H_a(e^{j\tilde{\omega}})$ represents the frequency response on the mel-frequency scale $\tilde{\omega}$, and $H_a(\tilde{z})$ represents the transfer function of the system with frequency response $H_a(e^{j\tilde{\omega}})$ on $\tilde{\omega}$.

3.2 Filter with exponential transfer function

Assume that the transfer function $H_a(\tilde{z})$ of the digital filter can be represented, using the transfer function $F_a(\tilde{z})$ of the digital filter called basic filter, as

$$H_a(\tilde{z}) = \exp(F_a(\tilde{z})) \quad (13)$$

If the basic filter is stable, its transfer function $F_a(\tilde{z})$ does not have a pole on or outside of the unit circle on the \tilde{z} -plane. Consequently, the transfer function $H_a(\tilde{z})$ of Eq. (13) takes a finite value in that region. Since the inside, circumference, and outside of the unit circle on the \tilde{z} -plane correspond, respectively, to the inside, circumference and outside of the unit circle of the z -plane, it follows that if the basic filter is stable, the filter of Eq. (13) with exponential transfer function is always stable and is a minimum-phase system.

The basic filter is arbitrary as far as it is stable. In the following, a nonrecursive-type filter is considered in which the transfer function $F_a(\tilde{z})$ can be represented as

$$F_a(\tilde{z}) = \sum_{m=0}^M c_a(m) \tilde{z}^{-m} \quad (14)$$

In such a case, the log magnitude response $\ln |H_a(e^{j\tilde{\omega}})|$ of the filter with exponential transfer function on the mel-frequency scale is given by Eqs. (13) and (14) as

$$\ln |H_a(e^{j\tilde{\omega}})| = \sum_{m=0}^M c_a(m) \cos(m\tilde{\omega}) \quad (15)$$

Letting $c_a(m)$ in Eqs. (14) and (15) be the mel cepstrum corresponding to the desired mel-log spectral envelope, it follows from the property of the Fourier series that the log magnitude response of the filter on the mel-frequency scale is the best approximation in the sense of least mean-square error to the desired mel-log spectral envelope. It is actually impossible directly to realize the digital filter with the transfer function of Eq. (13). However, by approximating the transfer function $\exp(F_a(\tilde{z}))$ of exponential type by a polynomial or a rational function of the transfer function $F_a(\tilde{z})$ of the basic filter, the digital filter can be realized. The system with the transfer function of Eq. (13) or its adequate approximation is called mel-log spectrum approximation (MLSA) filter.

4. Construction of MLSA Filter

4.1 Padé-approximation of exponential function

MLSA filter is realized by approximating the transfer function $\exp(F_a(\tilde{z}))$ of exponential type by a rational function of the transfer function $F_a(\tilde{z})$ of the basic filter. The complex exponential function $\exp(w)$ can be approximated by the modified Padé-approximation of (L, L) -th order:

$$R_L(w) = \frac{P_L(w)}{P_L(-w)} \quad \left(\begin{array}{l} L = 1, 2, \dots; \\ |w| \leq \bar{\tau}_L \end{array} \right) \quad (16)$$

$$P_L(w) = P_L^I(w) = 1 + \sum_{l=1}^L p_{L,l}^I w^l \quad (17)$$

$$p_{L,l}^I = \frac{1 - \lambda_{L,l}}{l!} \binom{L}{l} \bigg/ \binom{2L}{l} \quad (18)$$

The coefficients $\lambda_{L,l}$ ($L = 1, 2, 3, 4; 1 \leq l \leq L$) are selected so that the approximation error on the logarithmic function scale

$$E_L(\bar{\tau}_L) = \max_{|w| \leq \bar{\tau}_L} |\ln R_L(w) - \ln(\exp w)| \quad (19)$$

is minimized [2]

$$\lambda_{1,1} = 0.13 \times 10^{-1} \bar{\tau}_1^4 \quad (\bar{\tau}_1 < W_1 = 2.00) \quad (20)$$

$$\left. \begin{aligned} \lambda_{2,1} &= 0.44 \times 10^{-5} \bar{\tau}_2^8 \\ \lambda_{2,2} &= 0.11 \times 10^{-2} \bar{\tau}_2^4 \end{aligned} \right\} (\bar{\tau}_2 < W_2 = 3.48) \quad (21)$$

$$\left. \begin{aligned} \lambda_{3,1} &= 0.40 \times 10^{-9} \bar{\tau}_3^{12} \\ \lambda_{3,2} &= 0.74 \times 10^{-6} \bar{\tau}_3^8 \\ \lambda_{3,3} &= 0.30 \times 10^{-3} \bar{\tau}_3^4 \end{aligned} \right\} (\bar{\tau}_3 < W_3 = 4.75) \quad (22)$$

$$\left. \begin{aligned} \lambda_{4,1} &= 0.49 \times 10^{-14} \bar{\tau}_4^{16} \\ \lambda_{4,2} &= 0.42 \times 10^{-10} \bar{\tau}_4^{12} \\ \lambda_{4,3} &= 0.75 \times 10^{-7} \bar{\tau}_4^8 \\ \lambda_{4,4} &= 0.11 \times 10^{-3} \bar{\tau}_4^4 \end{aligned} \right\} (\bar{\tau}_4 < W_4 = 6.25) \quad (23)$$

where W_L ($L=1,2,3,4$) represents the bounds for the variable for which the modified Padé approximation has neither pole nor zero.

In this paper, the modified Padé approximation is called simply Padé approximation. The approximation error $E_L(\bar{\tau}_L)$, on the logarithmic function scale, of the Padé approximation of (L, L) -th order is bounded by [2]

$$E_L(\bar{\tau}_L) \leq \frac{2L!(L+1)!}{(2L)!(2(L+1))!} \bar{\tau}_L^{2L+1} (\bar{\tau}_L < W_L) \quad (24)$$

To suppress the approximation error to ϵ or less, the upper bound $\bar{\tau}_L$ of the variable should be set as

$$\bar{\tau}_L \leq \left(\frac{(2L)!(2(L+1))!\epsilon}{2L!(L+1)!} \right)^{\frac{1}{2L+1}} \quad (25)$$

For example, if the approximation error is to be suppressed to 0.023 (0.2 dB) or less, the upper bound $\bar{\tau}_L$ of the variable should be set so as not to exceed 0.65 ($L=1$), 1.75 ($L=2$), 3.03 ($L=3$) or 4.37 ($L=4$).

The numerator polynomial $P_L(w)$ of the Padé approximation $R_L(w)$ can be written as

$$P_L(w) = P_L^{\text{II}}(w) = 1 + \sum_{i=1}^L p_{L,i}^{\text{II}} w^i \quad (26)$$

$$p_{L,i}^{\text{II}} = \frac{L-i+1}{2L-i+1} \frac{1-\lambda_{L,i}}{1-\lambda_{L,i-1}} (\lambda_{L,0} = 0) \quad (27)$$

or

$$\begin{aligned} P_L(w) &= P_L^{\text{III}}(w) \\ &= 1 + \sum_{i=1}^L p_{L,i}^{\text{III}} w^i 2^{-\frac{(i-1)(i+2)}{2}} \end{aligned} \quad (28)$$

$$p_{L,i}^{\text{III}} = \frac{2^i(L-i+1)}{2L-i+1} \frac{1-\lambda_{L,i}}{1-\lambda_{L,i-1}} (\lambda_{L,0} = -1) \quad (29)$$

These representations are more suitable from the viewpoint of the accuracy in computation. Padé approximation $R_L(w)$, in which the numerator polynomial is represented as $P_L^I(w)$ ($I=1, \text{II}, \text{III}$), is written as $R_L^I(w)$ ($I=1, \text{II}, \text{III}$).

4.2 Transfer function of MLSA filter.

To realize a MLSA filter, the transfer function $\exp(F_a(\tilde{z}))$ of exponential type is approximated by the Padé approximation of Eq. (16). The resulting transfer function $H_a(\tilde{z})$ of the MLSA filter is given by

$$H_a(\tilde{z}) = R_L(F_a(\tilde{z})) \quad (30)$$

If the maximum magnitude of the basic filter

$$\tau_L = \max_{\tilde{\omega}} |F_a(e^{j\tilde{\omega}})| \quad (31)$$

does not exceed the upper bound $\bar{\tau}_L$ for the variable range for which the Padé approximation is effective, the log magnitude response of the MLSA filter $\ln |H_a(e^{j\tilde{\omega}})|$ is given by

$$\begin{aligned} \ln |H_a(e^{j\tilde{\omega}})| &= F_a(e^{j\tilde{\omega}}) + e_L \\ &\simeq F_a(e^{j\tilde{\omega}}) \end{aligned} \quad (32)$$

$$|e_L| \leq E_L(\bar{\tau}_L) \quad (33)$$

where $E_L(\bar{\tau}_L)$ is the approximation error of Padé approximation given by Eq. (24).

When the transfer function $F_a(\tilde{z})$ of the basic filter is specified by the speech mel cepstrum $c_a(m)$ through Eq. (14) and if the amplitude $|F_a(e^{j\tilde{\omega}})|$ of the basic filter does not exceed the upper bound $\bar{\tau}_L$ for the variable range to make the Padé approximation effective, the log magnitude response $\ln |H_a(e^{j\tilde{\omega}})|$ of the MLSA filter is a good approximation to the mel-log spectral envelope of the speech.

4.3 Modification of transfer function of basic filter

The transfer function of $F_a(\tilde{z})$ ($=F(z)$) of the basic filter given by Eq. (14) contains terms which do not have the unit-delay variable z^{-1} as the factor. As itself, MLSA filter of transfer function $R_L(F_a(\tilde{z}))$ of Eq. (30) cannot be realized.

Applying Strube's method (called method I) to the transfer function $F_a(\tilde{z})$ ($=F(z)$) of the basic filter for the MLSA filter, Eq. (14) is modified as [3]

$$F_a(\tilde{z}) = F_a^I(\tilde{z}) = F(z) \\ = b_a^I(0) + \sum_{m=1}^M b_a^I(m) \left(\frac{z^{-1}}{1 - \alpha z^{-1}} \right)^m \quad (34)$$

$$b_a^I(m) = (1 - \alpha^2)^m \sum_{k=m}^M (-\alpha)^{k-m} \binom{k}{m} c_a(k) \quad (35)$$

In this way, the term which does not contain the unit-delay variable z^{-1} as a factor can be separated from the transfer function of the basic filter so that the feedback loop of the MLSA filter with the transfer function of Eq. (30) does not contain a path without a delay. As will be described later, it is difficult to apply the basic filter, obtained by method I to MLSA filter. In this paper, two methods of effective modification are proposed for the transfer function of the basic filter.

By Eqs. (1) and (10), the unit-delay variable z^{-1} can be represented, by the transfer function \tilde{z}^{-1} of the all-pass filter, as

$$z^{-1} = \frac{\alpha + \tilde{z}^{-1}}{1 + \alpha \tilde{z}^{-1}} \quad (36)$$

Consequently, separating the transfer function $F_a(\tilde{z}) (= F(z))$ of the basic filter into terms containing z^{-1} as a factor and the terms not containing z^{-1} , corresponds to the separation of $F_a(\tilde{z})$ into the terms containing $(\alpha + \tilde{z}^{-1})$ as a factor and the terms not containing it.

Method I represents the transfer function $F_a(\tilde{z})$ of the basic filter as a polynomial of $(\alpha + \tilde{z}^{-1})$. As method II, one may modify the transfer function $F_a(\tilde{z})$ of the basic filter as

$$F_a(\tilde{z}) = F_a^{II}(\tilde{z}) \\ = B_a^{II}(0) + (\alpha + \tilde{z}^{-1}) \sum_{m=1}^M B_a^{II}(m) \tilde{z}^{-(m-1)} \\ = \frac{b_a^{II}(0)}{1 - \alpha} + \frac{z^{-1}}{1 - \alpha z^{-1}} \sum_{m=1}^M b_a^{II}(m) \tilde{z}^{-(m-1)} \quad (37)$$

Comparing the coefficients of Eqs. (14) and (37), it is seen that the coefficient $b_a^{II}(m)$ of the filter can be obtained from the mel cepstrum $c_a(m)$ as

$$b_a^{II}(m) = \begin{cases} (1 - \alpha^2) B_a^{II}(m) & (1 \leq m \leq M) \\ (1 - \alpha) B_a^{II}(m) & (m = 0) \end{cases} \quad (38)$$

$$B_a^{II}(M + 1) = 0 \quad (39)$$

$$B_a^{II}(m) = c_a(m) - \alpha B_a^{II}(m + 1) \\ (m = M, M - 1, \dots, 1, 0) \quad (40)$$

As method III, the transfer function $F_a(\tilde{z})$ of the basic filter is modified as

$$F_a(\tilde{z}) = F_a^{III}(\tilde{z}) \\ = b_a^{III}(0) + \frac{\alpha + \tilde{z}^{-1}}{1 + \alpha \tilde{z}^{-1}} \sum_{m=1}^{M+1} b_a^{III}(m) \tilde{z}^{-(m-1)} \\ = b_a^{III}(0) + z^{-1} \sum_{m=1}^{M+1} b_a^{III}(m) \tilde{z}^{-(m-1)} \quad (41)$$

Comparing coefficients of Eqs. (14) and (41), the coefficient $b_a^{III}(m)$ of the filter is obtained from the mel cepstrum $c_a(m)$ as

$$b_a^{III}(M + 1) = \alpha c_a(M) \quad (42)$$

$$b_a^{III}(m) = c_a(m) + \alpha (c_a(m - 1) - b_a^{III}(m + 1)) \\ (m = M, M - 1, \dots, 3, 2) \quad (43)$$

$$b_a^{III}(1) = (c_a(1) - \alpha b_a^{III}(2)) / (1 - \alpha^2) \quad (44)$$

$$b_a^{III}(0) = c_a(0) - \alpha b_a^{III}(1) \quad (45)$$

In the calculation of the coefficient of the basic filter $b_a^K(m)$ ($K = I, II, III$; $m = 0, 1, \dots, M_K$; $M_I = M_{II} = M$, $M_{III} = M + 1$) from the mel cepstrum $c_a(m)$, the amount of computations is $(M + 1)^2/2$ multiplication-addition in method I, $2(M + 1)$ multiplications and $(M + 1)$ additions in method II, and $(M + 3)$ multiplications and $2M$ additions in method III. Also, it is possible to perform the computation in method I by a recursive formula as in methods II and III, but the amount of computation does not change much.

Usually, the order M of the speech mel cepstrum is set between 8 and 15. Consequently, methods II and III require less computation than method I. When the fixed-point arithmetic is used, where addition and subtraction are easier than multiplication, method III is the simplest.

5. Coefficient Sensitivity of MLSA Filter

5.1 Variation of log magnitude due to coefficient error

Consider the case where the transfer function $F_a(\tilde{z})$ of the basic filter is given by

$$F_a(\tilde{z}) = F_a^K(\tilde{z}) \quad (K = I, II, III) \quad (46)$$

and the Padé approximation $P_L(w)$ approximating the exponential function $\exp u$ ($u = F_a(\tilde{z})$) of $F_a(\tilde{z})$ is represented as

$$R_L(w) = R_L^J(w) \quad (J = I, II, III) \quad (47)$$

Then consider the variation in the log magnitude response characteristic due to the error in the coefficients of MLSA filter.

Padé approximation has the coefficients $p_{L,i}^J$ ($J=I, II, III$; $L=1, 2, 3, 4$; $1 \leq i \leq L$), and the transfer function of the basic filter has the coefficients $b_a^K(m)$ ($K=I, II, III$; $m=0, 1, \dots, M_K$; $M_I=M_{II}=M$, $M_{III}=M+1$). Let the errors in $p_{L,i}^J$ and $b_a^K(m)$ be $\Delta p_{L,i}^J$ and $\Delta b_a^K(m)$, respectively. Then the variation ΔS^{JK} of the log magnitude of MLSA filter is given by

$$\begin{aligned} \Delta S^{JK} \approx & \sum_{i=1}^L \left(\frac{\partial}{\partial p_{L,i}^J} \ln |R_L^J(F_a^K(e^{j\tilde{\omega}}))| \right) \Delta p_{L,i}^J \\ & + \sum_{m=0}^{M_K} \left(\frac{\partial}{\partial b_a^K(m)} \ln |R_L^J(F_a^K(e^{j\tilde{\omega}}))| \right) \Delta b_a^K(m) \\ & (J, K=I, II, III, \quad M_I=M_{II}=M, \quad M_{III}=M+1) \end{aligned} \quad (48)$$

The discussion of the variation of Eq. (48) is meaningful only in the range where the exponential type transfer function $\exp(F_a(\tilde{z}))$ can well be approximated by Padé approximation $R_L(F_a(\tilde{z}))$, and the following discussion is made only for the case where the logarithmic error of Padé approximation is 0.023 (0.2 dB) or less, i.e., the variable does not exceed the upper bound \tilde{r}_L given by Eq. (25).

5.2 Sensitivity to Padé approximation coefficients

Padé approximation $R_L(w)$ is characterized by its numerator polynomial $P_L(w)$, and $P_L(w)$ is represented by Eqs. (17), (26) or (28). The maximum sensitivity $S(p_{L,i}^J)$ of the log magnitude of Padé approximation to the coefficient $p_{L,i}^J$ ($J=I, II, III$) of numerator polynomials is given by

$$S(p_{L,i}^J) = \max_{|w| \leq \tilde{r}_L} \left| \frac{\partial}{\partial p_{L,i}^J} \ln |R_L^J(w)| \right| \quad (w = F_a(e^{j\tilde{\omega}})) \quad (49)$$

Calculating numerically the maximum sensitivity $S(p_{L,i}^J)$ of the log magnitude of Padé approximation to the coefficient $p_{L,i}^J$ ($J=I, II, III$), Table 1 is obtained. For MLSA filter, the amount of computation in filtering is almost the same, independently of the type of Padé approximation $R_L^J(w)$ ($J=I, II, III$). Consequently, it is seen from Table 1 that the Padé approximation $R_L^I(w)$ is the best suited since the coefficient sensitivity is the lowest.

Table 1. Maximum Sensitivity $S(p_{L,i}^J)$ ($J=I, II, III$) of Logarithmic Amplitude of Padé Approximation $R_L^J(w)$ to Coefficient $p_{L,i}^J$

$S(p_{L,i}^J)$	$L \backslash i$	1	2	3	4
$S(p_{L,i}^I)$	1	1.5			
	2	5.5	6.7		
	3	18.6	50.1	168.7	
	4	55.7	236.0	1153.9	4380.2
$S(p_{L,i}^{II})$	1	1.5			
	2	4.4	3.3		
	3	11.1	18.1	16.8	
	4	24.7	68.5	92.1	50.6
$S(p_{L,i}^{III})$	1	1.5			
	2	4.4	0.8		
	3	11.1	4.5	2.1	
	4	24.8	17.7	11.8	3.4

5.3 Sensitivity to basic filter coefficients

The maximum sensitivity $S(b_a^K(m))$ of the log magnitude of MLSA filter to the coefficient $b_a^K(m)$ ($K=I, II, III$) of the basic filter is given by

$$\begin{aligned} S(b_a^K(m)) &= \max_{|w_K| \leq \tilde{r}_L} \left| \frac{\partial}{\partial b_a^K(m)} \ln |R_L^J(w_K)| \right| \\ &\leq \max_{|w_K| \leq \tilde{r}_L} \left| \frac{\partial}{\partial b_a^K(m)} \ln |R_L^J(w_K)| \right. \\ &\quad \left. + j \frac{\partial}{\partial b_a^K(m)} \arg R_L^J(w_K) \right| \\ &= \max_{|w_K| \leq \tilde{r}_L} \left| \frac{\partial}{\partial b_a^K(m)} \ln R_L^J(w_K) \right| \\ &\approx \max_{|w_K| \leq \tilde{r}_L} \left| \frac{\partial w_K}{\partial b_a^K(m)} \right| = \bar{S}(b_a^K(m)) \\ &\quad (w_K = F_a^K(e^{j\tilde{\omega}})) \end{aligned} \quad (50)$$

The maximum sensitivity $S(b_a^K(m))$ of the log magnitude of MLSA filter does not exceed the maximum sensitivity $\bar{S}(b_a^K(m))$ of the logarithmic frequency response; $\bar{S}(b_a^K(m))$ ($K=I, II, III$) are given by Eqs. (34), (37), or (41), (50) as

$$\bar{S}(b_a^I(m)) = \max_{\tilde{\omega}} \left| \frac{\alpha + e^{-j\tilde{\omega}}}{1 - \alpha^2} \right|^m = \frac{1}{(1 - \alpha)^m} \quad (51)$$

$$\bar{S}(b_a^{II}(m)) = \max_{\tilde{\omega}} \left| \frac{\alpha + e^{-j\tilde{\omega}}}{1 - \alpha^2} \right| \cdot |e^{-j(m-1)\tilde{\omega}}| = \frac{1}{1 - \alpha} \quad (52)$$

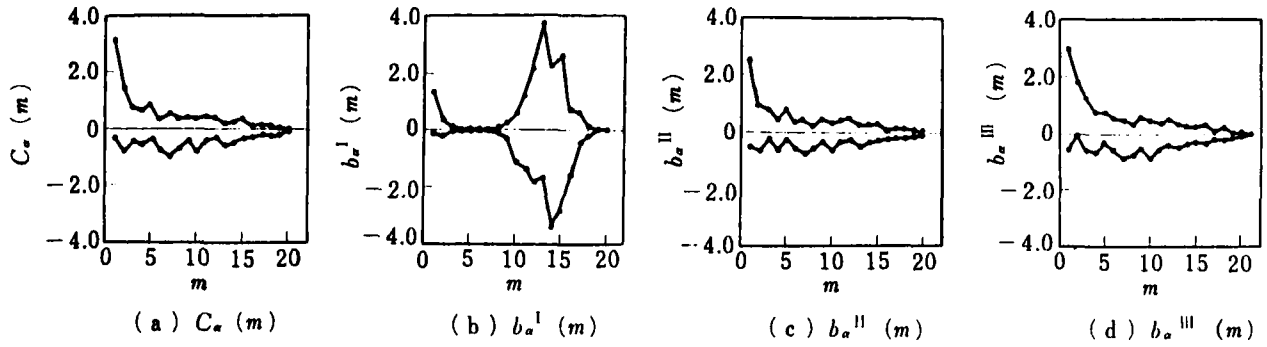


Fig. 1. Maximal and minimal values of mel cepstrum $c_\alpha(m)$ and filter coefficients $b_\alpha^K(m)$ for utterance "nambu dewa higashi no kaze" of a male speaker ($\alpha=0.4$, $K=1, II, III$).

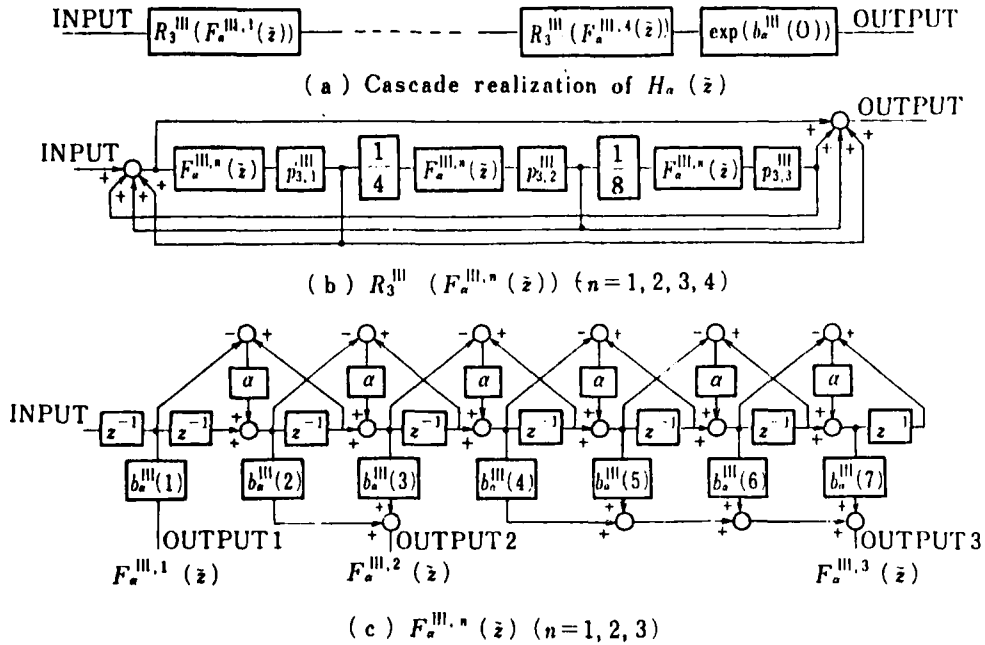


Fig. 2. MLSA filter for speech synthesis: (a) cascade realization; (b) partial filter $R_3^{III}(F_\alpha^{III,n}(z))$; (c) basic filter $F_\alpha^{III,n}(z)$.

$$\bar{S}(b_\alpha^{III}(m)) = \max_{\tilde{\omega}} \left| \frac{\alpha + e^{-j\tilde{\omega}}}{1 + \alpha e^{-j\tilde{\omega}}} \right| \cdot |e^{-j(m-1)\tilde{\omega}}| = 1 \quad (53)$$

Restricting the magnitude response of the basic filter so that the logarithmic error of Padé approximation is 0.023 (0.2 dB) or less, the maximum sensitivity $\bar{S}(b_\alpha^K(m))$ of the logarithmic frequency response of MLSA filter is calculated numerically for the frequency compression parameter of 0.2 and 0.4. Then, it is seen that Eqs. (51) and (52)

coincide within the error of 2%. The values of the coefficients $\bar{S}(b_\alpha^K(m))$ can be determined when the frequency compression parameter α is small and the order m of the coefficient $b_\alpha^K(m)$ is small. Moreover, in that case, they are considerably large.

The maximum sensitivity $\bar{S}(b_\alpha^K(m))$ of the log magnitude response to the basic filter coefficient $b_\alpha^K(m)$ ($K=1, II, III$) of MLSA filter does not exceed the maximum sensitivity $\bar{S}(b_\alpha^K(m))$ of the logarithmic frequency response. Thus, one may conclude that the basic filter

by method III, for which $\delta(b_k^M(m))$ is the smallest, has the best quantization characteristic for the coefficient.

6. Design Example of MLSA Filter

An MLSA filter for speech synthesis is designed. The filter for speech synthesis must approximate naturally well the speech spectral envelope. Other requirements are that the filter coefficients should easily be derived from the parameters representing the spectral envelope, it should have good quantization and time interpolation characteristics of the spectral envelope parameters or filter coefficients, and the filter should have a simple structure with less filtering computation. If an adequate design is made, the MLSA filter satisfies almost all these requirements.

When the speech spectral envelope is represented by the mel-log spectrum, the mel cepstrum can be used as the parameter to represent the spectral envelope, which, however, cannot be used directly as the coefficient of the MLSA filter. The coefficients of the basic filter in MLSA filter, rather than the mel cepstrum, can be used in speech synthesis in a more direct and suitable way.

To apply the filter coefficients as the parameter of the speech spectral envelope, the data compression ratio must be high in its utilization. Consequently, it is required that the sensitivity of the log magnitude of MLSA filter to the filter coefficients is low and the distributions of the filter coefficients are narrow. Figure 1 shows the case of frequency compression parameter α of 0.4, where a male utters "Nanbu dewa higashi no kaze (East wind in southern part)" and the maxima and the minima of the mel cepstrum $c_a(m)$ and the filter coefficient $b_k^M(m)$ ($K=1, \dots, M$) are shown. The sampling frequency is 10 kHz and the speech is digitized with the word length of 12 bits. The frequency analysis is made using 256 point Blackman window and the frame period of 5 ms. The condition for the analysis is the same in the following.

As is shown in Fig. 1, the maximum and the minimum of the filter coefficients $b_k^M(m)$ and $b_k^L(m)$ are almost the same as those of the mel cepstrum $c_a(m)$. Considering the coefficient sensitivity discussed in section 5, it is seen that the most suitable filter coefficient to use is $b_k^M(m)$. The filter coefficient $b_k^L(m)$ has wider distribution and larger coefficient sensitivity (which is not practical).

The mel cepstrum $c_a(m)$ is given as in Eq. (9), by the cosine Fourier coefficients of the mel-log spectral envelope. Conse-

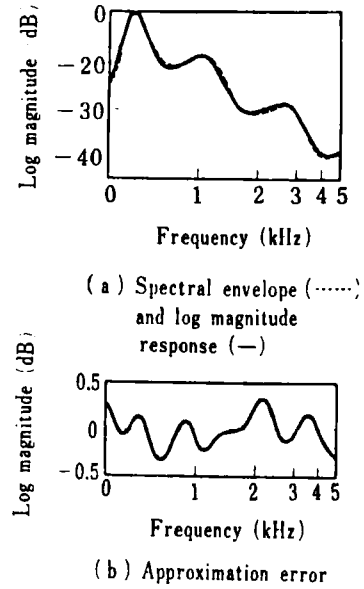


Fig. 3. Mel-log spectral envelope and log magnitude response of 8-bit coefficient MLSA filter, and approximation error for a vowel portion of utterance /na/ of a male speaker ($\alpha = 0.4$, $M = 12$).

quently, the sensitivity of the mel-log spectral envelope to the mel cepstrum $c_a(m)$ is 1, which is the same as the filter coefficient $b_k^M(m)$. Thus, one may consider that the data compression ratio of speech by $b_k^M(m)$ is nearly the same as that of $c_a(m)$.

As Padé approximation $R_L(w)$ to represent the transfer function of MLSA filter, $R_L^M(w)$ is the best suited with the smallest coefficient sensitivity as was shown in section 5. Thus, the MLSA filter for speech synthesis can be represented as a system with the transfer function, in which the exponential-type function $\exp(F_a^M(\tilde{z}))$ is approximated by Padé approximation $P_L^M(w)$.

To apply Padé approximation $P_L^M(w)$ in the range of good approximation, the transfer function $F_a^M(\tilde{z})$ of the basic filter given by Eq. (41) is partitioned into the sum of $b_a^M(0)$ and $F_a^{M,n}(\tilde{z})$ ($n=1, 2, \dots, N$) as

$$F_a^M(\tilde{z}) = b_a^M(0) + \sum_{n=1}^N F_a^{M,n}(\tilde{z}) \quad (54)$$

Consequently, the transfer function $H_a(\tilde{z})$ of MLSA filter is given by

$$H_a(\tilde{z}) = \exp(b_a^M(0)) \prod_{n=1}^N R_L^{M,n}(F_a^{M,n}(\tilde{z})) \quad (55)$$

In the above, care is taken so that the maximum of the frequency response of the term $F_a^{(n)}(z)$ obtained by partitioning the transfer function $F_a^{(n)}(z)$ of the basic filter:

$$r_{L(n)} = \max_{\omega} |F_a^{(n)}(e^{j\omega})| \quad (56)$$

does not exceed the upper bound \bar{r}_L for which the Padé approximation $R_L(\omega)$ is effective.

Let the terms $F_a^{(n)}(\tilde{z})$ ($n = 1, 2, 3, 4$) obtained by partitioning the transfer function $F_a^{(n)}(\tilde{z})$ of the basic filter of Eq. (41) be

$$F_a^{(1)}(\tilde{z}) = z^{-1} b_a^{(1)}(1) \quad (57)$$

$$F_a^{(2)}(\tilde{z}) = z^{-1} (b_a^{(2)}(2) \tilde{z}^{-1} + b_a^{(2)}(3) \tilde{z}^{-2}) \quad (58)$$

$$F_a^{(3)}(\tilde{z}) = z^{-1} (b_a^{(3)}(4) \tilde{z}^{-3} + b_a^{(3)}(5) \tilde{z}^{-4} + b_a^{(3)}(6) \tilde{z}^{-5} + b_a^{(3)}(7) \tilde{z}^{-6}) \quad (59)$$

$$F_a^{(4)}(\tilde{z}) = z^{-1} (b_a^{(4)}(8) \tilde{z}^{-7} + b_a^{(4)}(9) \tilde{z}^{-8} + \dots + b_a^{(4)}(M+1) \tilde{z}^{-M}) \quad (M \leq 20) \quad (60)$$

The frequency compression parameter α is set between 0 and 0.4, and male speech "Nanbu dewa higashi no kaze (East wind in southern part)" and female speech "Dare ka kite kudasai (Somebody should come)" are analyzed. In the worst condition, where quantization width is 0.5 so that the filter coefficient $b_a^{(n)}(m)$ is quantized very roughly, the maximum magnitude $r_{L(n)}$ ($n = 1, 2, 3, 4$) of the term $F_a^{(n)}(m)$ of the transfer function $F_a^{(n)}(z)$ is obtained as

$$r_{L(n)} \lesssim \begin{cases} 3.00 & (n = 1, 2) \\ 2.75 & (n = 3, 4) \end{cases} \quad (61)$$

The value of Eq. (61) is in the range of good approximation for the third-order Padé approximation, and the required transfer function $H_a(\tilde{z})$ of the MLSA filter is given, for $F_a^{(n)}(\tilde{z})$ of Eqs. (57) to (60), as

$$H_a(\tilde{z}) = \exp(b_a^{(0)}(0)) R_3^{(1)}(F_a^{(1)}(\tilde{z})) R_3^{(2)}(F_a^{(2)}(\tilde{z})) \times R_3^{(3)}(F_a^{(3)}(\tilde{z})) R_3^{(4)}(F_a^{(4)}(\tilde{z})) \quad (62)$$

The MLSA filter with the transfer function $H_a(\tilde{z})$ of Eq. (62) can be constructed as in Fig. 2. The frequency compression parameter α can practically be set as 2 to 3 bits. Setting the word length of the coefficient $p_{L,i}^{(n)}$ of Padé approximation $R_L^{(n)}(\omega)$ as 7 bits, and the maximum word lengths of the basic filter

coefficient $b_a^{(n)}(m)$ as 8 bits including the sign bit, the approximation error in the log magnitude of MLSA filter is 0.5 dB or less for all frames in male "Nanbu dewa higashi no kaze," and female "Dare ka kite kudasai." When the basic filter coefficients are not quantized, the approximation error is 0.25 dB or less. The mel spectral envelope of speech is obtained from the improved cepstral method [1]. An example of the approximate mel-log spectral envelope by 8-bit MLSA filter is shown in Fig. 3.

7. Conclusions

This paper presented a construction method for mel-log spectral approximation (MLSA) filter with relatively simple structure and low coefficient sensitivity, together with a design example for MLSA filter for speech synthesis. The transfer function of MLSA filter is given as a rational function with the transfer function of a first-order, all-pass filter as the variable. Consequently, the path without delay must be deleted from the feedback loop in the realization of the filter. Using the construction method for MLSA filter shown in this paper, the path without delay can easily be deleted from the feedback loop. The obtained MLSA filter has a relatively simple structure and low coefficient sensitivity.

Although not stated in this paper, it is verified that the proposed MLSA filter is very useful in synthesis by analysis and synthesis by rule of speech. This will be reported in another paper.

Acknowledgement. A part of this work was supported by the Scientific Research Grant, Min. Education, General B 57460116.

REFERENCES

1. Kitamura and S. Imai. Spectrum distortion and quality of synthesized speech in cepstrum vocoder, Trans. (A) I.E.C.E., Japan, J65-A, 5, pp. 478-484 (May 1982).
2. Imai, S. Logarithmic amplitude approximation (LMA) filter, Trans. (A) I.E.C.E., Japan, J63-A, 12, pp. 886-893 (Dec. 1980).
3. Strube, H. W. Linear prediction on a warped frequency scale, J. Acoust. Soc. Am., 68, 4, pp. 1071-1076 (Oct. 1980).