

람버트 W 함수를 사용한 라플라스 신호의 최소 평균제곱오차 양자화

송현정 · 나상신
아주대학교 전자공학과
전화 : 031-219-2376

The Lambert W Function in the Design of Minimum Mean Square-Error Quantizers for a Laplacian Source

Hyunjung Song and Sangsin Na
Department of Electronic Engineering, Ajou University
E-mail : bananapi@hanmail.net

Abstract

This paper reports that the Lambert W function applies to a non-iterative design of minimum mean square-error scalar quantizers for a Laplacian source. The contribution of the paper is in the reduction of the time needed for the design and the increased accuracy in resulting quantization points and thresholds, because the algorithm is non-iterative and the Lambert W function can be evaluated as accurately as desired.

I. 서론

람버트 W 함수는 $f(x) = xe^x$ 의 역함수로 정의된다[1]. 따라서, 이 정의에 의하면, $xe^x = a$ 의 해는, $x = W(a)$ 이다. 그런데, 실제로 $f(x) = xe^x$ 는 함수 값이 음일 때, 즉 x 값이 음일 경우에는 다대일 함수이므로, 엄밀한 의미에서의 역함수는 존재하지 않는다. 따라서, 이 논문에서는 x 값이 $[-1, \infty)$ 로 한정된 경우의 $f(x) = xe^x$ 의 역함수를 람버트 W 함수로 정의하기로 한다. 이 경우의 람버트 W 함수를, 흔히 람버트 W 함수의 주가지(principal branch)라 한다[1]. 그림1에 이 람버트 W 함수를 도시하였다.

이 함수의 정의역은 $[-e^{-1}, \infty)$ 이며, 연속인 함수이다. 또, 이 람버트 W 함수는 0에서 해석적이고, 따라서 급수로 표현 가능하다. 구체적인 급수 표현식은 다음과

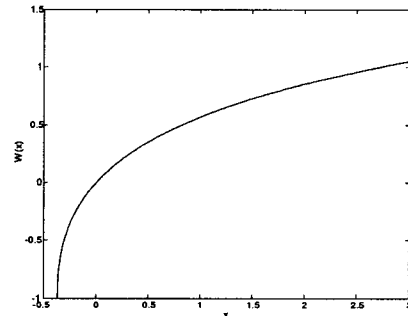


그림1. 람버트 W 함수의 주가지

같다.

$$W(x) = \sum_{n=1}^{\infty} \frac{(-n)^{n-1}}{n!} x^n, \quad |x| < \frac{1}{e}.$$

람버트 W 함수는 Matlab이나 Maple같은 소프트웨어로 구현되어 있어서, 쉽게 계산될 수 있다.

람버트 W 함수는 참고문헌 [1]에 제시된 것처럼 다양한 분야에 응용되고 있다. 이 중 다음과 같은 초월함수를 포함하는 비선형 방정식의 풀이에 사용될 수 있다.

$$a^x = x + b. \quad (1)$$

이 방정식의 풀이는 다음과 같이 람버트 W 함수로 나타내어진다.

$$x = -b - \frac{1}{\ln a} W\left(-\frac{\ln a}{a^b}\right). \quad (2)$$

이 논문에서는 램버트 W 함수가 식(2)의 형태로 라플라스 신호원(Laplacian source)에 최소 평균제곱오차를 갖는 양자기 설계에 쓰인다는 사실을 발견하여 보고하고, 더 나아가 이를 사용하여 평균제곱오차의 의미에서 이 신호원에 최적인 양자기를 설계하는 법을 다룬다.

N -점 홀양자기(scalar quantizer) Q_N 은 N 개의 양자점 y_1, y_2, \dots, y_N 과 경계값 x_1, x_2, \dots, x_{N+1} , 그리고 다대일 함수인 양자함수 $Q_N(\cdot)$ 로 표현할 수 있다. 일반적으로, 경계값은 $x_2 < x_3 < \dots < x_{N-1} < x_N$ 이 되도록 배열하며 $x_1 = -\infty, x_{N+1} = \infty$ 이다. 양자점 y_i 는 구간 $(x_i, x_{i+1}]$ 에 위치하는 대표값이다. 또, 이들과 양자함수의 상호관계는 $x_i < x < x_{i+1}$ 이면, $Q_N(x) = y_i$ 로 나타낼 수 있다. 제 i 번째 양자 구간의 크기는 그 양자구간의 크기(step size)이고, 이를 Δ_i 로 표시하면, $\Delta_i = x_{i+1} - x_i$ 이다.

신호원의 확률밀도함수가 $p(x)$ 일 때, 이 신호원의 N -점 최적 홀양자기 Q_N^* 은 아래 식으로 주어지는 평균제곱오차를 최소화시키는 양자기 Q_N 이다.

$$D(Q_N) = E_p\{(X - Q_N(x))^2\} \\ = \sum_{i=1}^N \int_{x_i}^{x_{i+1}} (x - y_i)^2 p(x) dx.$$

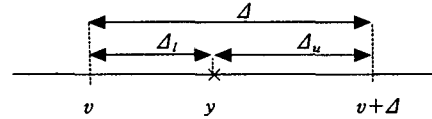
어떤 홀양자기가 최적 양자기가 되기 위해서는 중점 조건과 무게중심 조건을 만족해야 한다. 중점 조건은, 각 경계값이 인접한 두 양자점의 중점이 되어야 한다는 것이고, 무게중심 조건은, 각 양자점이 경계값으로 주어지는 해당 양자구간의 $p(x)$ 에 대한 무게중심이 되어야 한다는 것이다.

신호원의 확률밀도함수에 대한 N -점 최적 홀양자기는 이 두 최적조건에 의해 설계되며, 이러한 조건을 만족시키기 위해서 일반적으로 점화적이며, 반복적인 설계방식을 사용해야 한다. 예를 들어, Lloyd-Max 설계법[3,4]은 최외곽 양자점 y_1 이나 y_N 을 임의의 초기값으로 선택하여 이 두 필요 조건을 만족하도록 반복적이며 점화적으로 인접하는 경계값과 양자점을 결정하도록 함으로써 양자기를 설계한다.

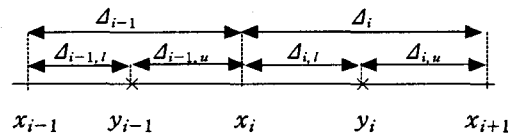
신호원의 분포가 라플라스 확률밀도함수인 경우에는 일반적인 경우와 달리, 비반복적 설계가 가능하다[5,6,7]. 따라서, 라플라스 신호원에 대한 최적 양자기는 설계상의 비반복적 성질 때문에 보다 빠르게 설계될 수 있다. 이러한 비반복적 설계방법을 이 논문에서 재고찰 하였다. 최적 양자기의 양자점과 해당 양자구간의 크기는 램버트 W 함수를 사용한 닫힌 식으로 표현된다는 것을 발견하였다. 램버트 W 함수를 사용한 설계방법의 또 하나의 이점은, 램버트 W 함수의 값을 원하는 정확도로 계산함으로써, 설계된 양자점과 경계값의 정확도를 원하는 만큼 증가시킬 수 있다는 것이다. 이러한 비반복적 방법은 라플라스 확률밀도함수의 특성을 직접적으로 사용하였기 때문에 다른 확률밀도함수에 대해서는 적용되지 않는다.

II. 라플라스 신호원에 대한 최적 양자기의 설계

그림2(a)에 임의의 양자구간 $(v, v + \Delta)$ 을 도시하였다.



(a) 최적 양자구간



(b) 서로 인접해 있는 최적 양자구간
그림2. 양자구간

편의상, $v \geq 0$ 인 영역만을 가정하면, \times 로 표시된 양자점 y 는 확률밀도함수 $p(x)$ 에 대하여 그 해당 양자구간의 무게중심이 되는 점이다. 또, 양자구간의 크기 Δ 는 하위 구간의 크기 Δ_l 과 상위 구간의 크기 Δ_u 로 구분하여 나타내었다. 무게중심 조건에 의하면, y 와 v 의 관계는 다음과 같다.

$$y = \frac{\int_v^{v+\Delta} xp(x) dx}{\int_v^{v+\Delta} p(x) dx}.$$

분산이 1인 라플라스 확률밀도함수 $p(x) = \frac{1}{\sqrt{2}} e^{-\sqrt{2}|x|}$ 가 주어지면, 위 식으로부터 직접적인 계산에 의해 좀더 구체적인 식을 얻을 수 있다. 즉,

$$y = v + \frac{1}{\sqrt{2}} - \frac{\Delta e^{-\sqrt{2}\Delta}}{1 - e^{-\sqrt{2}\Delta}}. \quad (3)$$

그림2(a)에서 $\Delta_u = \Delta - \Delta_l$ 이고, $\Delta_l = y - v$ 이므로, 상위 구간의 크기 $\Delta_u = \Delta - (y - v)$ 이며, 여기에 식(3)을 대입하여 풀면, 다음 식을 얻는다.

$$\Delta_u = \frac{\Delta}{1 - e^{-\sqrt{2}\Delta}} - \frac{1}{\sqrt{2}}.$$

또는, Δ 에 대하여 정리하여,

$$\Delta = \left(\Delta_u + \frac{1}{\sqrt{2}}\right)(1 - e^{-\sqrt{2}\Delta}). \quad (4)$$

이제 방정식(4)를 Δ 에 대하여 풀이하여 궁극적으로 Δ_l 을 Δ_u 로 표현하고자 한다. 즉, 상위 구간의 크기 Δ_u 를 사용하여 하위 구간의 크기 Δ_l 을 나타냄으로써 Δ_u 와

램버트 W 함수를 사용한 라플라스 신호의 최소 평균제곱오차 양자화

Δ_i 의 관계식을 유도해 내하고자 하는 것이다. 이를 위해, 식(4)를 다음과 같이 정리하여 보자.

$$(\Delta_u + \frac{1}{\sqrt{2}})e^{-\sqrt{2}\Delta} = -\Delta + (\Delta_u + \frac{1}{\sqrt{2}}). \quad (5)$$

그리고, 간단히 식을 변형시키면,

$$e^{-\sqrt{2}\Delta - \frac{1}{\sqrt{2}} \ln(\Delta_u + \frac{1}{\sqrt{2}})} = -[\Delta - \frac{1}{\sqrt{2}} \ln(\Delta_u + \frac{1}{\sqrt{2}})] + [\Delta_u + \frac{1}{\sqrt{2}} - \frac{1}{\sqrt{2}} \ln(\Delta_u + \frac{1}{\sqrt{2}})]. \quad (6)$$

이제 식(6)의 각 항을 다음과 같이 치환한다.

$$a = e^{\sqrt{2}},$$

$$t = -(\Delta - \frac{1}{\sqrt{2}} \ln(\Delta_u + \frac{1}{\sqrt{2}})),$$

$$b = (\Delta_u + \frac{1}{\sqrt{2}} - \frac{1}{\sqrt{2}} \ln(\Delta_u + \frac{1}{\sqrt{2}})).$$

그러면, 식(6)은 $a^t = t + b$ 의 형태가 되고, 이는 곧 식(1)과 같은 형태이다. 따라서, 식(6)의 해는 식(2)에 의해서, $t = -b - \frac{1}{\ln a} W(-\frac{\ln a}{a^b})$ 가 되고, 이를 계산하여 정리하면, 최종적으로 다음 식을 얻는다.

$$\Delta = \Delta_u + \frac{1}{\sqrt{2}} + \frac{1}{\sqrt{2}} W[-(\sqrt{2}\Delta_u + 1)e^{-\sqrt{2}\Delta_u + 1}]. \quad (7)$$

또는, $\Delta = \Delta_u + \Delta_i$ 이므로

$$\Delta_i = \frac{1}{\sqrt{2}} + \frac{1}{\sqrt{2}} W[-(\sqrt{2}\Delta_u + 1)e^{-\sqrt{2}\Delta_u + 1}]. \quad (8)$$

이 논문의 주 결과물인, 식(8)은 하위 구간의 크기 Δ_i 이 상위 구간의 크기 Δ_u 와 어떠한 관계를 가지고 있는지를 보여준다. 바로 이 식과 중점 조건에 의하여 최적 라플라스 양자기의 양자구간의 크기를 발생해 낸다. 그림2(b)에 최적 라플라스 양자기의 서로 인접한 두 양자영역을 도시하였다. 중점 최적조건에 의해, 각 $i=2, \dots, N$ 에 대해서 $\Delta_{i-1,u} = \Delta_{i,l}$ 이다. 따라서, 일단 $\Delta_{N,u}$ 가 주어지면, 식(8)에 의해 $\Delta_{N,l}$ 값을 구할 수 있고, 중점 조건에 의해 $\Delta_{N-1,u} = \Delta_{N,l}$ 이며, 식(8)에 의해 다시 $\Delta_{N-1,l}$ 을 얻는다. 이러한 점화적인 방법으로 모든 양자구간의 크기를 구함으로써, 최적 양자기를 설계한다. 이에 대한 좀 더 구체적인 절차를 아래에 제시하였다.

최적 양자구간 크기 발생 라플라스 확률밀도함수에 대하여 최적 양자기가 대칭이라는 것은 참고문헌 [8]에 제시되었다. 따라서, N -점 최적 라플라스 양자기는 음이 아닌 양자점과 경계점만으로 표현할 수 있다. 이제 식(8)과 중점-경계조건을 이용하여 0에서부터 오른쪽 영역에 대한 최적 양자기가 설계되는 과정을 보이고자 한다.

분산이 1인 라플라스 신호원에 대한 N 이 짝수인 최적 양자기 Q_N^* 를 설계하기로 한다.

(a) $\Delta_{N,u} = \infty$.

$$W(0) = 0, \text{ 식(8)에 의하여 } \Delta_{N,l} = \frac{1}{\sqrt{2}}.$$

(b) 중점 조건에 의해 $\Delta_{N-1,u} = \Delta_{N,l}$ 이므로,

$$\Delta_{N-1,u} = \frac{1}{\sqrt{2}}.$$

이제 식(8)에 Δ_u 대신 $\Delta_{N-1,u}$ 를 대입하면,

$$\begin{aligned} \Delta_{N-1,l} &= \frac{1}{\sqrt{2}} + \frac{1}{\sqrt{2}} W(-2e^{-2}) \\ &= \frac{1}{\sqrt{2}} + \frac{1}{\sqrt{2}} (-0.40637573995996) \\ &= 0.41975573975116. \end{aligned}$$

따라서, $(N-1)$ 번째 최적 양자구간의 크기는

$$\Delta_{N-1} = \Delta_{N-1,u} + \Delta_{N-1,l} \text{ 즉}$$

$$0.41975573975116 + \frac{1}{\sqrt{2}} = 1.12686252093771.$$

(c) 위 과정(b)를 필요한 횟수만큼 반복하여 $\frac{N}{2}$ 개의 양자구간의 크기를 구한다. ■

표1에 몇 개의 최적 양자구간의 크기를 제시하였다.

표1. 최적 양자구간의 크기

Δ_N	∞
Δ_{N-1}	1.126862520937
Δ_{N-2}	0.719536169588
Δ_{N-3}	0.533181346531
Δ_{N-4}	0.424567779085
Δ_{N-5}	0.353072423538
Δ_{N-6}	0.302332174845
Δ_{N-7}	0.264412609116
Δ_{N-8}	0.234981296462
Δ_{N-9}	0.211466160906
Δ_{N-10}	0.192241620772
Δ_{N-11}	0.176228966312
Δ_{N-12}	0.162683806207
Δ_{N-13}	0.151075656025
Δ_{N-14}	0.141016144537
Δ_{N-15}	0.132214352273

이 값들은 [5,7]의 결과와 같을 뿐 아니라, 훨씬 높은 정확도를 갖는다. 여기서 특기할 것은, N 값에 상관없이 최적 양자구간의 크기 $\Delta_{N,u} = \infty$, $\Delta_{N,l} = \frac{1}{\sqrt{2}}$ 이고, 따라

서 $\Delta_{N-1,u} = \frac{1}{\sqrt{2}}$ 이며, 즉 $\Delta_{N-1,l} = 0.41975573975116$ 이라는 결론을 보인다는 것이다. 이것은, N 값에 상관없이 최적 양자구간의 크기 $\Delta_N, \Delta_{N-1}, \dots$ 은 변하지 않고 일정한 값을 갖게 된다는 것을 의미한다. 이 점을 고려하면, 표1에서의 양자구간의 크기 값들은 32-점 최적 양자기의 양의 영역, 즉 $\Delta_{32}, \Delta_{31}, \dots, \Delta_{17}$ 인 16 개의 양자구간의 크기 값들을 나타낸 것인데, 이 값들은 64-점

최적 양자기에서 양의 영역 중 가장 오른쪽의 16 개의 양자구간의 크기 값들과 같다.

최적 라플라스 양자기의 생성 편의상, 짝수인 N 에 대해 고찰해 보자. 먼저 몇 개의 양자구간의 크기를 구할 것인지를 결정하고, 0에서부터 첫 번째 양자구간의 크기를 $\Delta_{\frac{N}{2}+1}$ 로 시작하여, $\Delta_{\frac{N}{2}+2}, \Delta_{\frac{N}{2}+3} \dots$ 과 같이 오

른쪽으로 연속적으로 나열함으로써 양자기의 양의 영역을 만들어 낸다. 음의 영역은 대칭성으로부터 구한다.

부록I의 프로그램은 Matlab을 사용하여 N -점 양자기를 구현한 것이다. Matlab에는 프로그램에서 보인 바와 같이 lambertw(\cdot)로 램버트 W 함수를 제공하고 있다.

III. 결론

이 논문에서는 라플라스 신호원에 대한 최소 평균제곱 오차 양자기의 설계를 다루었다. 이 논문의 기여점은, 이 설계법에서 비선형 방정식을 풀어야 하는 대신 램버트 W 함수를 계산하여서 양자기의 설계가 가능한 것을 보인 점이다. 램버트 W 함수는 원하는 만큼의 정확도를 얻도록 계산할 수 있기 때문에, 이 식을 사용하여 양자기 설계시 정확도를 높일 수 있다.

이 연구의 결과는 최적 라플라스 양자기의 여러 성질을 조명하는데 기여할 것으로 사료된다. 첫째, 양자기의 설계 측면에서는, N 이 아주 큰 경우에도 수치적 오류를 줄인 양자기를 설계할 수 있었다. 둘째, 설계된 최적 양자기의 중요변수를 고찰하여 양자기 불일치의 연구에도 도움이 될 수 있을 것이다.

참고문헌

- [1] R.M. Corless, G.H. Gonnet, D.E.G. Hare, D.J. Jeffery, and D.E. Knuth, "On the Lambert W function," *Advances in Computational Mathematics*, vol. 5, no. 4, pp. 329-359, 1996.
- [2] E.M. Lemaray, "Racines de quelques equations transcendante. Integration d'une equation aux differences meeles. Racines imaginaries," *Nouvelles Annales de Mathematiques*, pp. 540-546, 1897.
- [3] S.P. Lloyd, "Least squares quantization in PCM," *IEEE Trans. on Inform. Theory*, pp. 129-136, Mar. 1982.
- [4] J. Max, "Quantizing for minimum distortion," *IRE Trans. on Inform. Theory*, IT-6, pp. 7-12, Mar. 1960.
- [5] K. Nitadori, "Statistical analysis of Δ PCM," *Electron. Commun. in Japan*, vol. 48, pp. 17-26, Feb. 1965.
- [6] H. Lanfer, "Maximum signal-to-noise-ratio

quantization for Laplacian-distributed signals," *Information and System Theory in Digital Communications*, NTG-Report vol. 65, VDE-Verlag GmbH Berlin, Germany, p. 52, 1978.

- [7] P. Noll and Zelinski, "Comments on 'quantizing characteristics for signals having Laplacian amplitude probability density function'," *IEEE Trans. on Comm.*, vol. COM-27, no. 8, pp. 1259-1260, Aug. 1979.
- [8] P.E. Fleischer, "Sufficient conditions for achieving minimum distortion in a quantizer," *Int. Conv. Rec.*, Part 1, pp. 104-111, 1964.

부록1.

%양자점의 개수 N에 대하여, 라플라스 신호원에 대한 양자기 생성.

% 0을 중심으로 오른쪽 영역에 대한 경계값과 양자점출력.

```

clear;
N=16;
for m=1:N
    ntp=2^m;
    npoints=ntp/2; % 오른쪽 영역 양자점의 개수
    dl(npoints)=1/sqrt(2); % 가장 오른쪽 양자구간의 크기
    for i=npoints:-1:2 % 하위 구간의 크기
        tmp=-(sqrt(2)*dl(i)+1);
        dl(i-1)=sqrt(2)/2+lambertw(tmp*exp(tmp))/sqrt(2);
    end
    du(npoints)=Inf; % 상위 구간의 크기 발생

    for i=1:(npoints-1)
        du(i)=dl(i+1); % 중점 조건의 적용
    end
    d=dl+du;

    x(1)=0; % 양자점과 경계값의 발생
    y(1)=dl(1);
    for i=2:npoints
        x(i)=sum(d(1:i-1));
        y(i)=dl(1)+2*sum(dl(2:i));
    end
    result=[];
    result=[x ; y]
end
    
```