*T 125 AUTOMATIZAVIMAS, ROBOTOTECHNIKA*

# Asymptotic Analysis of Optimal Uniform Two-Dimensional Quantization for the Laplace Source

## Z. H. Peric, Z. B. Nikolic
*Faculty of Electronic Engineering, University of Nis Beogradska 14, 18000 Nis, Serbia; e-mail: peric@elfak.ni.ac.vu*
## D. B. Drajic
*Belgrade University, Yugoslavia*

## Indroduction

The asymptotic optimal quantization problem, even for the simplest case - uniform scalar quantization, is very actual nowadays, [1]. The importance of using the rectangular cells and the optimal density (number) of points for product quantization and Gaussian source is considered in [2-3]. In [4] the granular gain (due to cell shape, being 1.53 dB at the maximum) as well as the boundary gain (due to the increase of the dimensions number) was defined showing that the boundary gain dominates at higher dimensions. In [5] the uniform cubic quantization (only the boundary gain) is considered for 8 and 16 dimensions. In this letter, quantizers are designed and analysed under additional constraint – each scalar quantizer is a uniform one.

The optimization of two-dimensional Laplace source quantization is analysed and the existence of a single minimum, depending on the number of points for various levels, is proven. The gain over the optimum uniform scalar quantizer [5] is about (2.8-6.8) dB for rates from 4 to 8 bits per sample (see Fig.1). The resulting gain (obtained using rectangule cells ) can even be compared to boundary gain in highdimensional space.

## Description and optimization

The 2-D (two-dimensional) probability density function for independent identically distributed Laplace random variables (source) with the zero mean and the unity variance is given as

$$f(\mathbf{x}) = \frac{1}{2} e^{-\sqrt{2}\left(|x_1|+|x_2|\right)}, \qquad (1)$$

**x** is the source vector with elements $x_1$ and $x_2$. To simplify the vector quantizer, the Helmert transformation is applied on the source vector giving contours with constant probability densities. The transformation is defined as:

$$r = \frac{1}{\sqrt{2}}\left(|x_1|+|x_2|\right), \quad u = \frac{1}{\sqrt{2}}\left(|x_1|-|x_2|\right). \quad \text{The obtained}$$

probability density function is $f(r,u) = \frac{1}{2} e^{-2r}$. The quantizing cells are rectangular $S_{ij}\left([r_i, r_{i+1}); [u_{ij} u_{ij+1})\right)$, and the representation vector is ($m_i, \hat{u}_{ij}$). For a uniform quantization having $L$ concentric domain:

$$r_i = (i-1)\Delta_L, \ 1 \le i \le L+1; m_i = (i-1/2)\Delta_L,$$
$$1 \le i \le L; \quad \Delta_L = r_{\max} / L$$
$$u_{ij} = (j-1)\Delta_u(i), \ -L \le j \le L+1; \hat{u}_{ij} = (j-1/2)\Delta_u(i),$$
$$-L \le j \le L; \quad \Delta_u(i) = 2m_i / N_i.$$

The number of cells in $i$-th concentric domain is $N_u(m_i)$. Every concentric domain can be subdiveded in four equivalent subdomains, i.e. $N_u(m_i) = 4N_i$ . The number $N_i$ being the same for all the subdomain at $i$-th domainis. The Helmart transform is ortogonal. The MSE (mean-square error) per dimension can be separated into a granular distortion $D_g$ and an overload distortion $D_o$:

$$D_g + D_o = \frac{1}{2}\sum_{i=1}^{L}\sum_{j=1}^{N_u(m_i)}\int_{r_i}^{r_{i+1}}\int_{u_{i,j}}^{u_{i,j+1}}\left[(r-m_i)^2 + (u-\hat{u}_{ij})^2\right]\frac{1}{2}e^{-2r}du\,dr$$

$$+\frac{1}{2}\sum_{j=1}^{N_u(m_L)}\int_{r_{\max}}^{\infty}\int_{u_{L,j}}^{u_{L,j+1}}\left[(r-m_L)^2 + (u-\hat{u}_{Lj})^2\right]\frac{1}{2}e^{-2r}du\,dr \ . \ (2)$$

By using aymptotic approximation, the following is obtained $D_g = \sum_{i=1}^{L}\left[m_i \frac{\Delta_L^2}{12} + \frac{m_i^3}{3N_i^2}\right] \cdot P_r(m_i)$,

where $\qquad P_r(m_i) = 2e^{-2m_i}\Delta_L$ . $\qquad (3)$

The minimization of the function $D_g(\mathbf{N})$ (vector $\mathbf{N}=(N_1,\dots, N_L)$) for fixed number of magnitude levels $L$ constrained by the total number of reconstruction points $N/4$ is formulated in this way: minimize $D_g(\mathbf{N})$ under the constraints

$$g_0(P_1, P_2, \cdots P_L) = N - \sum_{i=1}^{L} P_i \ge 0$$

$g_i(P_i) = P_i \geq 0; i = 1, 2, \cdots L$ .

We prove that the problem of minimization of the $D_g(\mathbf{N})$ is a convex programming problem. This follows directly from Lemma 1.

*Lemma 1*: Function $D_g(\mathbf{N})$ is convex and constraints $g_0(\mathbf{N})$ and $g_i(N_i)$ form the convex set.

*Proof of Lemma 1*

To prove that the function $D_g(\mathbf{N})$ is convex and that constraints $g_0(\mathbf{N})$ and $g_i(N_i)$ form the convex set it is sufficient to prove that Hessian matrices of the following functions: $D_g(\mathbf{N})$, $-g_0(\mathbf{N})$, $-g_i(N_i)$ $1 \leq i \leq L$ are positive semi-definite [6,p27].

We find partial derivative for $\mathbf{N}$ from (3) and it is:

$$\frac{\partial D_g}{\partial N_i} = -\frac{2P_r(m_i)}{3N_i^3} m_i^3 ,$$

while the second partial derivative is:

$$\frac{\partial^2 D_g}{\partial N_i \partial N_j} = \begin{cases} \dfrac{2P_r(m_i)}{N_i^4} m_i^3 , & i = j, \\ 0, & i \neq j, \end{cases}$$

from which it follows that

$$\Rightarrow \frac{\partial^2 D_g}{\partial N_i N_j} \geq 0.$$

For Hessian matrix it obviously holds

$$\begin{bmatrix} \dfrac{\partial^2 D_g^e}{\partial P_1 \partial P_1} \cdots \dfrac{\partial^2 D_g^e}{\partial P_1 \partial P_L} \\ \cdots\cdots\cdots\cdots\cdots \\ \dfrac{\partial^2 D_g^e}{\partial P_L \partial P_1} \cdots \dfrac{\partial^2 D_g^e}{\partial P_L \partial P_L} \end{bmatrix} \geq 0,$$

while for the constraints we have

$$-\frac{\partial^2 g_0}{\partial P_i \partial P_j} = 0 , \quad -\frac{\partial^2 g_i}{\partial P_i \partial P_j} = 0; i = 1, 2, \cdots, L$$

$$\begin{bmatrix} \dfrac{\partial^2 g_0}{\partial P_1 \partial P_1} \cdots \dfrac{\partial^2 g_0}{\partial P_1 \partial P_L} \\ \cdots\cdots\cdots\cdots \\ \dfrac{\partial^2 g_0}{\partial P_L \partial P_1} \cdots \dfrac{\partial^2 g_0}{\partial P_L \partial P_L} \end{bmatrix} = 0 , \quad -\frac{\partial^2 g_i}{\partial P_i \partial P_j} = 0; i = 1, 2, \cdots, L.$$

This completes the proof.

Optimal solution is found applying the method of Lagrange multipliers.

$$J = \sum_{i=1}^{L} \left[ m_i \frac{\Delta_L^2}{12} + \frac{m_i^3}{3N_i^2} \right] \cdot P_r(m_i) + \lambda \sum_{i=1}^{L} N_i ,$$

$$\frac{\partial J}{\partial N_i} = -\frac{2}{3N_i^3} m_i^3 \cdot P_r(m_i) + \lambda = 0$$

yielding : $N_i = \dfrac{N}{4} \dfrac{m_i \sqrt[3]{P_r(m_i)}}{\sum_{j=1}^{L} m_j \sqrt[3]{P_r(m_j)}}$ for fixed $N$. For the uniform quantization, after some manipulation the following is obtained $N(r) \approx \dfrac{N r_{max} r \sqrt[3]{g(r)}}{4L \int_0^{r_{max}} r \sqrt[3]{g(r)} dr}$ where the sum is approximated by the integral, and $g(r) = e^{-2r}$. Returning to (3) the final expression for $D_g$ is:

$$D_g = \frac{r_{max}^2}{6L^2} \int_0^{r_{max}} r g(r) dr + \frac{32L^2}{3N^2 r_{max}^2} \left( \int_0^{r_{max}} r \cdot \sqrt[3]{g(r)} dr \right)^3 . \quad (4)$$

In such a way a pretty simple solution for granular distortion is :

$$D_g = \frac{r_{max}^2}{6L^2} I_0 + \frac{32}{3} \frac{L^2}{N^2 r_{max}^2} I^3 , \quad (5)$$

where $I_0 = \dfrac{1}{4} - \dfrac{1}{4} e^{-2r_{max}} - \dfrac{1}{2} r_{max} e^{-2r_{max}}$ and

$$I = \frac{9}{4} - e^{-\frac{2r_{max}}{3}} \left( \frac{9}{4} + \frac{3r_{max}}{2} \right) .$$

The optimal number of levels problem can be solved analytically only for the asymptotical analysis from the condition $\dfrac{\partial D_g}{\partial L} = 0$. The optimal solution is $L_{opt} = r_{max} \sqrt[4]{\dfrac{I_0 N^2}{64 I^3}}$. Substituting this expression into (5), granular distortion becomes:
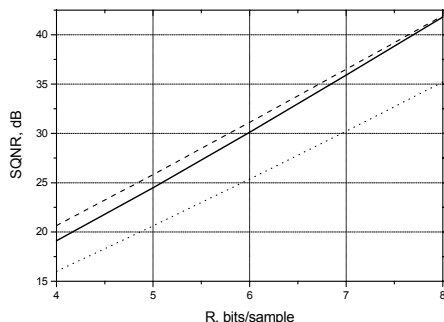
$$D_g^{opt} = \frac{8}{3N} I \sqrt{I \cdot I_0} . \quad (6)$$

This expression is similar to the optimal expressions for the distortion of nonuniform asymptotic quantizers [3]. The expression for $D_g^{opt}$ can be interpreted as Zador-Gersho formula for the uniform quantizer and a 2-D Laplace source.

For comparing the obtained results to the previous ones, $r_{max}$ from [5] is used, being obtained for 1-D approach. The corresponding $D_g^{scal}$ is compared to the obtained result using the following gain definition $G = 10 \log(D_g^{scal} / D_g^{opt})$. The performance gain obtained by our method over the uniform scalar quantization for different rates can be presented in this manner: for R=4. G=2.68dB; for R=5, G= 3.31dB and for R=8, G=6dB.

Exactly optimal value for $r_{max}$ is obtained by repeating our optimisation method for different $r_{max}$ and choosing the values for which mse=$D_g$+$D_o$ is minimal. In Fig.1 signal-to-quantization noise ratio

$SQNR = 10\log(1/mse)$ vs. the number of bits per sample R is shown. It can be concluded that the propsed 2-D quantizer is even comparable to 16-D uniform quantizer from [5], also shown in the figure.



**Fig. 1.** SQNR vs. the number of bits per sample for optimum scalar quantization [5], 16-D uniform quantization [5] and 2-D uniform quantizationan on Laplace source, ____2-D uniform quantization, -----16-D uniform quantization, …… optimum scalar quantization

## Conclusion

The optimization of 2-D Laplace source uniform quantization is carried out and the existence of a single minimum depending on the number of points on various levels is proven. Simple expression for granular distortion in closed form is obtained. The results obtained by the asymptotic analysis demonstrate the significant performance gain over the uniform scalar quantization (even 6.8dB for R=8). The obtained gain using rectangule cells can even be compared to boundary gain in highdimensional space. That automatically provides lower complexity and easier realization.

## References

1. **Hui D., Neuhoff D.L.** Asymmptotic Analysis of Optimal Fixed-Rate Uniform Scalar Quantization // IEEE Trans. – 2001. – **IT-47**(3). – P. 957-977.
2. **Gray R.M., Neuhoff D.L.** Quantization // EEE Trans. – 1998 – **IT-44**(6). – P. 2325-2384.
3. **Na S., Neuhoff D.L.** Bennett's Integral for Vector Quantizers // IEEE Trans. – 1995. – **IT-41**(4). – P. 886-900.
4. **Eyuboglu M.V., Forney G.D.** Lattice and Trellis Quantization with Lattice-and Trellis-Bounded Codebooks-High-Rate Theory for MemorylessSources // IEEE Trans. – 1993. – **IT-39**(1) – P. 46-60.
5. **Jeong D.G., Gibson J.** Uniform and Piecewise Uniform Lattice Vector Quantization for Memoryless Gaussian and Laplacian Sources // IEEE Trans. – 1993. – **IT-39**(3). – P. 786-804.
6. **Himmelblau D.M.** Applied Nonlinear Programming. - McGraw-Hill, Inc., USA, 1972.

**Z. H. Peric, Z. B. Nikolic, D. B. Drajic. Nekintančio dvimačio gausinio signalo kvantizacijos asimptotinė analizė // Elektronika ir elektrotechnika. – Kaunas: Technologija, 2004. – Nr. 2(51). – P. 5-7.**

Pateikiama paprasta ir kompleksinė dvimačio kvantizatoriaus analizė. Įrodoma, kad Laplaso šaltinių signalams būdinga kvadratinė priemonės klaida. Optimizuojant kvantizatoriaus parametrus gaunama 6,8 dB klaida esant 8 bitų santykiui. Nustatyta, kad, naudojant stačiakampę ląstelę gaunamas didelis stiprinimas aukšto lygio dimensiniuose paviršiuose. Pateikiamas teorinių ir eksperimentinių rezultatų palyginimas. Il. 1, bibl.6 (anglų kalba; santraukos lietuvių, anglų ir rusų k.).


**Z. H. Peric, Z. B. Nikolic, D. B. Drajic. Asymptotic Analysis of Optimal Uniform Two-Dimensional Quantization for the Laplace Source // Electronics and Electrical Engineering. – Kaunas: Technologija, 2004. – No. 2(51). – P. 5-7.**

The simple and complete asymptotical analysis is given for a uniform two-dimensional quantizer for Laplace source with the respect to the mean-square error. The significant gain over the optimum uniform scalar quantization is obtained (about 6.8 dB for the rate of 8 bits per sample). The obtained gain using rectangular cells can even be compared to the boundary gain in highdimensional space. Ill. 1, bibl. 6 (in English; summaries in Lithuanian, Russian, English).


**З.Г. Перич, З.Б. Николич, Д.Б. Драич. Асимптотический анализ постоянного квантизации постоянного двухмерного гаусовского сигнала // Электроника и электротехника. – Каунас: Технология, 2004. – № 2(51). – С. 5-7.**

Представлен простой и сложный анализ двухмерного квантизатора. Показано, что для сигналов источников Лапласа характерна квадратическая ошибка устройства. При оптимизации параметров квантизатора получается ошибка порядка 6,8 дБ, когда соотношение 8 битов. При применении прямоугольной цели устройства появляется выигрыш для высокоуровной димензионной поверхности. Приводится сравнение теоретических и экспериментальных результатов. Ил. 1, библ. 6 (на английском языке; рефераты на литовском, английском и русском яз.).