

Effects of signal quantization in digital filtering

Lab 11, SDP

Table of contents

1 Objective	1
2 Theoretical notions	1
2.1 Binary representation of fractionary numbers	1
3 Theoretical exercises	3
4 Practical exercises	3
5 Final questions	3

1 Objective

Students should observe the effects of internal quantization errors on the output signal of a digital filter.

2 Theoretical notions

2.1 Binary representation of fractionary numbers

TBD

①

$$\begin{array}{cccccccc} & 4 & 3 & 2 & 1 & 0 & -1 & -2 & -3 & -4 \\ & 2^4 & 2^3 & 2^2 & 2^1 & 2^0 & 2^{-1} & 2^{-2} & 2^{-3} & 2^{-4} \\ 1 & 1 & 0 & 1 & 1 & . & 0 & 1 & 0 & 1 \end{array} = 26.3125$$

$$16 + 8 + 2 = 26 \quad \frac{2^{-2}}{2^2} + \frac{2^{-4}}{2^4} = \frac{1}{4} + \frac{1}{16} = 0.25 + 0.0625 = 0.3125$$

② 273.21875

$$273 = 256 + 16 + 1$$

$$\quad \quad \quad 2^8 \quad \quad 2^4 \quad 2^0$$

$$100010001$$

$$\begin{array}{lcl} 0.21875 \times 2 & = & 0.43750 \\ 0.4375 \times 2 & = & 0.87500 \\ 0.875 \times 2 & = & 1.75 \\ 0.75 \times 2 & = & 1.5 \\ 0.5 & = & 1.0 \end{array}$$

↓

273	2	1	↑
136	2	0	
68	2	0	
34	2	0	
17	2	1	
8	2	0	
4	2	0	
2	2	0	
1	2	1	
0			

$$0.21875 : 0.00111$$

$$273.21875 : 100010001.00111$$

Figure 1: Binary representation of fractionary numbers

3 Theoretical exercises

1. Consider the system with the following difference equation:

$$y[n] = \frac{1}{2}y[n-1] + x[n]$$

Compute the first 6 samples of the response to the input signal $x[n] = \left(\frac{1}{4}\right)^n$, in three different ways:

- a. Computations with infinite precision
- b. Computations with fixed-point 1S0I4F format, quantize by truncation
- c. Computations with fixed-point 1S0I4F format, quantize by rounding

4 Practical exercises

1. In Matlab, write a script file to study the quantization of the `mtlb` signal on $N = 8$ bits.
 - a. Load the predefined `mtlb` signal (use `load mtlb`);
 - b. Figure out if a sign bit is needed or not;
 - c. Find the maximum absolute value of the signal, and figure out the number of bits required for the integer part. The remaining bits will be allocated to the fractionary part;
 - d. Use the function `fixdt()` to create the corresponding fixed-point data type;
 - e. Use the function `num2fixpt` to convert the signal `mtlb` to the fixed-point data type, using all the three quantization methods;
 - f. For all the three quantization methods, visualize the quantized signal, the quantization error, and compute the total energy of the quantization errors. Which quantization error produces minimum errors?
 - g. Play the quantized signal. Can you hear the difference from the original signal?
2. In Matlab, create a function to implement the system from exercise 1. The values shall be quantized after each multiplication and addition. Apply at the system input the quantized signal from exercise 2, and display the output.

5 Final questions

1. TBD