

# REDUCEREA ERORII DE ROTUNJIRE ÎN ARITMETICA CU VIRGULĂ FIXĂ

conf. dr. ing. Ovidiu Radu  
dr. ing. Alexandru Vasile

Universitatea Politehnica București

**Rezumat:** Una din tehnicile utilizate pentru reducerea erorilor de rotunjire la filtrele digitale recursive o constituie reacția de eroare (RE). Aceasta corespunde faptului că evaluarea unui polinom prin metoda Horner este echivalentă cu un filtru recursiv, de ordinul întâi. RE este o tehnică utilizată pentru a reduce erorile de rotunjire, care apar în evaluarea unei funcții prin aproximare polinomială.

**Cuvinte cheie:** reacție de eroare, eroare de rotunjire, aproximare polinomială, filtru recursiv, funcție exponențială, funcție logaritmică.

## 1. Introducere

Când se evaluează unele funcții, ca cele de tip exponențial, se utilizează în mod uzual aproximarea polinomială. Un polinom este aproximat în mod curent prin metoda Horner, similar funcționării unui filtru recursiv de ordinul întâi.

Pe de altă parte, aprecierea erorilor spectrului (ES) prin reacția de eroare (RE) este o tehnică importantă, care poate reduce erorile de rotunjire, provocate de efectele lungimii finite a cuvântului asupra filtrelor digitale recursive cu virgulă fixă.

În consecință, se întrevide posibilitatea reducerii erorii prin tehnica RE, când se evaluează o funcție prin aproximarea polinomială, folosind aritmetica cu virgulă fixă.

În continuare, se arată în partea a II a, relația dintre metoda Horner și un filtru recursiv, de ordinul întâi cu RE.

În partea a III a, se prezintă câteva exemple de evaluare a unor funcții exponențiale și logaritmice cu RE.

## 2. Metoda Horner și filtrul recursiv cu RE

Un polinom,  $f(x)$ , de ordinul  $N$ , poate fi evaluat prin metoda Horner, astfel:

$$f(x) = (\dots((c_N x + c_{N-1})x + c_{N-2})x + \dots + c_1)x + c_0, \quad (1)$$

Ecuția (1) poate fi evaluată prin algoritmul recursiv, astfel:

$$f[-1] := 0$$

$$\text{pentru } k=0 \text{ la } N, \quad f[k] := x * f[k-1] + c[N-k].$$

În finalul acestei iterații, se obține  $f[N]$ . În acest algoritmul, dacă se consideră că  $c[N-k]$  este semnalul de intrare,  $f[k]$  este semnalul de ieșire și  $x$  este un coeficient, atunci el reprezintă ecuația unui filtru recursiv de ordinul întâi. Adică, ieșirea de ordinul  $N$  a acestui filtru reprezintă evaluarea unui polinom  $f(x)$ , de ordinul  $N$ .

Se știe că, dacă se implementează un filtru recursiv prin aritmetică cu virgulă fixă, pot apărea erori, dacă rezultatele multiplicărilor sunt cuantizate, [1], [2].

RE reprezintă una din tehnicile utilizate pentru a reduce erorile de cuantizare.

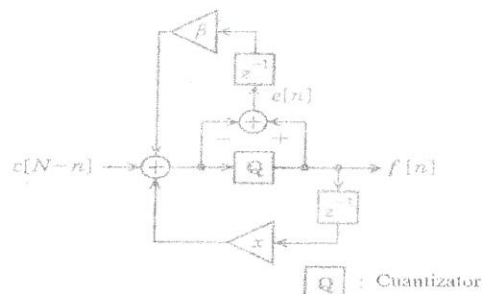


Figura 1. Filtru recursiv de ordinul întâi cu reacție de eroare

În figura 1, se dă schema bloc a unui filtru recursiv de ordinul întâi, cu RE, unde  $e[n]$  este eroarea de cuantizare și  $\beta$  reprezintă coeficientul reacției de eroare.

Dacă se limitează  $\beta$  la partea întreagă, pentru a nu crește numărul de operații aritmetice, se arată în [3] că cele mai bune valori sunt :

$$\beta = \begin{cases} \dots 0 \dots\dots\dots |x| \leq 0.5 \\ -1 \dots\dots\dots 0.5 < x < 1 \\ \dots 1 \dots\dots\dots -1 < x \leq -0.5 \end{cases} \quad (2)$$

În consecință, algoritmul pentru evaluarea RE este următorul:

```
f[-1]: =0
e[-1]: =0
pentru k: =0 la N
begin
w: =x* f [ k-1 ] + c [N- k ] +  $\beta$  * e [ k - 1 ] ;
f [ k ] : = Q ( w ) ;
e [ k ] : = f [ k ] - w ;
end
```

unde  $Q( )$  reprezintă o funcție de cuantizare.

### 3. Exemple pentru evaluarea funcțiilor cu RE

Se vor prezenta două exemple de evaluare a funcțiilor, prin polinoame cu RE:

- exponențială,  $2^{x-1}$ , cu  $0 \leq x < 1$ ;
- logaritmică,  $\log_2 (x + 1)$ , cu  $0 \leq x < 1$ ;

#### A. Aproximări polinomiale

Se presupune că evaluarea este efectuată de un procesor de 16 biți cu virgulă fixă, că partea fracționară a datelor este de 15 biți și că lungimea rezultatului operațiilor de tipul  $Ax + B \rightarrow C$  din multiplicator este  $16b \times 16b + 32b \rightarrow 32b$ .

Deoarece partea fracționară a datelor este de 15 biți, se presupune că mărimea maximă a erorii absolute tolerate este valoarea cifrei 1 a bitului cel mai puțin semnificativ (LSB), adică  $2^{-15} = 3,0518 \cdot 10^{-5}$

Setul de coeficienți pentru aproximări polinomiale se obține astfel:

1. Se calculează setul de coeficienți prin minimizarea aproximării maxime, utilizând calculul cu virgulă mobilă și precizie dublă; maximul erorii amplitudinii poate fi inferior valorii  $3,0518 \cdot 10^{-5}$ .
2. Se cuantizează setul de coeficienți la 15 biți. Această cuantizare se efectuează prin rotunjire sus sau jos, astfel încât să minimizeze eroarea maximă a amplitudinii, când se evaluează polinoame, utilizând calculul cu virgulă mobilă și precizie dublă.
3. Dacă eroarea maximă a amplitudinii pentru un set de coeficienți cuantizați este mai mare decât  $3,0518 \cdot 10^{-5}$ , se mărește ordinul polinoamelor de aproximare și se reia punctul 1.

În continuare, sunt prezentate polinoame de aproximare și coeficienții lor; acești coeficienți sunt dați prin fracții pentru a putea fi mai precis memorati în procesor.

1. Exponențial:  $2^{x-1}$ ,  $0 \leq x < 1$

$$2^{x-1} = \sum_{k=0}^4 c_k x^k$$

$$c_1 = 11354/32768$$

$$c_2 = 3959/32768$$

$$c_3 = 847/32768$$

$$c_4 = 224/32768$$

(3)

2. Logaritmice:  $\log_2(x+1)$ ,  $0 \leq x < 1$

$$\log_2(x+1) = \sum_{k=0}^6 c_k x^k$$

$$c_0 = 0/32768$$

$$c_1 = 47268/32768$$

$$c_2 = 23520/32768$$

$$c_3 = 14959/32768$$

$$c_4 = -9062/32768$$

$$c_5 = 3965/32768$$

$$c_6 = -842/32768$$

(4)

Eroarea absolută maximă a amplitudinii pentru (3) și (4) este  $0,609 \cdot 10^{-5}$ , respectiv  $0,869 \cdot 10^{-5}$ , când evaluarea se face prin calcul cu virgulă mobilă și precizie dublă.

### B. Exemple de evaluare a funcțiilor

Coeficientul reacției de eroare este ales  $-1$  sau  $0$ , conform cu (2). În consecință, se va utiliza rotunjirea, pentru cuantizare în figura 1.

În figura 2, se arată eroarea lui  $2^{x-1}$  funcție de  $x$ , fără RE, (2a) și cu RE, (2b). Erorile sunt calculate pentru  $2^{13} = 8192$  valori ale lui  $x$ , distribuite uniform, în intervalul  $0 \leq x < 1$ .

Din figura 2, rezultă că erorile amplitudinii depășesc valoarea  $2^{-15}$ , crescând o dată cu apropierea lui  $x$  de valoarea 1.

Din figura 2b, reiese că erorile amplitudinii nu depășesc valoarea  $2^{-15}$  și chiar se reduc, când  $x=1$ . Rezultă eficiența utilizării RE.

În figura #, este reprezentată distribuția erorii lui  $2^{x-1}$  pentru orice  $x$ , care este reprezentat prin cuvinte de 16 biți lungime, adică pentru  $2^{15} = 32768$  valori uniform distribuite în intervalul  $0 \leq x < 1$ .

Figura 3a, (fără RE), arată că există cuvinte ale căror erori în amplitudine sunt mai mari decât  $2^{-15}$ .

Figura 3b, (cu RE), arată că nu există cuvinte a căror erori în amplitudine să depășească  $2^{-15}$ .

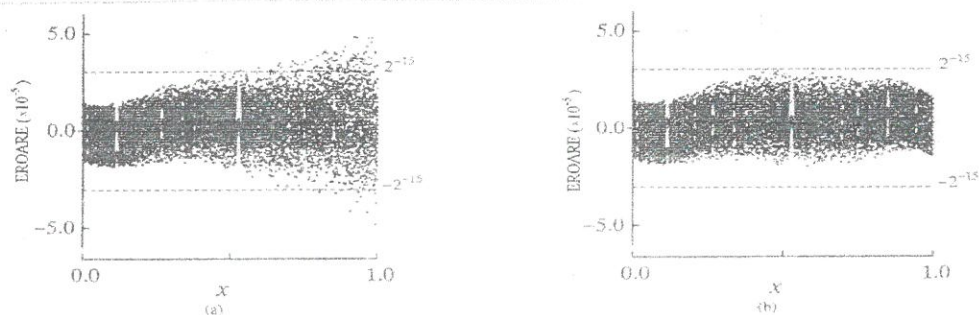


Figura 2. Eroarea funcției  $2^{x-1}$ : (a) fără RE; (b) cu RE

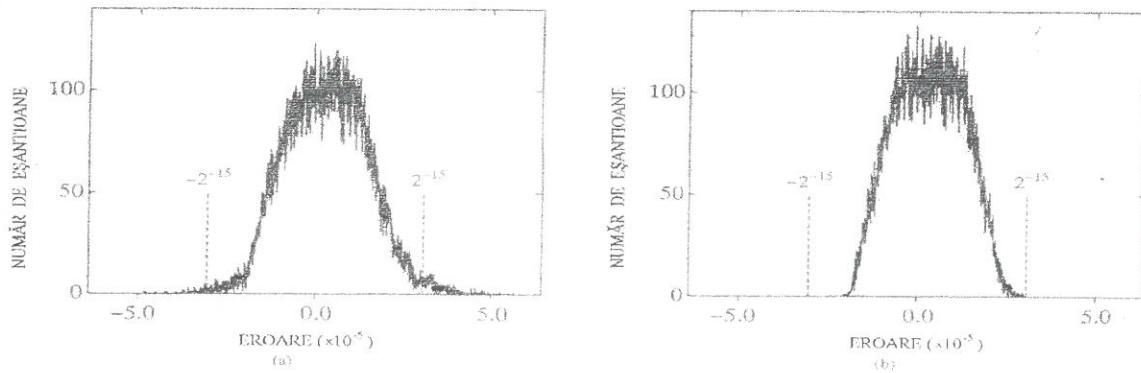


Figura 3. Distribuția erorii funcției  $2^{x-1}$ : (a) fără RE; (b) cu RE

Tabelul I

	$2^{x-1}$		$\log_2(x+1)$	
	Cu RE	Fără RE	Cu RE	Fără RE
Eroarea maximă a amplitudinii ( $\times 10^{-5}$ )	2,978	5,655	3,324	6,363
Abaterea relativă a datelor(%)	0,00	1,46	0,05	2,59

Tabelul I cuprinde eroarea maximă a amplitudinii și abaterea relativă a datelor față de valoarea  $2^{-15}$ .

Rezultatele prezentate mai sus arată că, atunci când se evaluează o funcție prin aproximare polinomială utilizând aritmetică cu virgulă fixă, RE poate fi o bună metodă pentru anumite funcții de reducere a erorilor care apar în operațiile aritmetice.

## Concluzii

1. Algoritmul de evaluare a unui polinom folosind metoda Horner este echivalent cu realizarea unui filtru recursiv de ordinul întâi.
2. Reacția de eroare, care este una din metodele utilizate pentru reducerea erorilor care apar în operațiile aritmetice ale filtrelor recursive, poate fi aplicată cu succes pentru reducerea erorilor de rotunjire, în evaluarea funcțiilor prin aproximarea polinomială.

## Bibliografie

1. RADU, O. : Filtre numerice . Aplicații, Editura Tehnică, București, 1979.
2. CAVICCHI, T. J.: Digital Signal Processing, John Wiley, New York, 2000.
3. ELLIOT, D. F. Handbook of Digital Signal Processing, New York Academic, 1987.