

CAPITOLUL 5

EFECTELE LUNGIMII FINITE A CUVINTELOR ÎN FILTRAREA DIGITALĂ

5. 1. Introducere

Teoria filtrelor digitale s-a bazat pe presupunerea că atât semnalele, cât și parametrii filtrelor pot avea orice valoare finită. În realitate, datorită limitării lungimilor cuvintelor din orice sistem digital, sunt permise numai valori discrete ale amplitudinii semnalelor, respectiv coeficienților. Luând în considerație aceste valori discrete în relațiile care caracterizează filtrele, vor rezulta ecuații neliniare, care, în general, nu vor putea fi riguros prelucrate.

Implementarea sistemelor discrete, fără a considera efectele lungimii finite a cuvintelor, inerente în orice implementare digitală, a condus la obținerea unor caracteristici liniare. De fapt, au fost analizate sisteme modelate liniar, dar ale căror realizări digitale sunt implicit neliniare. Această problemă reprezintă un dezavantaj major al filtrelor digitale și, prin urmare, analiza efectelor lungimii finite a cuvintelor asupra performanțelor filtrelor constituie o etapă importantă în proiectarea filtrelor digitale.

În cazul filtrelor recursive, caracteristicile neliniare rezultate din operația de cuantizare din multiplicatoare, pot cauza un comportament oscilatoriu la ieșirea filtrelor, chiar și în absența semnalului de intrare. Mai mult, în sumatoare poate apărea depășirea aritmetică care produce, de asemenea, oscilații la ieșire.

În cazul calculatoarelor care lucrează cu lungimi mari ale cuvintelor (adică au un număr mare de biți disponibili pentru reprezentarea numerelor), efectele cuantizării pot fi nesemnificative. Acestea cresc cu descreșterea numărului de biți. Din acest motiv sunt necesare modele matematice care să permită estimarea efectelor lungimii finite a cuvintelor asupra performanțelor filtrelor. Un model simplu este

cel care se bazează pe presupunerea că erorile de cuantizare sunt mici în comparație cu nivelul semnalului sau al parametrului, adică este o cuantizare „fină” în care erorile pot fi tratate ca zgomot și problema devine liniară [23].

Principalele tipuri de erori de cuantizare care apar în filtrarea digitală sunt:

1. Erori de cuantizare ale semnalului de intrare în conversia analog – digitală (A/D);
2. Erori rezultate din cuantizarea coeficienților filtrelor digitale;
3. Erori rezultate din rotunjirea produselor;
4. Depășirea aritmetică;
5. Oscilații cu cicluri limită.

Dintre aceste tipuri de efecte, erorile de cuantizare ale semnalului de intrare au loc în afara filtrului, înaintea calculelor interne, restul efectelor sunt interne filtrului și influențează metoda prin care sistemul va fi implementat.

De exemplu, pentru un filtru digital de ordinul întâi

$$y[n] = Ay[n-1] + x[n] \quad (5.1)$$

eroarea de tipul 1 se referă la cuantizarea intrării $x[n]$, eroarea de tipul 2 apare în reprezentarea parametrului A iar cea de tipul 3 apare la formarea produsului $Ay[n-1]$, necesar la fiecare iterație.

Elementul de bază dintr-un calculator numeric este circuitul cu două stări echiprobabile, căruia i se asociază o informație de 1 bit. N astfel de dispozitive pot fi cascade pentru a forma un registru care conține N biți de informație. Implementarea unui filtru digital recursiv de ordinul întâi descris de ecuația (5.1) și redată în figura 5.1, ilustrează cele mai importante operații ce trebuie efectuate.

Ieșirea anterioară $y[n-1]$ este stocată în registrul de ieșire sub forma unui număr pe N biți. Acesta este multiplicat cu numărul pe N biți care reprezintă coeficientul A care a fost stocat în registrul pentru coeficienți. Produsul $Ay[n-1]$ (după rotunjire la N biți) este adunat la intrarea curentă $x[n]$ (de asemenea un număr pe N biți) pentru a forma ieșirea actuală $y[n]$ care este stocată pentru multiplicare cu A în iterația următoare. Întreaga procedură începe cu o valoare inițială $y[-1]$ stocată în registrul de ieșire. Aceasta poate fi sau nu, egală cu zero. Filtrele de ordin superior pot fi implementate într-un mod similar.

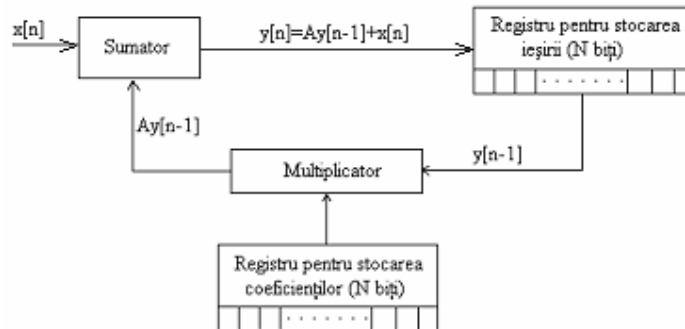


Figura 5.1. Implementarea unui filtru recursiv de ordinul întâi

Diferitele structuri de implementare ale unui sistem descris de ecuații cu diferențe cu coeficienți constanți sunt echivalente dacă furnizează aceeași ieșire pentru o intrare dată, presupunând calculele interne ca fiind efectuate cu precizie infinită. Acestea nu sunt echivalente când sunt realizate cu precizie finită.

Trei factori importanți contribuie la alegerea unei anumite realizări a filtrelor:

- complexitatea calculelor,
- necesarul de memorie,
- efectele lungimii finite a cuvintelor.

Efectul lungimii finite a cuvintelor reprezintă un factor important în implementarea sistemelor digitale de prelucrare a semnalelor și trebuie luat în calcul la realizarea filtrelor digitale, deoarece limitarea numărului de biți conduce la degradarea performanțelor filtrelor digitale. Înainte de a examina aceste efecte, se va prezenta o scurtă introducere în aritmetica digitală.

5.2. Reprezentarea numerelor

În procesarea digitală a semnalelor analogice, eșantioanele semnalului analogic sunt reprezentate în format digital. În principiu, procesul de conversie A/D implică eșantionarea semnalului analogic și reprezentarea eșantioanelor ca secvențe de biți care definesc amplitudinea cuantizată a semnalului. Principala caracteristică a aritmeticii digitale constă în numărul limitat (de obicei fix) de biți folosiți în reprezentarea numerelor. Această constrângere are ca rezultat precizia finită a

calculelor, care conduce la erori și efecte neliniare în comportamentul filtrelor digitale.

În cadrul reprezentării binare a numerelor reale sunt mai multe metode prin care un eșantion al unui semnal analogic poate fi reprezentat în format binar. Clasa reprezentărilor binare poate fi împărțită în reprezentările în virgulă fixă, virgulă mobilă și virgulă mobilă cu blocuri.

5.2.1. Reprezentarea numerelor în virgulă fixă

Reprezentarea numerelor în virgulă fixă este generalizarea reprezentării zecimale, în care numerele din stânga virgulei reprezintă partea întreagă a numărului, iar cele din dreapta virgulei, partea fracționară.

$$x = (b_{-a} \dots b_{-1} b_0, \dots b_b)_r = \sum_{i=-a}^b b_i r^{-i} \quad 0 \leq b_i \leq (r-1) \quad (5.2)$$

unde b_i reprezintă cifra, r – baza, $a+1$ – numărul de cifre ale părții întregi și b – numărul de cifre ale părții fracționare.

Datorită vitezei și costului scăzut al părții hard asociate, reprezentarea în virgulă fixă este deseori preferată în computere mai puțin performante și în circuite dedicate care lucrează în timp real. Cea mai cunoscută reprezentare este cea pentru care $r=2$, în care numerele b_i se numesc *numere binare* sau *biți* și pot lua valorile $\{0,1\}$, obținându-se *codul binar natural direct*. „Virgula binară” dintre b_0 și b_1 nu există fizic în calculator. Circuitele logice ale acestuia sunt proiectate astfel încât calculele să aibă ca rezultat numere ce corespund poziției virgulei binare. Totuși, în cele ce urmează, se va folosi virgula pentru a sublinia caracterul fracționar al numărului reprezentat.

Folosind un format întreg pe n biți ($a=n-1, b=0$), se pot reprezenta întregi fără semn cuprinși în domeniul $0 \div (2^n-1)$. De obicei se folosește formatul fracționar ($a=0, b=n-1$), cu virgula binară între b_0 și b_1 , care permite reprezentarea numerelor în domeniul $0 \div (1 - 2^{-n})$.

Indiferent dacă codul binar reprezintă o fracție, un întreg, sau ambele, primul bit din stânga este numit cel mai semnificativ bit (most significant bit, MSB) iar bitul cel mai din dreapta, cel mai puțin semnificativ bit (least significant bit, LSB). În reprezentarea unei fracții, MSB are o pondere de $2^{-1}=1/2$ iar LSB are o pondere de $2^{-b}=1/2^b$, unde b este numărul de biți pe care este reprezentată fracția. Ponderea $2^{-b}=1/2^b$ desemnată de LSB este numită și *rezoluție*.

Orice întreg sau număr cu parte întreagă și fracționară poate fi reprezentat în format fracționar prin factorizarea termenului r^a în relația (5.2). În această notație un cuvânt de cod de $a+1$ biți, cum ar fi 10011, corespunde numărului întreg

$$A = 1 \cdot 2^0 + 1 \cdot 2^1 + 0 \cdot 2^2 + 0 \cdot 2^3 + 1 \cdot 2^4 = 1 + 2 + 16 = 19$$

Pe de altă parte, numărul 0,10011 reprezintă o fracție corespunzătoare numărului zecimal

$$B = 1 \cdot 2^{-1} + 0 \cdot 2^{-2} + 0 \cdot 2^{-3} + 1 \cdot 2^{-4} + 1 \cdot 2^{-5} = \frac{1}{2} + \frac{1}{16} + \frac{1}{32} = \frac{19}{32}$$

Se observă că o deplasare a virgulei binare spre stânga cu n poziții corespunde unei împărțiri a numărului cu 2^n , iar o deplasare a virgulei binare spre dreapta cu n poziții corespunde unei înmulțiri a numărului cu 2^n .

Pentru a transforma un număr zecimal în corespondentul său binar, se procedează astfel: se divide în mod repetat numărul zecimal din stânga virgulei la 2, reținându-se restul. Acesta, scris în ordine inversă (de la dreapta spre stânga) este reprezentarea binară a părții întregi. Partea din dreapta virgulei se multiplică în mod repetat cu 2, înlăturând de fiecare dată partea zecimală și reținând partea întreagă. Scriind aceasta în ordine normală, (de la stânga la dreapta), se obține reprezentarea binară a părții fracționare.

Exemplul 5.1.

Să se transforme numărul zecimal 627,625 în format binar.

Soluție.

Partea întreagă		Partea zecimală	
627 : 2 = 313	1	0.625 x 2 = 1.250	1
313 : 2 = 156	1	0.250 x 2 = 0.500	0
156 : 2 = 78	0	0.500 x 2 = 1.000	1
78 : 2 = 39	0	0.000 x 2 = 0.000	0
39 : 2 = 19	1		
19 : 2 = 9	1		
9 : 2 = 4	1		
4 : 2 = 2	0		
2 : 2 = 1	0		
1 : 2 = 0	1		

Prin urmare $(627,625)_{10} = (1001110011,101)_2$

Operațiile cu numere binare se execută similar celor zecimale.

1. Adunarea
 - $0 + 0 = 0$
 - $0 + 1 = 1$
 - $1 + 0 = 1$
 - $1 + 1 = 0$ se transportă 1
2. Scăderea
 - $0 - 0 = 0$
 - $1 - 0 = 1$
 - $0 - 1 = 1$ se importă 1
 - $1 - 1 = 0$
3. Multiplicarea
 - $0 \times 0 = 0$
 - $1 \times 0 = 0$
 - $0 \times 1 = 0$
 - $1 \times 1 = 1$
4. Împărțirea
 - $1 : 1 = 1$
 - $0 : 1 = 0$
 împărțirea la 0 nu este definită.

Aritmetica în virgulă fixă este potrivită atât pentru operații cu numere întregi, cât și fracționare.

Dacă este necesară rotunjirea produsului a două numere, este mai bine a se limita reprezentarea în virgulă fixă a numerelor fracționare, decât a celor care au atât parte întreagă, cât și fracționară, deoarece reducerea numărului de biți ai părții întregi ar cauza erori mari.

În conversia semnalelor analogice bipolare, este necesar un bit adițional pentru a purta informația de semn. De obicei cel mai semnificativ bit este rezervat semnului numărului, cu convenția ca zero să indice un număr pozitiv, iar unu, un număr negativ. Rezultatul este un *cod bipolar*. Există mai multe posibilități de reprezentare a codurilor bipolare binare, alegerea dintre acestea făcându-se în funcție de avantajele și dezavantajele pe care le prezintă fiecare pentru aplicația respectivă. Patru metode sunt frecvent folosite pentru reprezentarea numerelor bipolare. În continuare se va considera că numerele sunt reprezentate pe $N=b+1$ biți, din care unul pentru semn.

Formatul mărime cu semn sau semn – valoare este cea mai simplă metodă pentru reprezentarea numerelor cu semn în format digital. Un zero în poziția MSB reprezintă un număr pozitiv, iar un unu în aceeași poziție

reprezintă un număr negativ. Restul de b biți reprezintă *modulul* sau *amplitudinea* numărului.

În cazul numerelor fracționare, reprezentarea mărime cu semn pentru un număr pozitiv $x \geq 0$ este de forma

$$(x)_{ms} = 0, b_1 b_2 \dots b_b, \quad (5.3)$$

iar pentru numărul negativ $x_N = -x = -0, b_1 b_2 \dots b_b$, de forma

$$(x_N)_{ms} = 1, b_1 b_2 \dots b_b, \quad (5.4)$$

Așa cum s-a precizat deja, virgula nu există fizic în reprezentarea numărului, dar, în cele ce urmează va fi utilizată pentru a specifica numerele fracționare. Se observă că în acest format zero are două reprezentări: $0,0\dots 0$ și $1,00\dots 0$.

Valoarea zecimală a unui număr fracționar pozitiv este

$$(x)_{ms} = \sum_{i=1}^b b_i 2^{-i}, \quad (5.5)$$

iar a unui număr fracționar negativ este

$$(x_N)_{ms} = -\sum_{i=1}^b b_i 2^{-i}. \quad (5.6)$$

Modulul unui număr fracționar reprezentat în formatul mărime cu semn este dat de

$$|x| = |x_N| = \sum_{i=1}^b b_i 2^{-i}. \quad (5.7)$$

Reprezentarea în *complement față de unu* este identică celei în reprezentarea mărime cu semn pentru numere pozitive, dar diferă prin modul cum sunt formate numerele negative. În acest format, un număr negativ este obținut prin complementarea numărului pozitiv corespunzător.

În cazul formatului fracționar, numerele pozitive se reprezintă ca în relația (5.3), iar cele negative $x_N = -x = -0, b_1 b_2 \dots b_b$ sub forma

$$(x_N)_{1C} = \overline{0, b_1 b_2 \dots b_b} = 1, \overline{b_1 b_2 \dots b_b} \quad (5.8)$$

Plecând de la relația (5.8), reprezentarea în complement față de unu a unui număr negativ fracționar mai poate fi exprimată în forma

$$(x_N)_{1C} = 1 \times 2^0 + \sum_{i=1}^b (1 - b_i) 2^{-i} = 2 - 2^{-b} - |x| \quad (5.9)$$

Se observă ambiguitate în reprezentarea lui zero, ca $0,0\dots 0$ sau $1,1\dots 1$.

Modulul numărului negativ $b_0, b_1 b_2 \dots b_b$ reprezentat în complement față de unu este

$$|x_N| = 1 - \sum_{i=1}^b b_i 2^{-i} - 2^{-b} \quad (5.10)$$

Valoarea zecimală a numărului negativ $b_0, b_1 b_2 \dots b_b$ reprezentat în complement față de unu este

$$((x_N)_{1C})_{10} = -1 + \sum_{i=1}^b b_i 2^{-i} + 2^{-b} \quad (5.11)$$

Spre exemplu, reprezentarea lui $-3/8$ este 1,100, care este complementul față de unu al lui 0,011 ($3/8$).

Reprezentarea în *complement față de doi* este identică cu formatul mărime cu semn în cazul numerelor pozitive. Prin urmare numerele pozitive sunt reprezentate cu un zero în poziția bitului de semn. Pentru a obține reprezentarea în complement față de doi a unui număr negativ, se scrie modulul acestuia în formatul mărime cu semn, se inversează biții acestei reprezentări și se adună o unitate logică în poziția LSB.

Similar, un număr fracționar pozitiv se reprezintă sub forma (5.3), iar numărul fracționar negativ $x_N = -x = -0, b_1 b_2 \dots b_b$, sub forma

$$(x_N)_{2c} = \bar{0}, \bar{b}_1 \bar{b}_2 \dots \bar{b}_b + 0,0 \dots 01 \quad (5.12)$$

Semnul “+” indică adunarea modulo 2 care ignoră bitul de transport, dacă acesta este prezent în MSB.

Plecând de la relația (5.12), reprezentarea în complement față de doi a unui număr fracționar negativ mai poate fi exprimată în forma

$$(x_N)_{2C} = 1 + \sum_{i=1}^b (1 - b_i) 2^{-i} + 2^{-b} = 2 - |x|, \quad (5.12')$$

adică, un număr fracționar negativ este complementul față de doi al numărului pozitiv corespunzător, care se obține scăzând numărul pozitiv din 2, reprezentat în binar. De aici provine denumirea formatului.

Din (5.9) și (5.12') rezultă

$$(x_N)_{2C} = (x_N)_{1C} + 2^{-b} \quad (5.13)$$

Valoarea zecimală a unui număr $b_0, b_1 b_2 \dots b_b$ reprezentat în complement față de doi, este

$$(x_{2C})_{10} = -b_0 2^0 + \sum_{i=1}^b b_i 2^{-i} \quad (5.14)$$

unde $b_0 = 0$, pentru numere pozitive și $b_0 = 1$, pentru numere negative.

Modului numărului negativ reprezentat în complement față de doi este

$$|x_N| = 1 - \sum_{i=1}^b b_i 2^{-i} \quad (5.15)$$

De exemplu, reprezentarea în complement față de doi a numărului $-3/8$ se obține din complementarea lui $0,011$ ($3/8$), rezultând $1,100$, și apoi adăugând $0,001$. Rezultatul final este $1,101$.

Codul binar deplasat sau offsetul binar este similar codului binar direct, obținându-se din acesta prin deplasarea în domeniul valorilor negative cu jumătate din întreaga scală. Cu $b+1$ biți se pot reprezenta 2^{b+1} numere. Pentru un cod bipolar există $2M$ numere, cu $M=2^b$, cuprinse în intervalul $-2^b \div (2^b - 1)$ pentru numere întregi și în intervalul $-1 \div (1-2^{-b})$ pentru numere fracționare. În acest format cel mai mic număr negativ este reprezentat de un număr format din $b+1$ biți de zero iar cel mai mare număr pozitiv este format din $b+1$ biți de unu. În acest caz zero are o singură reprezentare și, prin urmare, se evită ambiguitatea întâlnită la formatul mărime cu semn. Marele dezavantaj al acestei notații este dat de posibilele erori ce pot apărea la citirea MSB-ului, în loc de unu, zero sau invers, rezultând o eroare de amplitudine mare.

Dacă se compară formatul complement față de doi și offsetul binar, se constată că ele diferă prin MSB și, prin urmare, este ușor a se trece de la o reprezentare la alta.

În Tabelul 5.1 sunt date codurile bipolare prezentate pentru reprezentarea numerelor întregi pe 4 biți, dintre care unul pentru semn.

TABEL 5.1 Coduri bipolare

<i>Număr</i>	<i>Formatul mărime cu semn</i>	<i>Ofset binar</i>	<i>Complement față de doi</i>	<i>Complement față de unu</i>
7	0111	1111	0111	0111
6	0110	1110	0110	0110
5	0101	1101	0101	0101
4	0100	1100	0100	0100
3	0011	1011	0011	0011
2	0010	1010	0010	0010
1	0001	1001	0001	0001
0	0000	1000	0000	0000

0	1000	1000	0000	1111
-1	1001	0111	1111	1110
-2	1010	0110	1110	1101
-3	1011	0101	1101	1100
-4	1100	0100	1100	1011
-5	1101	0011	1011	1010
-6	1110	0010	1010	1001
-7	1111	0001	1001	1000
-8	-	0000	-	-

În Tabelul 5.2 sunt date, comparativ, diferite reprezentări ale numerelor fracționare pentru o lungime de 3 biți a cuvintelor.

Tabelul 5. 2

Număr binar	Echivalentul zecimal folosind reprezentarea		
	Mărime și semn	Complement față de 1	Complement față de 2
0,11	3 / 4	3 / 4	3 / 4
0,10	2 / 4	2 / 4	2 / 4
0,01	1 / 4	1 / 4	1 / 4
0,00	0	0	0
1,00	- 0	- 3 / 4	- 4 / 4 = - 1
1,01	-1 / 4	- 2 / 4	- 3 / 4
1,10	- 2 / 4	- 1 / 4	- 2 / 4
1,11	- 3 / 4	- 0	- 1 / 4

Din tabel se observă, așa cum s-a mai specificat, că există două reprezentări pentru zero în format mărime cu semn și complement față de 1 și nici o reprezentare pentru -1. Formatul complement față de 2 are o singură reprezentare pentru 0 și poate reprezenta numere cuprinse între -1 și $1 - 2^{-2}$ sau, în general, între -1 și $1 - 2^{-(N-1)}$ pentru un registru de N biți. Reprezentarea în *complement față de 2* este adesea utilizată în implementarea filtrelor digitale datorită ușurinței efectuării operațiilor de adunare și scădere, caz în care descăzutul se adună cu complementul față de doi a scăzătorului.

Diferența dintre numărul maxim și cel minim ce poate fi reprezentată se numește *domeniu dinamic*.

Exemplul 5.2.

Folosind reprezentarea în complement față de 2 pe 4 biți să se efectueze operațiile a) A - B și b) B - A unde A = 0,250 și B = 0,625

Soluție

a)	zecimal	complement față de 2
	0,250 -	0,010 +
	<u>0,625</u>	<u>1,011</u>
	-0,375	1,101 = - 0,375
b)	0,650 -	0,101 +
	<u>0,250</u>	<u>1,110</u>
	0,375	0,011 = 0,375

Se observă că în reprezentarea în complement față de 2 bitul de transport în poziția cea mai semnificativă este neglijat.

Adunarea și scăderea în complement față de 1 sunt similare, dar bitul de transport din poziția cea mai semnificativă este deplasat în poziția celui mai puțin semnificativ bit.

De exemplu, $\frac{4}{8} - \frac{3}{8} = \frac{1}{8}$. În formatul complement față de unu, transportul din MSB, dacă este prezent, este purtat spre LSB. Astfel, calculul $\frac{4}{8} - \frac{3}{8} = \frac{1}{8}$ devine $0,100 \oplus 1,100 = 0,000 \oplus 0,001 = 0,001$.

Adunarea și scăderea în sistemul mărime cu semn sunt mai complexe și, ca urmare, acesta este folosit mai mult la multiplicare, care se efectuează prin multiplicarea modulelor și stabilind semnul produsului.

Exemplul 5.3.

Să se multiplice numerele 0,625 și 0,250 folosind reprezentarea mărime cu semn.

Soluție.

Zecimal	Mărime cu semn
0,625	0,101
<u>0,250</u>	<u>0,010</u>
0000	000
3125	101
<u>1250</u>	<u>000</u>
0,156250	0,001010 = 0,156250

Multiplicarea în aritmetica complement față de 1 și față de 2 este mai dificilă și necesită un hard sau algoritmi speciali.

Dacă rezultatul unei operații aritmetice depășește numărul maxim ce poate fi reprezentat pe b biți, apare *depășirea*. În procesarea digitală se folosește, de obicei, formatul fracționar, numerele care reprezintă mărimile ce intervin în procesare și rezultatele operațiilor aritmetice sunt scalate, astfel încât modulul lor să nu depășească valoarea 1.

La multiplicarea numerelor fracționare, nu există probleme de depășire în cele trei aritmetici. Depășirea poate apărea numai când suma numerelor fracționare este mai mare decât 1. Dacă depășirea apare într-o etapă intermediară a adunării, în final nu va exista depășire, cu condiția ca valoarea absolută a rezultatului final să fie subunitară.

Exemplul 5.4.

Să se adune $0,3125 + 0,7500 + (-0,6250)$ folosind aritmetica în complement față de 1 pe cinci biți.

Soluție.

zecimal	complement față de 1	
0,3125	0,0101	
<u>+0,7500</u>	<u>0,1100</u>	
1,0625	1,0001	→ incorect, MSB = 1 implică număr negativ
<u>-0,6250</u>	<u>1,0101</u>	
0,4375	0,0111	→ ultimul 1 se datorează transportului

Exemplul 5.5.

Să se exprime fracțiile $\frac{7}{8}$ și $-\frac{7}{8}$ în formatele: mărime cu semn, complement față de 1 și complement față de 2.

Soluție. $x = \frac{7}{8}$, este reprezentat ca $2^{-1}+2^{-2}+2^{-3}$, care, în formatul mărime cu semn conduce la $x = 0,111$, iar $x = -\frac{7}{8}$ este reprezentat ca $x = 1,111$. Reprezentarea în complement față de unu și față de doi a lui $x = \frac{7}{8}$ este aceeași ca formatul mărime cu semn, adică $x = 0,111$. Reprezentarea în complement față de unu a lui $x = -\frac{7}{8}$ este $x_{1C} = 1,000$ și în complement față de doi este $x_{2C} = 1,000 + 0,001 = 1,001$.

Deși sunt posibile o mare varietate de alte reprezentări în virgulă fixă, cele descrise anterior sunt cele mai utilizate în practică. Cele mai multe procesoare de semnal în virgulă fixă folosesc aritmetica în complement față de doi. Aritmetica complementului față de doi este de fapt aritmetica modulo- 2^{b+1} (adică orice număr care depășește domeniul, este redus la acest domeniu, prin scăderea celui mai apropiat multiplu de 2^{b+1}).

La adunarea sau scăderea a două numere în virgulă fixă, fiecare de b biți lungime (cu un bit adițional de semn), rezultatul este un număr de b biți. Dacă rezultatul adunării depășește cel mai mare număr care poate fi reprezentat pe b biți, apare depășirea. Singura metodă pentru evitarea acestei probleme este creșterea numărului de biți din acumulator și, prin urmare, creșterea gamei dinamice care poate fi acoperită.

În general, înmulțirea a două numere în virgulă fixă, fiecare în lungime de b biți, are ca rezultat un produs de lungime $2b$ biți. În aritmetica cu virgulă fixă, produsul este de obicei trunchiat sau rotunjit la b biți, ceea ce conduce la o eroare de trunchiere sau rotunjire cauzată de eliminarea celor mai puțin semnificativi b biți.

Depășirea în cazul adunării numerelor în reprezentarea în aritmetica în virgulă fixă este un dezavantaj cauzat de domeniul dinamic redus. Aritmetica în virgulă mobilă nu prezintă acest dezavantaj.

5.2.2. Reprezentarea numerelor în virgulă mobilă

Reprezentarea în virgulă fixă a numerelor, permite acoperirea unui domeniu dinamic, $x_{\max}-x_{\min}$ cu o rezoluție

$$\Delta = \frac{x_{\max} - x_{\min}}{m - 1}, \quad (5.16)$$

unde $m=2^{b+1}$ este numărul de nivele, iar $b+1$ numărul de biți. O caracteristică de bază a reprezentării în virgulă fixă este că rezoluția este fixă. În plus, Δ crește direct proporțional cu creșterea domeniului dinamic.

Reprezentarea în virgulă mobilă poate fi folosită ca o metodă de acoperire a unui domeniu dinamic mai larg. Reprezentarea în virgulă mobilă cel mai des întâlnită în practică constă dintr-o mantisă M , care este partea fracționară a numărului și se încadrează în domeniul $1/2 \leq M < 1$, înmulțită cu factorul exponențial 2^E unde exponentul E este un întreg pozitiv sau negativ. Un număr X , este reprezentat ca: $X = M \cdot 2^E$.



Figura 5.2 Reprezentarea în virgulă mobilă

Mantisa și exponentul necesită fiecare câte un bit de semn pentru reprezentarea numerelor pozitive sau negative. Deoarece mantisa este o fracție cu semn, se poate folosi oricare din reprezentările în virgulă fixă descrise anterior.

De exemplu, numărul $X_1=5$ este reprezentat de următoarea mantisă și exponent:

$$M_1=0,101000$$

$$E_1=011$$

în timp ce numărul $X_2=\frac{3}{8}$ este reprezentat de următoarea mantisă și exponent:

$$M_2=0,110000$$

$$E_2=101$$

Dacă cele două numere se înmulțesc, mantisele sunt înmulțite și exponenții adunați. Prin urmare produsul celor două numere date mai sus este:

$$X_1 \cdot X_2 = M_1 \cdot M_2 \cdot 2^{E_1+E_2} = (0,011110) \cdot 2^{010} = (0,111100) \cdot 2^{001}$$

Împărțirea a două numere reprezentate în virgulă mobilă se efectuează prin împărțirea mantiselor și scăderea exponenților.

$$\frac{X_1}{X_2} = \frac{M_1}{M_2} \cdot 2^{(E_1-E_2)}$$

Adunarea a două numere în virgulă mobilă necesită ca exponenții să fie egali. Aceasta se poate obține deplasând virgula binară a mantisei celui mai mic număr spre stânga și compensând prin creșterea corespunzătoare a exponentului. Atunci numărul X_2 poate fi exprimat în forma

$$M_2=0,000011$$

$$E_2=011$$

Cu $E_1=E_2$, se pot aduna cele două numere X_1 și X_2 . Rezultatul este

$$X_1 + X_2 = (0,101011) \cdot 2^{011}$$

Se observă că operația de deplasare, impusă de egalarea exponenților lui X_2 și X_1 , poate conduce la o precizie mai mică în reprezentarea lui X_2 . În exemplul anterior, mantisa pe șase biți a fost suficient de lungă pentru a se face deplasarea a patru biți la dreapta pentru M_2 , fără a pierde nici unul. Totuși o deplasare a cinci biți va cauza pierderea unui singur bit iar deplasarea a șase biți va conduce la mantisa $M_2=0,000000$; de aceea aceasta va trebui rotunjită după deplasare astfel încât $M_2=0,000001$.

Eroarea de depășire apare la multiplicarea a două numere în virgulă mobilă când suma exponenților depășește domeniul dinamic al reprezentării în virgulă fixă a exponentului.

Comparând reprezentarea în virgulă fixă cu cea în virgulă mobilă, cu același număr total de biți, rezultă că reprezentarea în virgulă mobilă permite acoperirea unui domeniu mai larg prin varierea rezoluției în acel interval. Rezoluția scade odată cu creșterea mărimii numerelor succesive. Cu alte cuvinte, distanța succesivă dintre două numere reprezentate în virgulă mobilă crește odată cu creșterea numerelor în mărime. Astfel, pentru acoperirea aceluiași domeniu dinamic cu ambele reprezentări, în virgulă fixă și virgulă mobilă, reprezentarea în virgulă mobilă oferă rezoluție fină pentru numere mici, dar rezoluție slabă pentru numere mari, spre deosebire de reprezentarea în virgulă fixă, care oferă o rezoluție uniformă în reprezentarea numerelor.

De exemplu, pentru un calculator care lucrează pe 32 biți, este posibilă reprezentarea a 2^{32} numere. Dacă se dorește reprezentarea întregilor pozitivi începând cu zero, cel mai mare număr întreg ce poate fi reprezentat este: $2^{32}-1=4.294.967.295$. Distanța dintre două numere succesive (rezoluția) este 1. Altfel, se poate folosi bitul cel mai din stânga ca bit de semn și ceilalți 31 de biți rămași pentru valoare. Într-un astfel de caz reprezentarea în virgulă fixă permite acoperirea domeniului

$$-(2^{31}-1) = -2.147.483.647 \text{ la } (2^{31}-1) = 2.147.483.647$$

tot cu o rezoluție de 1. Dacă, însă, se alocă 10 biți pentru partea fracționară, 21 de biți pentru partea întreagă și un bit pentru semn, această reprezentare permite acoperirea domeniului dinamic:

$$-(2^{31}-1) \cdot 2^{-10} = -(2^{21}-2^{-10}) \text{ la } (2^{31}-1) \cdot 2^{-10} = 2^{21}-2^{-10} \text{ adică}$$

de la $-2.097.151,999$ la $2.097.151,999$

În acest caz, rezoluția este 2^{-10} . Prin urmare domeniul dinamic a fost scăzut cu un factor de aproximativ 1000 (2^{10} mai exact), în timp ce rezoluția a crescut cu același factor.

Pentru comparație, se presupune că cei 32 biți ai cuvântului sunt folosiți pentru a reprezenta numere în virgulă mobilă astfel: mantisa pe 23 de biți plus un bit de semn și exponentul cu 7 biți plus un bit de semn. Cel mai mic număr, în modul, va avea reprezentarea:

$$\begin{array}{cccc} \text{semn} & 23 \text{ biți} & \text{semn} & 7 \text{ biți} \\ 0, & 100\dots0 & 1 & 1111111 = \frac{1}{2} \times 2^{-127} \approx 0,3 \times 10^{-38} \end{array}$$

În cealaltă extremă, cel mai mare număr care poate fi reprezentat cu acest format în virgulă mobilă este:

$$\begin{array}{cccc} \text{semn} & 23 \text{ biți} & \text{semn} & 7 \text{ biți} \\ 0, & 11\dots1 & 0 & 1111111 = (1-2^{-23}) \times 2^{127} \approx 1,7 \times 10^{38} \end{array}$$

S-a obținut un domeniu dinamic de aproximativ 10^{76} , dar cu o rezoluție variabilă, adică rezoluție fină pentru numere mici și rezoluție slabă pentru numere mari.

5.2.3. Reprezentarea în virgulă mobilă pe bloc

Acest mod de reprezentare a numerelor este un hibrid între sistemele cu virgulă fixă și cele cu virgulă mobilă. În acest caz, în loc ca fiecare număr să fie reprezentat individual, ca în cazul sistemelor cu virgulă mobilă, un bloc sau un șir de numere are un exponent fix asociat. Acest exponent fix este obținut din examinarea tuturor numerelor din bloc și reprezentarea celui mai mare număr ca un număr cu virgulă mobilă cu o mantisă normalizată. Avantajul unui astfel de sistem constă în folosirea unui singur exponent pentru un bloc mare de numere. Astfel sistemul este potrivit pentru implementarea algoritmilor ce necesită un volum mare de calcule.

5.3. Efectele cuantizării în conversia A/D a semnalelor

Operațiile de bază îndeplinite de un convertor A/D sunt:

1. Să eșantioneze semnalul în mod periodic și cu rată de eșantionare suficient de mare pentru a evita eroarea alias;
2. Să cuantizeze amplitudinea eșantioanelor într-un set discret de nivele.

Prin urmare, dintr-un semnal analogic $x_a(t)$ eșantionat cu frecvența $F_s=1/T$, unde T este perioada de eșantionare, va rezulta o secvență $x[n]=x_a(nT)$, a cărei amplitudine este cuantizată, rezultând secvența

$$x_q[n] \equiv Q[x[n]] \quad (5.17)$$

unde $x_q[n]$ reprezintă semnalul cuantizat, iar $Q[\bullet]$ operația de cuantizare.

Dacă un semnal al cărui domeniu dinamic este R urmează a fi reprezentat pe $N=b+1$ biți, numărul nivelelor de cuantizare ce pot fi reprezentate este de 2^{b+1} . În reprezentarea în virgulă fixă b biți dau 2^b valori ale amplitudinii iar un bit dă informația de semn. Distanța dintre două nivele adiacente sau pasul de cuantizare este $\Delta = \frac{R}{2^{b+1}}$ [63].

În reprezentarea în virgulă fixă a numerelor fracționare, dacă domeniul dinamic depășește ± 1 , de multe ori este necesară scalarea semnalului, caz în care pasul de cuantizare al semnalului scalat este redus corespunzător la $\Delta_1 = \frac{2}{2^N} = 2^{-b}$.

Exemplul 5. 6.

Să se determine nivelele de cuantizare ale unui semnal continuu cu domeniul dinamic $\pm 20V$ după ce a fost eșantionat și apoi procesat cu un convertor A/D pe $N=4$ biți.

Soluție. Pasul de cuantizare pentru semnalul nescalat este $\Delta = \frac{40}{2^4} = 2,5V$. Pasul de cuantizare pentru semnalul scalat la domeniul ± 1 este $\Delta_1 = \frac{2}{2^4} = 0,125V$ care este $2^{-b} = 2^{-3}$, adică valoarea corespunzătoare unui 1 în poziția bitului cel mai puțin semnificativ.

5.3.1. Cuantizarea semnalului de intrare. Erori rezultate din rotunjire și trunchiere

În executarea calculelor folosind aritmetica în virgulă fixă sau mobilă, apare problema cuantizării numerelor prin trunchiere sau rotunjire de la o reprezentare pe un anumit număr de biți b_n (posibil a fi, la limită, și infinit în cazul unui eșantion al unui semnal analogic) la o alta, pe un număr mai mic de biți, b . Dacă valoarea semnalului se află între două nivele, aceasta poate fi aproximată fie prin cel mai apropiat nivel superior,

fie prin cel mai apropiat nivel inferior. Efectul cuantizării este că introduce o eroare a cărei valoare depinde de numărul de biți din numărul original și de numărul de biți de după cuantizare.

Sunt trei metode de cuantizare frecvent folosite:

- *Rotunjirea*, caz în care valoarea semnalului este aproximată de cel mai apropiat nivel de cuantizare.
- *Trunchierea*, caz în care valoarea semnalului este aproximată de cel mai mare nivel care este inferior sau egal valoric cu eşantionul semnalului.
- *Trunchierea semn – valoare*, care este asemănătoare cu trunchierea pentru numere pozitive, dar valorile negative ale semnalului sunt approximate de cel mai apropiat nivel de cuantizare mai mare sau egal cu semnalul.

Aceste descrieri se aplică cuantizării în aritmetica în virgulă fixă. Cele două metode de trunchiere rezultă din tratările diferite ale numerelor negative în reprezentările: mărime cu semn, complement față de 1, complement față de 2.

La un moment dat, nT , eroarea datorată cuantizării este

$$E_i = Q_i[x[n]] - x_a(nT) = x_{qi} - x_a \quad (5.18)$$

unde $i = r$ în cazul rotunjirii și $i = t$ în cazul trunchierii, $x_a = x_a(nT)$ reprezintă valoarea necuantizată a semnalului reprezentată pe $b_n + 1$ biți, iar $Q_i[x[n]] = x_{qi}$, valoarea cuantizată a semnalului reprezentată pe $b + 1$ biți.

Rotunjirea

În cazul rotunjirii

$$E_r = Q_r[x[n]] - x_a(nT) = x_{qr} - x_a \text{ și } -\frac{\Delta}{2} \leq E_r \leq \frac{\Delta}{2}, \Delta = 2^{-b} \quad (5.19)$$

Relația neliniară dintre x_{qr} și x_a este reprezentată în figura 5.3 unde x_a este un semnal cu amplitudine continuă ($b_n = \infty$).

În reprezentarea în virgulă fixă, eroarea de rotunjire satisface relația (5.19), indiferent de aritmetica folosită pentru reprezentarea numerelor negative, deoarece rotunjirea este independentă de semn, ea depinzând numai de mărimea numărului.

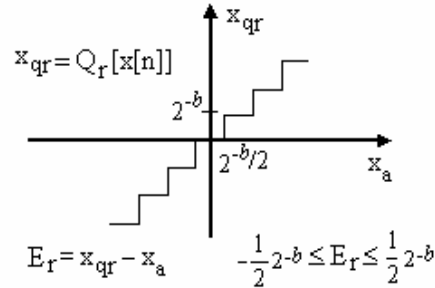


Figura 5.3 Relația dintre valorile cuantizate și necuantizate în cazul rotunjirii

În reprezentarea în virgulă mobilă, mantisa este cea trunchiată sau rotunjită.

Dacă
$$x_a = M_a \cdot 2^E \quad (5.20)$$

și
$$Q_r[x[n]] = M \cdot 2^E \quad (5.21)$$

atunci
$$E_r = Q_r[x[n]] - x_a = (M - M_a) 2^E \quad (5.22)$$

Dar pentru rotunjire
$$-\Delta/2 \leq M - M_a \leq \Delta/2 \quad (5.23)$$

și atunci din relația (5.19) rezultă

$$-2^E \Delta/2 \leq E_r \leq 2^E \Delta/2, \quad (5.24)$$

care dă eroarea absolută în virgulă mobilă datorată cuantizării mantisei.

Se definește eroarea relativă ε , astfel încât

$$Q_r[x[n]] = x_a(1 + \varepsilon) \quad (5.25)$$

Datorită rezoluției neuniforme, eroarea corespunzătoare reprezentării în virgulă mobilă este proporțională cu numărul, adică

$$E_r = \varepsilon \cdot x_a \quad (5.26)$$

și relația (5.24) devine

$$-2^E \Delta/2 \leq \varepsilon x_a \leq 2^E \Delta/2 \quad (5.27)$$

sau

$$-2^E \Delta/2 \leq \varepsilon M_a 2^E \leq 2^E \Delta/2 \quad (5.28)$$

adică

$$-\Delta/2 \leq \varepsilon M_a \leq \Delta/2 \quad (5.29)$$

Mantisa satisface relația

$$\frac{1}{2} \leq M_a < 1 \quad (5.30)$$

Dacă $M_a = \frac{1}{2}$, din (5.29) se obține domeniul maxim al erorii relative ca fiind

$$-\Delta \leq \varepsilon \leq \Delta \quad (5.31)$$

Trunchierea

Dacă metoda de cuantizare este trunchierea, numărul este aproximat în aritmetica în virgulă fixă, prin cel mai mare nivel care este mai mic sau egal cu valoarea semnalului. Trunchierea numerelor pozitive, negative și relația neliniară dintre x_{qt} și x_a sunt reprezentate în figura 5.4, unde x_a este un semnal cu amplitudine continuă.

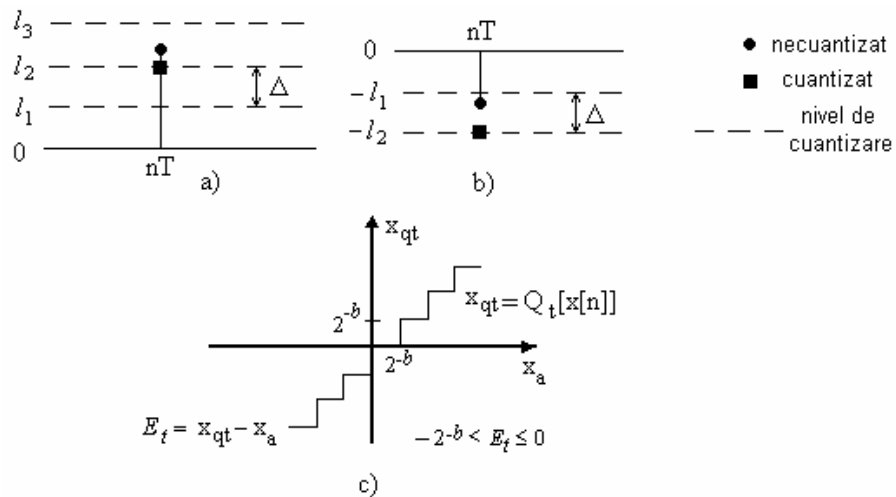


Figura 5.4. Relația dintre valorile cuantizate și necuantizate în cazul trunchierii a) pentru numere pozitive, b) pentru numere negative, c) caracteristica de trunchiere în complement față de 2

Eroarea de trunchiere $E_t = Q_t[x[n]] - x_a$ este negativă sau zero.

$$-\Delta < E_t \leq 0 \quad (5.32)$$

Acest lucru este valabil pentru toate numerele pozitive reprezentate în formatul mărime cu semn, complement față de 1 și complement față de 2.

În continuare se examinează trunchierea numerelor negative reprezentate în diverse formate. Fie întâi reprezentarea în *complement față de 2*. Se consideră că numărul ce urmează a fi trunchiat este reprezentat

pe $b_n + 1$ biți (la limită, se poate considera că $b_n = \infty$ pentru eșantioane ale unui semnal analogic). Modulul acestui număr negativ este

$$A_1 = 1 - \sum_{i=1}^{b_n} b_i \cdot 2^{-i} \quad (5.33)$$

Dacă acesta este trunchiat la b biți, modulul numărului devine

$$A = 1 - \sum_{i=1}^b b_i \cdot 2^{-i} \quad (5.34)$$

Diferența de mărime a modulului numărului negativ rezultată prin trunchiere este

$$A - A_1 = \sum_{i=1}^{b_n} b_i \cdot 2^{-i} - \sum_{i=1}^b b_i \cdot 2^{-i} = \sum_{i=b+1}^{b_n} b_i \cdot 2^{-i} \geq 0 \quad (5.35)$$

Deoarece modulul crește prin trunchiere, numărul negativ reprezentat în complement față de 2 devine mai mic. Valoarea maximă a modulului erorii se obține când toți coeficienții b_i sunt egali cu 1, caz în care

$$A - A_1 = 2^{-b} - 2^{-b_n} < \Delta, \quad (5.36)$$

deoarece $\Delta = 2^{-b}$. Prin urmare, în reprezentarea în complement față de 2, eroarea se situează în domeniul

$$-\Delta < E_i \leq 0 \quad (5.37)$$

Situația descrisă anterior este reprezentată în figura 5.4.

În cazul reprezentării numerelor negative în complement față de 1 pe $b_n + 1$ biți, modulul numărului negativ este

$$A_1 = 1 - \sum_{i=1}^{b_n} b_i 2^{-i} - 2^{-b_n} \quad (5.38)$$

Prin trunchierea la $b+1$ biți, modulul numărului negativ devine

$$A = 1 - \sum_{i=1}^b b_i 2^{-i} - 2^{-b}, \quad (5.39)$$

astfel încât diferența acestora este

$$\begin{aligned} A - A_1 &= \sum_{i=1}^{b_n} b_i 2^{-i} - \sum_{i=1}^b b_i 2^{-i} + 2^{-b_n} - 2^{-b} = \\ &= \sum_{i=b+1}^{b_n} b_i 2^{-i} - (2^{-b} - 2^{-b_n}) \leq 0 \end{aligned} \quad (5.40)$$

Modulul numerelor negative descrește prin trunchiere, adică, de fapt, acestea cresc. Situația este ilustrată în Figura 5.5. care reprezintă

trunchierea în reprezentarea semn - valoare. Prin urmare, domeniul în care poate lua valori eroarea ce apare prin trunchierea numerelor negative reprezentate în complement față de 1 este

$$0 \leq E_t < \Delta \quad (5.41)$$

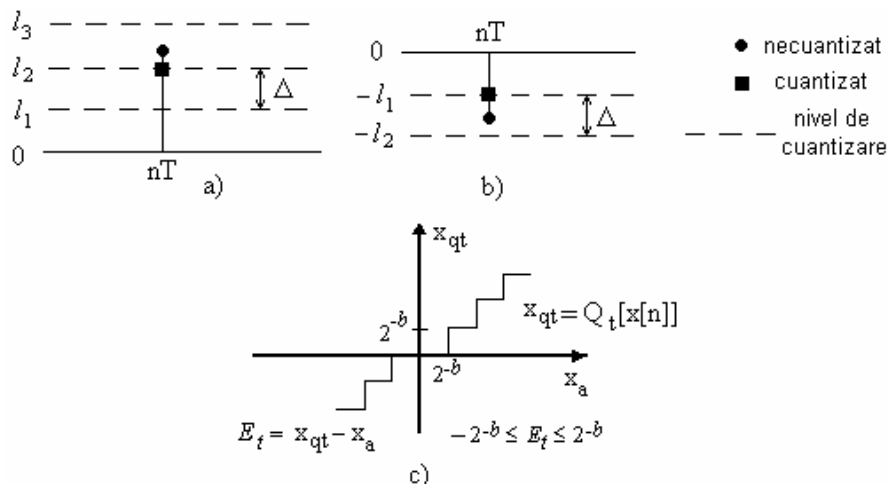


Figura 5.5. Relația dintre valorile cuantizate și necuantizate în cazul trunchierii semn valoare a) numere pozitive, b) numere negative, c) caracteristica de trunchiere în semn - valoare

În reprezentarea numerelor negative în formatul *mărime cu semn*, biții care reprezintă modulul numărului negativ sunt aceeași cu cei corespunzători numărului pozitiv, diferind numai bitul de semn. Aceasta înseamnă că prin trunchierea unui număr negativ modulul acestuia scade, iar valoarea trunchiată este dată de cel mai apropiat nivel de cuantizare care nu este mai mic decât numărul, situație reprezentată în Figura 5.5.

În continuare se va considera trunchierea mantisei în cazul reprezentării în *virgulă mobilă*.

$$E_t = Q_t[x[n]] - x_a = (M - M_a)2^E \quad (5.42)$$

În reprezentarea în *complement față de 2* a mantisei

$$-\Delta < M - M_a \leq 0 \quad (5.43)$$

sau
$$-2^E \Delta < E_t \leq 0 \quad (5.44)$$

Deoarece $E_t = \varepsilon x_a$, se obține

$$-2^E \Delta < \varepsilon x_a \leq 0 \quad (5.45)$$

$$\text{sau} \quad -2^E \Delta < \varepsilon M_a 2^E \leq 0 \quad (5.46)$$

$$\text{care implică} \quad -\Delta < \varepsilon M_a \leq 0 \quad (5.47)$$

Dacă $M_a = \frac{1}{2}$ se obține domeniul maxim al erorii relative ε , ca fiind

$$-2\Delta < \varepsilon \leq 0 \quad (5.48)$$

Dacă $M_a = -\frac{1}{2}$, domeniul erorii relative este

$$0 \leq \varepsilon < 2\Delta \quad (5.49)$$

În reprezentarea în *complement față de 1*, eroarea de trunchiere pentru valori pozitive ale mantisei este:

$$-\Delta < M - M_a \leq 0 \quad (5.50)$$

$$\text{sau} \quad -2^E \Delta < E_t \leq 0 \quad (5.51)$$

$$\text{Cu} \quad E_t = \varepsilon x_a = \varepsilon M_a 2^E \quad (5.52)$$

și $M_a = \frac{1}{2}$ se obține domeniul maxim al erorii relative pentru M_a pozitiv, ca fiind

$$-2\Delta < \varepsilon \leq 0 \quad (5.53)$$

Pentru valori negative ale mantisei, eroarea este

$$0 \leq M - M_a < \Delta \quad (5.54)$$

$$\text{sau} \quad 0 \leq E_t < 2^E \Delta \quad (5.55)$$

Pentru $M_a = -\frac{1}{2}$, domeniul maxim pentru eroarea relativă este

$$-2\Delta < \varepsilon \leq 0, \quad (5.56)$$

aceeași ca și pentru M_a pozitiv.

Acest lucru este valabil, de asemenea, și pentru cazul în care mantisa este reprezentată în formatul *mărime cu semn*.

5.3.2. Model statistic pentru cuantizarea fină

În calculele aritmetice ce implică cuantizare prin trunchiere sau rotunjire, este convenabil să se adopte o metodă statistică pentru caracterizarea erorilor rezultate. Cuantizorul poate fi modelat prin

introducerea unui zgomot aditiv $e[n]$ ce se suprapune peste semnalul $x[n]$, cu respectarea unor ipoteze ce vor fi specificate în cele ce urmează, adică

$$Q[x[n]] = x_q[n] = x[n] + e[n] \quad (5.57)$$

unde $e[n] = E_r$ pentru rotunjire și $e[n] = E_t$ pentru trunchiere, iar modelul este ilustrat în figura 5.6.

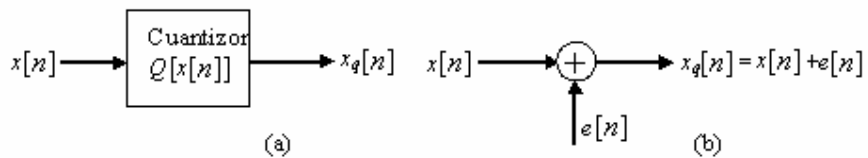


Figura 5.6. Modelul zgomotului aditiv pentru procesul liniar de cuantizare: (a) sistemul real; (b) model de cuantizare

Cum $x[n]$ poate fi orice număr care se încadrează în domeniul cuantizorului, eroarea de cuantizare este uzual modelată ca o variabilă aleatoare care se încadrează în limitele specificate anterior pentru erori. Mai mult, în practică, $b_n \gg b$, deci mărimea 2^{-b_n} poate fi neglijată în relațiile precedente. În aceste condiții, erorile de cuantizare ale numerelor reprezentate în virgulă fixă și virgulă mobilă se încadrează în intervalele prezentate în Tabelul 5.3.

Tabelul 5.3 Intervalele erorii de cuantizare

Tipul cuantizării	Tipul de aritmetică	Numere reprezentate cu virgulă fixă	Numere reprezentate cu virgulă mobilă
Rotunjire	-Semn-valoare -Complement față de 1 -Complement față de 2	$-2^{-b-1} \leq E_r \leq 2^{-b-1}$	$-2^{-b} \leq \varepsilon \leq 2^{-b}$
Trunchiere	Complement față de 2	$-2^{-b} < E_t \leq 0$	$-2^{-b+1} < \varepsilon \leq 0, x > 0$ $0 \leq \varepsilon < 2^{-b+1}, x < 0$
Trunchiere semn-valoare	-Complement față de 1 -Semn-valoare	$-2^{-b} < E_t \leq 0, x > 0$ $0 \leq E_t < 2^{-b}, x < 0$	$-2^{-b+1} < \varepsilon \leq 0$

În aceste condiții, funcțiile densitate de probabilitate pentru erorile de rotunjire și trunchiere pentru formatele de reprezentare în virgulă fixă prezentate sunt ilustrate în figura 5.7 [49]. Se observă că în cazul trunchierii în formatul complement față de doi, valoarea medie a erorii are o deplasare de $2^{-b}/2$, în timp ce pentru celelalte cazuri ilustrate anterior, eroarea are o valoare medie nulă.

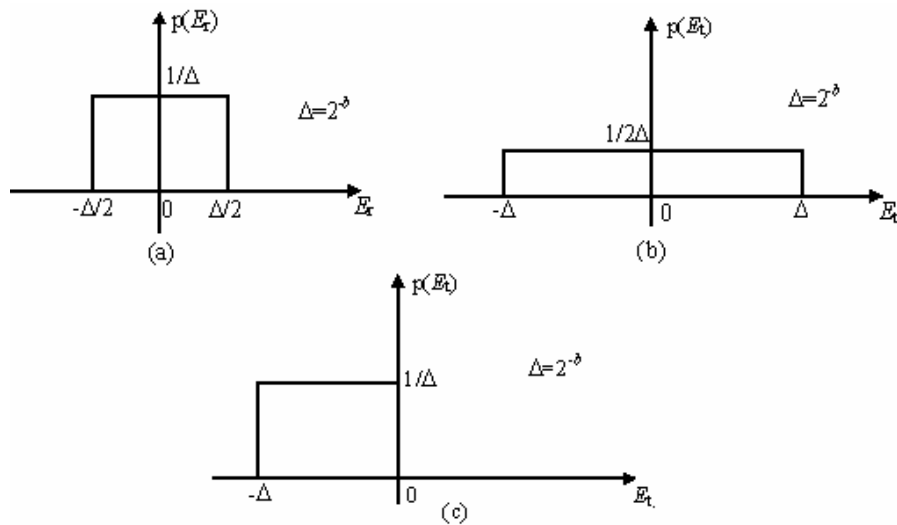


Figura 5.7 Caracterizarea statistică a erorilor de cuantizare. Funcțiile densitate de probabilitate ale (a) erorii de rotunjire; (b) erorii de trunchiere în formatul semn-valoare; (c) erorii de trunchiere în formatul complement față de doi

Analiza rezultatelor din Tabelul 5.3 și a expresiilor densităților de repartiție pentru erorile de rotunjire și trunchiere conduce la concluzia că rotunjirea este preferată altor metode de cuantizare, din următoarele motive[34]:

- semnalul de eroare este independent de tipul de aritmetică;
- media semnalului eroare este zero;
- nici o altă metodă de cuantizare nu conduce la o dispersie mai mică.

Cuantizarea reprezintă o operație neliniară și ireversibilă.

Efectele erorii de cuantizare datorate rotunjirii pot fi evidențiate dacă $e[n]$ se consideră o secvență aleatoare care satisface următoarele proprietăți:

1. Eroarea $e[n]$ este uniform distribuită în domeniul $[-\Delta/2, \Delta/2]$,
2. Secvența de eroare $\{e[n]\}$ este o secvență de zgomot alb staționar, pentru care $e[n]$ și $e[m]$, pentru $m \neq n$, sunt necorelate.
3. Secvența de eroare $\{e[n]\}$ este necorelată cu semnalul $x[n]$.

Ipotezele de mai sus sunt îndeplinite când pasul de cuantizare este mic și semnalul $x[n]$ traversează mai multe nivele de cuantizare între două eșantioane succesive. Efectul zgomotului aditiv, $e[n]$, asupra semnalului dorit poate fi studiat evaluând raportul semnal-zgomot (SNR) care, pe scară logaritmică (în decibeli), este

$$SNR = 10 \cdot \log_{10} \frac{P_x}{P_n} \quad (5.58)$$

unde P_x este puterea semnalului, iar P_n este puterea zgomotului de cuantizare.

Dacă eroarea de cuantizare este uniform distribuită în domeniul $(-\Delta/2, \Delta/2)$, așa cum este reprezentat în figura 5.7a, valoarea medie a erorii este zero și dispersia (puterea zgomotului de cuantizare) este

$$P_n = \sigma_e^2 = \int_{-\Delta/2}^{\Delta/2} e^2 p(e) de = \frac{1}{\Delta} \int_{-\Delta/2}^{\Delta/2} e^2 de = \frac{\Delta^2}{12} = \frac{2^{-2b}}{12} \quad (5.59)$$

Prin urmare, SNR este

$$SNR = 10 \cdot \log_{10} \frac{P_x}{P_n} = 10 \cdot \log_{10} P_x + 10 \cdot \log_{10} (12 \times 2^{2b}) \quad (5.60)$$

$$SNR = 10 \cdot \log_{10} P_x + 10,8 + 6b \quad (5.61)$$

Această expresie pentru SNR indică faptul că fiecare bit folosit în convertorul A/D sau cuantizor, mărește raportul semnal/zgomot de cuantizare cu 6 dB sau reduce puterea zgomotului de cuantizare cu 6 dB.

De exemplu, dacă se stabilește nivelul puterii zgomotului de cuantizare la -70 dB față de nivelul puterii semnalului, trebuie folosit un cuantizor pe 10 biți (sau convertor pe 10 biți).

Pentru a analiza efectul zgomotului de cuantizare asupra răspunsului unui sistem discret, liniar, invariant în timp, se consideră un astfel de sistem caracterizat de funcția pondere $h[n]$, la intrarea căruia se aplică semnalul cuantizat $x_q[n] = x[n] + e[n]$. Datorită liniarității sistemului, ieșirea sa este suma răspunsurilor sistemului la semnalul necuantizat $x[n]$ și la eroarea de cuantizare $e[n]$. Notând semnalul de

ieșire datorat zgomotului sau erorii de cuantizare cu $z[n]$, conform figurii 5.8, se poate scrie

$$z[n] = \sum_{k=0}^n h[k]e[n-k] \quad (5.62)$$

relație din care poate fi determinată dispersia zgomotului de ieșire cauzat de eroarea de cuantizare.

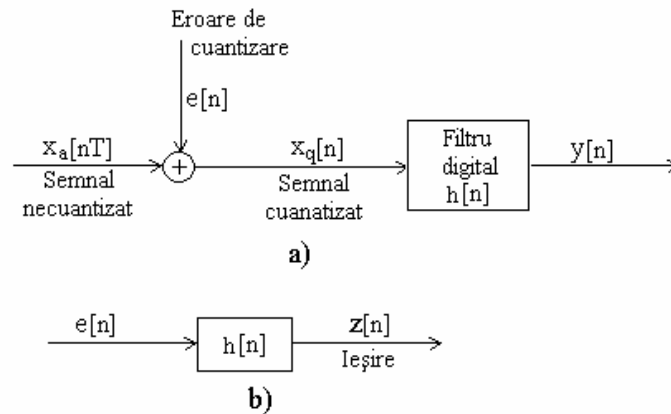


Figura 5.8. Model pentru eroarea datorată cuantizării semnalului de intrare
 a) Modelul de eroare, b) ieșirea datorată zgomotului de cuantizare

În cazul cuantizării prin rotunjire, ținând seama de ipotezele asumate pentru eroare și de relația (5.59), dispersia fiecărui termen din suma (5.62) este

$$\sigma_e^2 h^2[k] = \frac{\Delta^2}{12} h^2[k] \quad (5.63)$$

Deoarece dispersia unei sume de variabile aleatoare independente este egală cu suma dispersiilor lor, rezultă că, în ipoteza că erorile de cuantizare s-au presupus independente la diferite momente de timp, dispersia ieșirii $z[n]$ este

$$\sigma_{0z}^2[n] = \frac{\Delta^2}{12} \sum_{k=0}^n h^2[k] \quad (5.64)$$

Dispersia crește până la o valoare de regim permanent cu condiția ca filtrul să fie stabil. Dispersia de regim permanent se calculează cu relația

$$\sigma_{ozss}^2 = \lim_{n \rightarrow \infty} \sigma_{0z}^2[n] = \frac{\Delta^2}{12} \sum_{k=0}^{\infty} h^2[k] \quad (5.65)$$

O altă formă pentru expresia dispersiei de regim permanent a ieșirii poate fi obținută cu ajutorul funcției de sistem a filtrului, $H(z)$, în felul următor:

$$H(z) = \sum_{k=0}^{\infty} h[k] z^{-k} \quad (5.66)$$

$$H(z^{-1}) = \sum_{m=0}^{\infty} h[m] z^m \quad (5.67)$$

Prin urmare,

$$H(z)H(z^{-1}) = \sum_{k=0}^{\infty} \sum_{m=0}^{\infty} h[k]h[m] z^{m-k} \quad (5.68)$$

Multiplcând ambii membri cu z^{-1} și integrând după z pe un contur închis ce conține originea în planul z , rezultă

$$\int_c H(z)H(z^{-1})z^{-1} dz = \int_c \sum_{k=0}^{\infty} \sum_{m=0}^{\infty} h[k]h[m] z^{m-k-1} dz \quad (5.69)$$

Când conturul c este în regiunea de convergență pentru $H(z)$ și $H(z^{-1})$, se poate schimba ordinea de sumare și integrare din membrul drept. Se observă că cercul unitate este inclus în domeniul rezultat din intersecția regiunilor de convergență pentru $H(z)$ și $H(z^{-1})$, cu condiția ca $H(z)$ să fie stabil. Astfel se justifică alegerea cercului unitate drept contur de integrare. Relația (5.69) devine

$$\int_c H(z)H(z^{-1})z^{-1} dz = \sum_{k=0}^{\infty} \sum_{m=0}^{\infty} h[k]h[m] \int_c z^{m-k-1} dz \quad (5.70)$$

Deoarece conturul de integrare conține originea planului Z , conform teoremei lui Cauchy [48]

$$\int_c z^{m-k-1} dz = \begin{cases} 2\pi j & m = k \\ 0 & m \neq k \end{cases} \quad (5.71)$$

Cu (5.71), relația (5.70) devine

$$\int_c H(z)H(z^{-1})z^{-1} dz = 2\pi j \sum_{k=0}^{\infty} h^2[k] \quad (5.72)$$

și, deci,

$$\sum_{k=0}^{\infty} h^2[k] = \frac{1}{2\pi j} \int_c H(z)H(z^{-1})z^{-1} dz \quad (5.73)$$

Din (5.65) și (5.73) rezultă următoarea expresie pentru dispersia de regim permanent

$$\sigma_{ozss}^2 = \frac{\Delta^2}{12} \sum_{\substack{\text{polii din} \\ \text{cercul unitate}}} \text{reziduurile lui } H(z)H(z^{-1})z^{-1}, \quad (5.74)$$

expresie care, de multe ori, este mai ușor de evaluat decât (5.65).

Exemplul 5. 7.

Să se determine dispersia de regim permanent a zgomotului de la ieșirea unui sistem causal, stabil, de ordinul întâi, datorat cuantizării semnalului de intrare.

Soluție. Ecuația cu diferențe care caracterizează sistemul este $y[n] = Ay[n-1] + x[n]$, cu $|A| < 1$. Răspunsul la impuls al acestui sistem este $h[n] = A^n u[n]$. Din (5.65) rezultă dispersia zgomotului de ieșire

$$\sigma_{oz}^2[n] = \frac{\Delta^2}{12} \sum_{k=0}^n A^{2k} = \frac{\Delta^2}{12} \frac{1 - A^{2(n+1)}}{1 - A^2}$$

Dispersia de regim permanent, când $n \rightarrow \infty$, este $\sigma_{0zss}^2 = \frac{\Delta^2}{12(1 - A^2)}$.

$H(z) = \frac{1}{1 - Az^{-1}}$, cu un pol în $z = A$, și $H(z^{-1}) = \frac{1}{1 - Az}$ cu un pol în

$z = \frac{1}{A}$ în afara cercului unitate. Conform relației (5.74) rezultă

$$\sigma_{0zss}^2 = \frac{\Delta^2}{12} \sum \left(\text{reziduurile lui } \frac{z}{z-A} \cdot \frac{1}{1-Az} \cdot z^{-1} \Big|_{z=A} \right) = \frac{\Delta^2}{12(1 - A^2)}$$

identică, evident, cu expresia obținută anterior.

Pentru sisteme de ordin superior este mai ușor a se folosi relația (5.74) decât (5.65) din cauza complexității expresiei răspunsului la impuls.

5.4. Erori cauzate de cuantizarea coeficientilor filtrelor

5.4.1. Efectul cuantizării parametrilor filtrului asupra stabilității. Analiza sensibilității la cuantizarea coeficienților filtrelor IIR

Pentru a asigura stabilitatea unui filtru recursiv causal, toți polii acestuia trebuie să fie în interiorul cercului unitate din planul Z . În multe

cazuri este de dorit ca un pol sau o pereche de poli să fie în apropierea cercului unitate. Dacă în acest caz pasul de cuantizare este atât de mare încât reprezentarea polilor să fie pe sau în afara cercului unitate, filtrul astfel implementat devine instabil.

Fie, de exemplu, un filtru de ordinul întâi

$$y[n] = A y[n-1] + x[n] \quad (5.75)$$

și fie $N = b+1$, numărul biților disponibili reprezentării coeficientului A care, pentru un filtru stabil, este cuprins în domeniul $-1 < A < 1$.

Mărimea pasului de cuantizare este $\Delta = 2^{-b}$. Dacă $\varepsilon = 1 - A$ este distanța de la pol la cercul unitate, cea mai mică valoare a lui ε care poate fi precis reprezentată este $\Delta = 2^{-b}$. Pentru asigurarea stabilității trebuie ca pasul de cuantizare să fie mai mic sau egal cu distanța de la pol la cercul unitate, $\Delta \leq \varepsilon$, adică $2^{-N+1} \leq (1 - A)$, de unde rezultă

$$N \geq -\frac{\log_{10} \varepsilon}{\log_{10} 2} + 1 = -\frac{\log_{10}(1 - A)}{\log_{10} 2} + 1 \quad (5.76)$$

Exemplul 5. 8.

a) Fie $A = e^{-aT}$, unde $a = 1 \text{ rad/s}$, $T = 10^3$ secunde. Dacă se folosește trunchierea ca metodă de cuantizare, să se determine numărul minim de biți, N , necesar reprezentării lui A , astfel încât să nu rezulte instabilitate.

b) Dacă sunt disponibili 9 biți și $T = 10^3$ secunde, să se găsească a , astfel încât filtrul să fie stabil.

Soluție. a) $1 - A = 1 - e^{-aT} \approx aT$, prin urmare, $N \geq -\frac{\log_{10} aT}{\log_{10} 2} + 1 = 11$ biți

$$\text{b) } 9 = -\frac{\log_{10}(10^{-3} \cdot a)}{0.3} + 1 \text{ care necesită } a = 4 \text{ rad/secundă.}$$

Pentru filtrele de ordin superior localizarea polilor depinde, în general, de mai mulți coeficienți. Pentru a ilustra efectul cuantizării coeficienților asupra localizării polilor și, implicit, asupra caracteristicii de frecvență, fie un filtru IIR cu funcția de sistem

$$H(z) = \frac{\sum_{k=0}^M b_k z^{-k}}{1 + \sum_{k=1}^N a_k z^{-k}} \quad (5.77)$$

Filtrul IIR cu coeficienți cuantizați are funcția de sistem

$$\bar{H}(z) = \frac{\sum_{k=0}^M \bar{b}_k z^{-k}}{1 + \sum_{k=1}^N \bar{a}_k z^{-k}} \quad (5.78)$$

unde coeficienții cuantizați $\{\bar{b}_k\}$ și $\{\bar{a}_k\}$ pot fi exprimați în funcție de coeficienții necuantizați $\{b_k\}$ și $\{a_k\}$ prin relațiile

$$\begin{aligned} \bar{a}_k &= a_k + \Delta a_k & k = 1, 2, \dots, N \\ \bar{b}_k &= b_k + \Delta b_k & k = 0, 1, \dots, M \end{aligned} \quad (5.79)$$

$\{\Delta b_k\}$ și $\{\Delta a_k\}$ reprezentând erorile de cuantizare ale coeficienților.

Numitorul lui $H(z)$ poate fi exprimat în forma

$$D(z) = 1 + \sum_{k=0}^N a_k z^{-k} = \prod_{k=1}^N (1 - p_k z^{-1}) \quad (5.80)$$

unde $\{p_k\}$ sunt polii lui $H(z)$. Similar, se poate descompune numitorul lui $\bar{H}(z)$ în forma

$$\bar{D}(z) = \prod_{k=1}^N (1 - \bar{p}_k z^{-1}) \quad (5.81)$$

unde $\bar{p}_k = p_k + \Delta p_k$, $k=1, 2, \dots, N$, și Δp_k este eroarea sau perturbația care rezultă din cuantizarea coeficienților filtrului.

În continuare, se urmărește a se exprima perturbația totală Δp_i a polului p_i , în funcție de eroarea de cuantizare $\{\Delta a_k\}$ a coeficienților. Perturbația Δp_i poate fi exprimată ca [48]

$$\Delta p_i = \sum_{k=1}^N \frac{\partial p_i}{\partial a_k} \Delta a_k \quad (5.82)$$

unde $\frac{\partial p_i}{\partial a_k}$ reprezintă variația poziției polului p_i determinată de variația

coeficientului a_k . Astfel, eroarea totală este exprimată ca o sumă a erorilor datorate schimbărilor în fiecare din coeficienții $\{a_k\}$.

Derivatele parțiale $\partial p_i / \partial a_k$, $k=1, 2, \dots, N$, pot fi obținute diferențiind $D(z)$ în funcție de fiecare $\{a_k\}$, după cum urmează [48]:

$$\left(\frac{\partial D(z)}{\partial a_k} \right)_{z=p_i} = \left(\frac{\partial D(z)}{\partial z} \right)_{z=p_i} \left(\frac{\partial p_i}{\partial a_k} \right) \quad (5.83)$$

Din (5.83) rezultă

$$\frac{\partial p_i}{\partial a_k} = \frac{(\partial D(z) / \partial a_k)_{z=p_i}}{(\partial D(z) / \partial z)_{z=p_i}} \quad (5.84)$$

Numărătorul relației (5.84) este

$$\left(\frac{\partial D(z)}{\partial a_k} \right)_{z=p_i} = z^{-k} \Big|_{z=p_i} = p_i^{-k} \quad (5.85)$$

Numitorul relației (5.84) este

$$\begin{aligned} \left(\frac{\partial D(z)}{\partial z} \right)_{z=p_i} &= \left\{ \frac{\partial}{\partial z} \left[\prod_{l=1}^N (1 - p_l z^{-1}) \right] \right\}_{z=p_i} = \\ &= \left\{ \sum_{k=1}^N \frac{p_k}{z^2} \prod_{\substack{l=1 \\ l \neq k}}^N (1 - p_l z^{-1}) \right\}_{z=p_i} = \frac{1}{p_i} \prod_{\substack{l=1 \\ l \neq i}}^N (p_i - p_l) \end{aligned} \quad (5.86)$$

Prin urmare, relația (5.84) poate fi exprimată sub forma

$$\frac{\partial p_i}{\partial a_k} = \frac{p_i^{N-k}}{\prod_{\substack{l=1 \\ l \neq i}}^N (p_i - p_l)} \quad (5.87)$$

Înlocuind rezultatul din (5.87) în (5.82) rezultă eroarea totală de perturbație Δp_i în forma

$$\Delta p_i = \sum_{k=1}^N \frac{p_i^{N-k}}{\prod_{\substack{l=1 \\ l \neq i}}^N (p_i - p_l)} \Delta a_k \quad (5.88)$$

Această expresie oferă o măsură a sensibilității polului p_i la o schimbare a coeficienților $\{a_k\}$.

Un rezultat analog se poate obține pentru sensibilitatea zerourilor la erorile cauzate de cuantizarea parametrilor $\{b_k\}$.

Termenii $(p_i - p_l)$ din numitorul relației (5.88) reprezintă vectori, în planul Z , orientați de la polii $\{p_l\}$ la polul $\{p_i\}$. Dacă polii sunt foarte grupați, ca în cazul unui filtru de bandă îngustă reprezentat în figura 5.9, lungimile $|p_i - p_l|$ vor fi mici pentru polii din vecinătatea lui p_i . Aceste lungimi mici vor contribui la erori mari și va rezulta o perturbație Δp_i mare. Eroarea Δp_i poate fi minimizată prin maximizarea lungimii $|p_i - p_l|$.

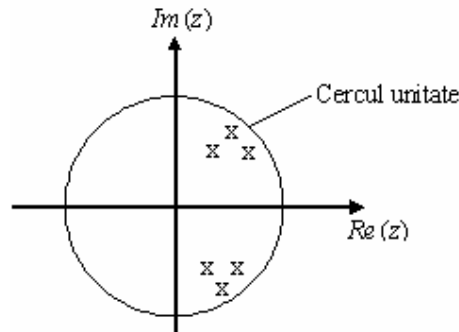


Figura 5.9 Poziții ale polilor unui filtru IIR trece bandă

Acest lucru se poate realiza prin implementarea filtrelor de ordin mare cu celule cu un singur pol sau cu doi poli. Filtrele cu un singur pol (și un singur zero) au valori complexe pentru coeficienți și necesită operații aritmetice în complex pentru realizarea lor. Această problemă poate fi evitată combinând polii și zerourile complex conjugate, pentru a forma secțiuni de filtru de ordin doi cu coeficienți reali. Deoarece polii complex conjugați sunt suficient de depărtați, eroarea de cuantizare Δp_i este minimizată și, în consecință, filtrul cu coeficienți cuantizați rezultă aproximativ mai bine caracteristica răspunsului în frecvență a filtrului cu coeficienți necuantizați.

Exemplul 5. 9.

Un filtru digital de ordinul doi are polii reali p_1 și p_2 . Acesta este implementat în forma directă. Se cere:

- a) Din relația generală (5.82) să se scrie o relație pentru modificarea poziției polilor datorată modificărilor coeficienților ecuației cu diferențe corespunzătoare.
- b) Dacă $p_1 = 0,98$ și $p_2 = 0,94$, care este numărul minim de biți necesar ca filtrul să rămână stabil în urma cuantizării coeficienților? Metoda de cuantizare se presupune a fi rotunjirea.

Soluție. a) Din (5.87) rezultă
$$\frac{\partial p_i}{\partial a_k} = \frac{p_i^{2-k}}{\prod_{\substack{l=1 \\ l \neq i}}^2 (p_i - p_l)}, k = 1, 2 \text{ și } i = 1, 2.$$

$$\frac{\partial p_1}{\partial a_1} = \frac{p_1}{p_1 - p_2} \quad \frac{\partial p_1}{\partial a_2} = \frac{1}{p_1 - p_2}$$

$$\frac{\partial p_2}{\partial a_1} = \frac{p_2}{p_2 - p_1} \quad \frac{\partial p_2}{\partial a_2} = \frac{1}{p_2 - p_1}$$

variația totală în poziția polilor este

$$\Delta p_i = \sum_{k=1}^2 \frac{\partial p_i}{\partial a_k} \Delta a_k$$

$$\text{adică } \Delta p_1 = \frac{\partial p_1}{\partial a_1} \Delta a_1 + \frac{\partial p_1}{\partial a_2} \Delta a_2 = \frac{1}{p_1 - p_2} [p_1 \Delta a_1 + \Delta a_2]$$

$$\text{și } \Delta p_2 = \frac{\partial p_2}{\partial a_1} \Delta a_1 + \frac{\partial p_2}{\partial a_2} \Delta a_2 = \frac{1}{p_2 - p_1} [p_2 \Delta a_1 + \Delta a_2]$$

b) Este necesar a determina Δa_1 și Δa_2 . Numitorul funcției de transfer a filtrului are forma $(z - p_1)(z - p_2) = z^2 - a_1 z + a_2$ unde $a_1 = p_1 + p_2$ și $a_2 = p_1 p_2$. Pentru asigurarea stabilității trebuie ca $-2 < a_1 < 2$ și $-1 < a_2 < 1$ [63]. În aritmetica în virgulă fixă coeficientul a_1 poate fi scalat pentru a se obține un număr fracționar, deși pentru coeficienții filtrului virgula binară este adesea mutată spre dreapta pentru a adapta coeficienții la mărimi mai mari ca unitatea. În orice caz se poate calcula pasul de cuantizare și numărul de biți, $N = b + 1$.

$$\text{Pentru } a_1, \Delta = \frac{4}{2^N} \text{ și pentru rotunjire } \Delta a_1 = \frac{\Delta}{2} = \frac{2}{2^N}$$

S-ar putea alege același pas de cuantizare și pentru a_2 , caz în care ar fi necesari $N - 1$ biți deoarece domeniul lui a_2 este jumătate din cel pentru a_1 . În schimb, s-ar putea adopta N biți pentru ambele registre, pentru a_1 și a_2 și pasul de cuantizare pentru a_2 să fie $\frac{2}{2^N} = \frac{\Delta}{2}$, astfel

$$\text{încât, pentru rotunjire } \Delta a_2 = \Delta/4 = \frac{1}{2^N}.$$

Pentru ultima alegere, din expresia menționată anterior pentru schimbarea poziției polului rezultă

$$\Delta p_1 = \frac{1}{0,98 - 0,94} [(0,98)2 + 1,0]/2^N = 74/2^N \text{ și}$$

$$\Delta p_2 = \frac{1}{0,94 - 0,98} [(0,94)2 + 1,0]/2^N = -72/2^N$$

Polul p_1 , fiind mai apropiat de cercul unitate este posibil să cauzeze instabilitatea filtrului, dacă nu este reprezentat adecvat. Pentru stabilitate, trebuie să fie îndeplinită relația $1 - p_1 = 0,02 > \Delta p_1 = 74/2^N$ sau $2^N > 3700$, care implică $N=12$ biți lungimea minimă a registrului.

Pentru a completa analiza, este necesar a considera și cazul polilor complex conjugați în expresia funcției de transfer (5.77). Numitorul acesteia se poate scrie

$$1 + \sum_{k=1}^N a_k z^{-k} = \prod_{i=1}^q (1 - p_i z^{-1}) \prod_{k=1}^s [1 - 2r_k (\cos \theta_k) z^{-1} + r_k^2 z^{-2}] \quad (5.89)$$

unde $s = \frac{N-q}{2}$, cu q poli simpli și s perechi de poli complex conjugați.

Diferențiind (5.89) în raport cu a_l , cu $1 \leq l \leq N$ se determină sensibilitatea la cuantizarea coeficienților $\frac{\partial p_m}{\partial a_l}$, $1 \leq m \leq q$ și $\frac{\partial r_g}{\partial a_l}$,

$\frac{\partial \theta_g}{\partial a_l}$, $1 \leq g \leq s$. După câteva prelucrări matematice rezultă pentru polii

simpli p_m , $1 \leq m \leq q$ [58]

$$\frac{\partial p_m}{\partial a_l} = \frac{p_m^{-l+1}}{\prod_{\substack{i=1 \\ i \neq m}}^q (1 - p_i p_m^{-1}) \prod_{k=1}^s [1 - 2r_k (\cos \theta_k) p_m^{-1} + r_k^2 p_m^{-2}]}, \quad (5.90)$$

și pentru polii complecși $r_g e^{\pm \theta_g}$, $1 \leq g \leq s$

$$\frac{\partial r_g}{\partial a_l} = \frac{-r_g^{-l+1} \sin[(l-1)\theta_g]}{2C_g \sin \theta_g} \quad (5.91)$$

$$\frac{\partial \theta_g}{\partial a_l} = \frac{r_g^{-l} \{ \sin[(l-2)\theta_g] - \cos \theta_g \sin[(l-1)\theta_g] \}}{2C_g \sin^2 \theta_g} \quad (5.92)$$

unde

$$C_g = \prod_{i=1}^q (1 - p_i z^{-1}) \prod_{\substack{k=1 \\ k \neq g}}^s (1 - 2r_k \cos \theta_k z^{-1} + r_k^2 z^{-2}) \Big|_{z=r_g e^{j\theta_g}} \quad (5.93)$$

Deviațiile totale sunt

$$\Delta p_i = \sum_{l=1}^N \frac{\partial p_m}{\partial a_l} \Delta a_l \quad l = 1, \dots, q \quad (5.94)$$

$$\Delta r_g = \sum_{l=1}^N \frac{\partial r_g}{\partial a_l} \Delta a_l \quad g = 1, \dots, s \quad (5.95)$$

$$\Delta \theta_g = \sum_{l=1}^N \frac{\partial \theta_g}{\partial a_l} \Delta a_l \quad g = 1, \dots, s \quad (5.96)$$

Din nou se observă că, dacă polii sunt grupați, ca în cazul filtrelor de bandă îngustă, polii realizării în forma directă sunt sensibili la erorile de cuantizare a coeficienților și, cu cât este mai mare numărul de poli grupați, cu atât și sensibilitatea este mai mare.

Este interesant de observat modul în care influențează structura de implementare a filtrului erorile cauzate de cuantizarea coeficienților. Pentru a ilustra acest lucru, fie un filtru cu doi poli complex conjugați, caracterizat de funcția de sistem

$$H(z) = \frac{1}{1 - (2r \cos \theta)z^{-1} + r^2 z^{-2}} \quad (5.97)$$

Filtrul are polii la $z_{1,2} = re^{\pm j\theta}$. Când este realizat ca în figura 5.10, există doi coeficienți: $a_1 = -2r \cos \theta$ și $a_2 = r^2$. Cu precizie infinită este posibil să obținem un număr infinit de poziții ale polilor. Evident, cu precizie finită (adică a_1 și a_2 cuantizați), pozițiile posibile ale polilor sunt în număr finit.

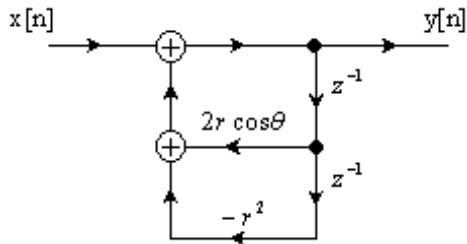


Figura 5.10. Realizare directă a unui filtru cu doi poli

De exemplu, pentru $b=3$, sunt posibile 7 valori nenule pentru a_1 și a_2 . În figura 5.11 sunt reprezentate pozițiile posibile ale polilor, numai pentru primul cadran al planului z . Sunt posibile 40 de poziții ale polilor în acest caz. Neuniformitatea în poziția polilor este datorată faptului că se cuantizează r^2 iar polii se găsesc pe un arc de cerc de rază r . Pentru o anumită cuantizare a coeficienților, polii se află pe o grilă din planul z

definită de intersecția cercurilor concentrice corespunzătoare cuantizării lui r^2 și liniilor verticale corespunzătoare cuantizării lui $2r\cos\theta$. De importanță particulară este setul rar de poli, pentru θ apropiat de zero și, datorită simetriei, pentru θ în apropierea lui π . Această situație va fi critic nefavorabilă pentru filtrele trece jos și filtrele trece sus care au în mod normal polii grupați în jurul frecvenței unghiulare $\theta=0$ și, respectiv, $\theta=\pi$.

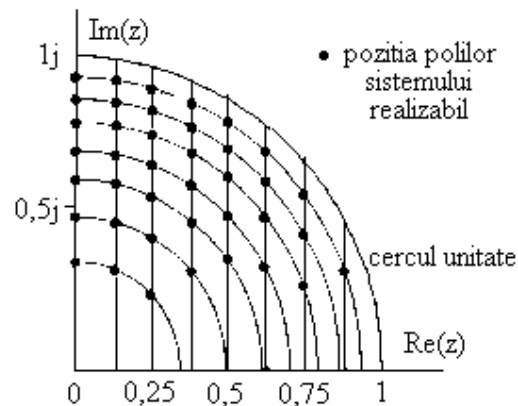


Fig. 5.11 Poziții posibile ale polilor structurii de ordinul doi în planul Z , pentru cuantizarea pe trei biți

O alternativă în realizarea filtrelor cu doi poli este forma cuplată, reprezentată în figura 5.12.

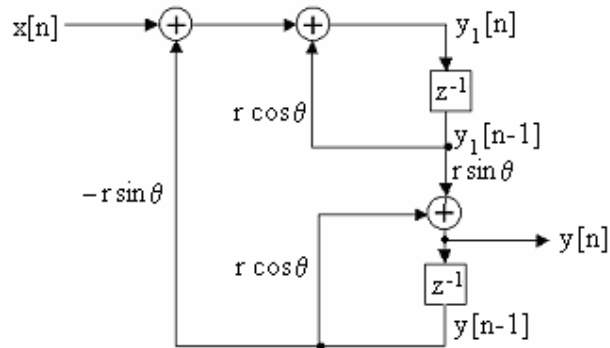


Figura 5.12. Realizare în forma cuplată a filtrului IIR cu doi poli

Cele două ecuații cuplate sunt:

$$\begin{aligned} y_1[n] &= x[n] + (r \cos \theta) \cdot y_1[n-1] - (r \sin \theta) \cdot y[n-1] \\ y[n] &= (r \sin \theta) \cdot y_1[n-1] + (r \cos \theta) \cdot y[n-1] \end{aligned} \quad (5.98)$$

Transformând aceste ecuații în domeniul Z, se poate scrie

$$\frac{Y(z)}{X(z)} = H(z) = \frac{(r \sin \theta)z^{-1}}{1 - (2r \cos \theta)z^{-1} + r^2 z^{-2}} \quad (5.99)$$

În forma cuplată se observă că sunt de asemenea doi coeficienți, $\alpha_1 = r \sin \theta$ și $\alpha_2 = r \cos \theta$. Deoarece ambii sunt liniari în r , pozițiile posibile ale polilor sunt acum puncte egal spațiate pe un caroiaj dreptunghiular, ca în figura 5.13.

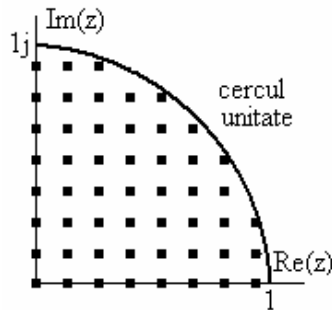


Figura 5.13. Poziții posibile ale polilor filtrului cu doi poli, realizat în forma cuplată din figura 5.12

Ca urmare, poziționarea polilor este acum uniform distribuită în interiorul cercului, lucru mult mai favorabil decât realizarea precedentă, mai ales pentru filtrele trece jos. Prețul plătit pentru această distribuire uniformă a poziției polilor este o creștere a volumului de calcule. Realizarea în formă cuplată necesită patru multiplicări, câte două pentru fiecare ieșire, în timp ce realizarea din figura 5.10 necesită doar două multiplicări. Este interesant de observat faptul că pentru o anumită lungime a coeficienților, forma directă permite o plasare mai adecvată a polilor cu r apropiat de unitate și θ mare, pe când forma cuplată este mai avantajoasă pentru θ mic.

Deoarece sunt diverse metode de a realiza secțiunile de ordin doi ale filtrelor, este, de asemenea, clar că sunt multe posibilități pentru localizarea polilor în cazul coeficienților cuantizați. Ideal ar fi să se selecteze o structură care conduce la un set dens de puncte în regiunea unde se află polii. Din nefericire nu există o metodă simplă și sistematică pentru determinarea realizării filtrului care să ducă la rezultatul dorit.

Având dat un filtru IIR de ordin înalt care trebuie implementat ca o combinație de secțiuni de ordinul doi, va trebui să se decidă între o structură în cascadă și una în paralel, adică între realizarea

$$H(z) = \prod_{k=1}^K \frac{b_{k0} + b_{k1}z^{-1} + b_{k2}z^{-2}}{1 + a_{k1}z^{-1} + a_{k2}z^{-2}} \quad (5.100)$$

și realizarea

$$H(z) = \sum_{k=1}^K \frac{c_{k0} + c_{k1}z^{-1}}{1 + a_{k1}z^{-1} + a_{k2}z^{-2}} \quad (5.101)$$

Dacă filtrul IIR are zerouri pe cercul unitate, cum este cazul filtrelor eliptice și Cebyshev de tipul doi, fiecare secțiune de ordin doi din configurația în cascadă din (5.100) conține o pereche de zerouri complex conjugate. Coeficienții $\{b_{ki}\}$ din (5.100) determină în mod direct pozițiile acestor zerouri, iar cuantizarea lor tinde să le deplaseze de pe cercul unitate. Sensitivitatea răspunsului sistemului la eroarea de cuantizare este ușor și direct controlabilă prin alocarea unui număr suficient de biți pentru reprezentarea coeficienților cuantizați $\{b_{ki}\}$ cu o precizie specificată. Astfel va exista control direct asupra polilor și zerourilor care rezultă din procesul de cuantizare. De fapt, se poate evalua efectul perturbării rezultate din cuantizarea coeficienților $\{b_{ki}\}$, cu o anumită precizie cerută.

Realizarea în paralel a lui $H(z)$, conform relației (5.101), asigură un control direct doar asupra polilor sistemului. Coeficienții numărătorului $\{c_{k0}\}$ și $\{c_{k1}\}$ sunt obținuți prin descompunerea în fracții simple a lui $H(z)$. Prin urmare polii influențează indirect localizarea zerourilor, prin combinarea tuturor termenilor din descompunerea în fracții simple a lui $H(z)$ și, în consecință, este mult mai dificil a se determina efectul erorii de cuantizare datorat coeficienților $\{c_{ki}\}$, în localizarea zerourilor sistemelor.

Cuantizarea parametrilor $\{c_{ki}\}$ poate produce o perturbație semnificativă a pozițiilor zerourilor și, de obicei, va fi suficient de mare în implementările cu virgulă fixă pentru a deplasa zerourile de pe cercul unitate. Aceasta este o situație foarte neplăcută, care poate fi însă remediată folosind o reprezentare în virgulă mobilă. În orice caz, structura în cascadă este mult mai robustă în prezența cuantizării coeficienților și trebuie să fie alegerea preferată în aplicații practice, mai ales unde este folosită reprezentarea în virgulă fixă.

5.4.2. Cuantizarea coeficienților filtrelor FIR

Așa cum s-a arătat și în secțiunea precedentă, analiza sensibilității aplicată polilor unui sistem se aplică direct și zerourilor filtrelor IIR. Prin urmare, o expresie asemănătoare cu relația (5.88) se poate obține pentru zerourile unui filtru FIR. Pentru a minimiza sensibilitatea la cuantizarea coeficienților, va trebui ca filtrul FIR cu un număr mare de zerouri să fie implementat ca o cascadă de secțiuni de ordinul unu și doi.

Un aspect important în practică îl reprezintă filtrele FIR cu răspuns liniar de fază. Realizările directe ale unor astfel de filtre mențin proprietatea de fază liniară chiar și în cazul cuantizării coeficienților. Aceasta rezultă din observația că funcția de sistem a unui filtru FIR de fază liniară satisface proprietatea

$$H(z) = \pm z^{-(M-1)} H(z^{-1}), \quad (5.102)$$

indiferent dacă coeficienții sunt sau nu, cunțizați.

Prin urmare, cuantizarea coeficienților filtrului FIR afectează doar caracteristica de amplitudine.

Din practică se știe că pentru a reprezenta coeficienții unui filtru FIR de fază liniară de lungime moderată ($M=32 \div 256$) sunt necesari cel puțin 10 biți, dar, dacă este posibil, se preferă a se folosi 12 până la 14 biți. Cu creșterea lungimii filtrului trebuie să crească și numărul de biți pentru reprezentarea coeficienților, pentru a menține aceeași eroare în răspunsul în frecvență al filtrului. Se presupune, de exemplu, că fiecare coeficient al filtrului este rotunjit la $(b+1)$ biți. Prin urmare, eroarea de rotunjire se încadrează în domeniul: $-2^{-b}/2 < e_r[n] < 2^{-b}/2$.

Valoarea cuantizată a răspunsului la impuls poate fi reprezentată ca $h_q[n] = h[n] + e_r[n]$ și eroarea în răspunsul în frecvență este

$$E_M(\omega) = \sum_{n=0}^{M-1} e_r[n] \cdot e^{-j\omega n} \quad (5.103)$$

Presupunând că $e_r[n]$ este o variabilă aleatoare uniform distribuită în intervalul $[-2^{-b}/2, 2^{-b}/2]$ cu valoarea medie zero, $E_M(\omega)$ va fi, de asemenea, de medie zero. Presupunând, în continuare, că $e_r[n]$ poate fi modelată ca o secvență de zgomot alb staționar, secvența erorilor $e_r[n]$, $0 \leq n \leq M-1$, are eșantioanele necorelate. Prin urmare, dispersia erorii în răspunsul în frecvență $E_M(\omega)$ este suma dispersiilor celor M termeni $e_r[n]$

$$\sigma_E^2 = \frac{2^{-2b}}{12} M \quad (5.104)$$

Ecuția (5.104) subliniază faptul că dispersia erorii crește liniar cu lungimea filtrului M . Deviația standard a erorii $E_M(\omega)$ este

$$\sigma_E = \frac{2^{-b}}{\sqrt{12}} \sqrt{M} \quad (5.105)$$

Prin urmare, pentru fiecare creștere de patru ori a lui M , precizia în reprezentarea coeficienților filtrului trebuie crescută cu un bit, pentru a menține deviația standard fixă. Din practică se constată că pentru a avea o deviație standard acceptabilă se folosesc 12, 13 biți. Dacă lungimea filtrului, M , este mai mare decât 256 sau numărul de biți folosiți pentru reprezentarea coeficienților este mai mic de 12, atunci filtrul trebuie implementat ca o cascadă de secțiuni de filtre de lungimi mai mici.

Într-o realizare în cascadă, de forma

$$H(z) = G \cdot \prod_{k=1}^K H_k(z) \quad (5.106)$$

secțiunile de ordinul doi sunt:

$$H_k(z) = 1 + b_{k1}z^{-1} + b_{k2}z^{-2}. \quad (5.107)$$

Coeficienții au forma $b_{k1} = -2r_k \cos\theta_k$ și $b_{k2} = r_k^2$. Cuantizarea lui b_{k1} și b_{k2} conduce la localizarea zerourilor ca în figura 5.11, cu excepția faptului că grid-ul se extinde în afara cercului unitate.

Ecuția (5.102) arată că zerourile lui $H(z^{-1})$ sunt identice cu cele ale lui $H(z)$. Dacă $H(z)$ are un zero complex $z = r_k \cdot e^{j\theta_k}$ atunci $H(z)$ trebuie să aibă și o “image–oglină” a acestuia, adică zeroul $z^{-1} = (1/r_k) \cdot e^{-j\theta_k}$. Pe de altă parte, dacă răspunsul la impuls este real, zerourile complexe ale lui $H(z)$ apar în perechi conjugate. Problema care apare în acest caz este menținerea proprietății de fază liniară, deoarece perechea de zerouri cuantizate $z_{3,4} = (1/r_k) \cdot e^{\pm j\theta_k}$ poate să nu fie imaginea în oglindă a perechii de zerouri cuantizate $z_{1,2} = r_k \cdot e^{\pm j\theta_k}$.

Această problemă poate fi evitată prin rearanjarea termenilor corespunzători imaginii în oglindă. Se pot scrie astfel coeficienții imaginii în oglindă, sub forma

$$\left(1 - \frac{2}{r_k} \cos\theta_k z^{-1} + \frac{1}{r_k^2} z^{-2}\right) = \frac{1}{r_k^2} \left(r_k^2 - 2r_k \cos\theta_k z^{-1} + z^{-2}\right) \quad (5.108)$$

Factorul $\{1/r_k^2\}$ poate fi combinat cu câștigul total G , sau poate fi distribuit în secțiunile de filtru de ordin doi. Termenul din (5.108) conține exact aceeași parametri ca și factorul $(1 - 2r_k \cos \theta_k z^{-1} + r_k^2 z^{-2})$ și, prin urmare, zerourile apar acum în perechi imagine-oglină chiar dacă coeficienții sunt cuantizați.

5.5. Erori cauzate de cuantizarea produselor. Caracterizarea statistică a efectelor cuantizării în realizarea în virgulă fixă a filtrelor digitale

Multiplicarea a două numere reprezentate pe b biți fiecare, exceptând bitul de semn, are ca rezultat un număr reprezentat pe $2b$ biți. În practică, datorită lungimii finite a registrelor cu care se lucrează, se impune exprimarea produselor prin b biți semnificativi, astfel încât, inevitabil, cuantizarea este asociată cu formarea produsului. Indiferent de tipul de cuantizare folosit, s-a încetățenit ca acesta să se numească *rotunjirea produsului*. Efectul acestei cuantizări asupra performanțelor filtrului depinde de modul de implementare a acestuia.

Se presupune că eroarea de rotunjire asociată formării produsului este independentă de la o iterație la alta, astfel încât poate fi folosit modelul cuantizării fine, sursele de zgomot fiind introduse în sistem după multiplicatoare. Astfel, multiplicatorul este modelat cu o operație în precizie infinită urmată de o sursă de zgomot aditiv $e[n]$, așa încât rezultatul final să fie egal cu un nivel de cuantizare, exact cum s-a procedat la caracterizarea erorii de cuantizare la conversia A/D a unui semnal analogic.

Se începe cu caracterizarea zgomotului de rotunjire într-un filtru cauzal, cu un singur pol, care este implementat în aritmetica cu virgulă fixă și este descris de ecuația neliniară cu diferențe

$$v[n] = Q_r[av[n-1]] + x[n] \quad (5.109)$$

Efectul rotunjirii produsului $av[n-1]$ este modelat cu o secvență de zgomot $e[n]$ adunată la produsul necuantizat $av[n-1]$, care este

$$Q_r[av[n-1]] = av[n-1] + e[n] \quad (5.110)$$

Cu acest model pentru eroarea de cuantizare, sistemul considerat este descris de ecuația liniară cu diferențe

$$v[n] = av[n-1] + x[n] + e[n] \quad (5.111)$$

Sistemul corespunzător este ilustrat în diagrama bloc din figura 5.14.

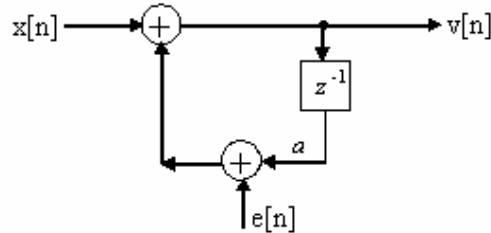


Figura 5.14. Modelul zgomotului aditiv pentru eroarea de cuantizare a produsului pentru un filtru cu un singur pol

Secvența de ieșire a filtrului $v[n]$, poate fi separată în două componente. Prima este răspunsul sistemului, $y[n]$, la secvența de intrare $x[n]$, iar a doua este răspunsul sistemului, $z[n]$, la zgomotul aditiv de cuantizare $e[n]$. Secvența de ieșire se exprimă ca o sumă a acestor două componente, adică

$$v[n] = y[n] + z[n] \quad (5.112)$$

Înlocuind $v[n]$ din (5.112) în (5.111), se obține

$$y[n] + z[n] = ay[n-1] + az[n-1] + x[n] + e[n] \quad (5.113)$$

Pentru a simplifica analiza, se fac următoarele presupuneri în legătură cu eroarea $e[n]$:

1. Pentru orice n , secvența de eroare $\{e[n]\}$ este uniform distribuită în intervalul $\left(-\frac{1}{2}2^{-b}, \frac{1}{2}2^{-b}\right)$. Aceasta implică valoarea medie a lui $\{e[n]\}$ egală cu zero, și dispersia

$$\sigma_e^2 = \frac{2^{-2b}}{12} \quad (5.114)$$

2. Eroarea $\{e[n]\}$ este o secvență staționară de zgomot alb și, ca urmare, $e[n]$ și $e[m]$ sunt necorelate pentru $n \neq m$.
3. Secvența de eroare $\{e[n]\}$ este necorelată cu semnalul $\{x[n]\}$.

Ultima presupunere permite separarea ecuației cu diferențe (5.113) în două ecuații independente:

$$y[n] = ay[n-1] + x[n] \quad (5.115)$$

$$z[n] = az[n-1] + e[n] \quad (5.116)$$

Ecuația cu diferențe (5.115) reprezintă relația de intrare-ieșire pentru sistemul dorit, iar cea din (5.116) reprezintă relația pentru eroarea de cuantizare la ieșirea sistemului.

Pentru a completa analiza se face apel la două relații importante. Prima este relația pentru valoarea medie a ieșirii $z[n]$ pentru un filtru liniar, invariant în timp, cu răspunsul la impuls $h[n]$, când este excitat de o secvență $e[n]$ cu media m_e . Rezultatul este [48]

$$m_z = m_e \sum_{n=0}^{\infty} h[n] \quad (5.117)$$

sau, echivalent,

$$m_z = m_e H(0) \quad (5.118)$$

unde $H(0)$ valoarea răspunsului în frecvență $H(\omega)$ la $\omega = 0$.

Deoarece eroarea de cuantizare datorată rotunjirii are media zero, valoarea medie a erorii la ieșire este $m_z=0$.

A doua relație importantă este expresia pentru secvența de autocorelație a ieșirii $z[n]$ a unui filtru cu răspunsul la impuls $h[n]$ la secvența aleatoare de intrare $e[n]$. Aceasta este [63]

$$\gamma_{zz}[n] = \sum_{k=0}^{\infty} \sum_{l=0}^{\infty} h[k]h[l]\gamma_{ee}[k-l+n] \quad (5.119)$$

unde $\gamma_{ee}[n]$ este funcția de autocorelație a secvenței de intrare $e[n]$.

În cazul particular când secvența aleatoare este zgomot alb, secvența de autocorelație $\gamma_{ee}[n]$ este un impuls scalat cu dispersia σ_e^2 , adică [34]

$$\gamma_{ee}[n] = \sigma_e^2 \delta[n] \quad (5.120)$$

După substituția relației (5.120) în (5.119), se obține secvența de autocorelație de la ieșirea filtrului excitat cu zgomot alb

$$\gamma_{zz}[n] = \sigma_e^2 \sum_{k=0}^{\infty} h[k]h[k+n] \quad (5.121)$$

Dispersia σ_z^2 a zgomotului de ieșire este obținută evaluând $\gamma_{zz}[n]$ la $n = 0$, adică [34]

$$\sigma_z^2 = \sigma_e^2 \sum_{k=-\infty}^{\infty} h^2[k] \quad (5.122)$$

sau, cu ajutorul teoremei lui Parseval [63], expresia alternativă

$$\sigma_z^2 = \frac{\sigma_e^2}{2\pi} \int_{-\pi}^{\pi} |H(\omega)|^2 d\omega \quad (5.123)$$

În cazul filtrului cu un singur pol, răspunsul la impuls este

$$h[n] = a^n u[n] \quad (5.124)$$

Dispersia erorii la ieșirea filtrului rezultă

$$\sigma_z^2 = \sigma_e^2 \sum_{k=0}^{\infty} a^{2k} = \frac{\sigma_e^2}{1-a^2} \quad (5.125)$$

Se observă că puterea zgomotului σ_z^2 la ieșirea filtrului este mărită față de puterea zgomotului de la intrare, σ_e^2 , cu factorul $1/(1-a^2)$. Acest factor crește odată cu apropierea polului de cercul unitate.

Fie, în continuare, un filtru recursiv de ordinul doi:

$$y[n] = -a_1 y[n-1] - a_2 y[n-2] + b_0 x[n] + b_1 x[n-1] \quad (5.126)$$

În calculul ieșirii sunt implicate patru multiplicări, dacă a_1, a_2, b_0 și b_1 nu sunt egali cu unitatea. Zgomotul de rotunjire asociat cu fiecare multiplicare este $e_i[n]$, $i = \overline{0,3}$.

Se consideră întâi realizarea în forma directă I, ca în figura 5.15.

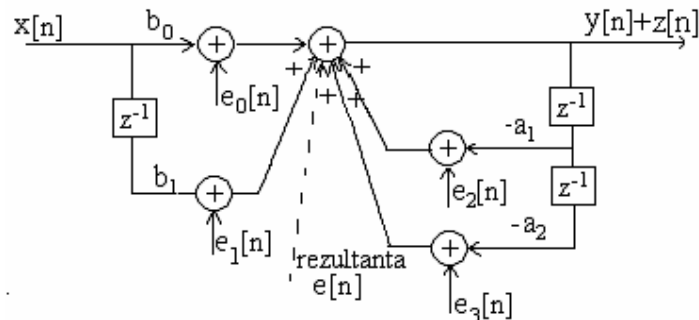


Fig. 5.15. Zgomotul de rotunjire la multiplicare pentru un filtru de ordinul doi în forma directă I

Deoarece toate sursele de zgomot se adună în același punct, acestea pot fi înlocuite cu o sursă de zgomot echivalentă

$$e[n] = \sum_{i=0}^3 e_i[n] \quad (5.127)$$

Se observă că în implementarea în forma directă I, zgomotul trece numai prin partea de sistem ce conține numai poli, adică zerourile nu au nici un efect asupra zgomotului din ieșire.

În cazul rotunjirii, când pasul de cuantizare este constant, dispersia unei surse de zgomot este

$$\sigma_{e_i}^2 = \frac{\Delta^2}{12}, \quad i = 0, 1, 2, 3. \quad (5.128)$$

Presupunând erorile de cuantizare independente, dispersia zgomotului rezultat este suma dispersiilor componentelor

$$\sigma_{e_i}^2 = \sum_{i=0}^3 \sigma_{e_i}^2 = \frac{\Delta^2}{3} \quad (5.129)$$

Pentru cazul general al formei directe I, când sistemul are $M+1$ multiplicări pentru zerouri și N multiplicări pentru poli cu coeficienți diferiți de 0 și 1, dispersia surselor de zgomot este

$$\sigma_e^2 = (M + N + 1) \frac{\Delta^2}{12} \quad (5.130)$$

Porțiunea din filtru prin care trece zgomotul de rotunjire este arătată în figura 5.16. Ieșirea $z[n]$ datorată zgomotului formează o parte a ieșirii cuantizate.

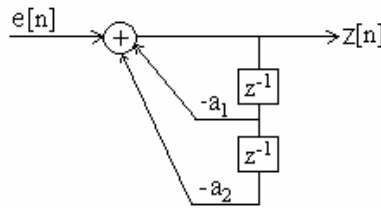


Figura 5.16. Porțiunea din filtrul recursiv afectată de zgomotul de rotunjire pentru realizarea în forma directă I.

Pentru figura 5.16 se poate scrie

$$\frac{Z(z)}{E(z)} = H'(z) = \frac{1}{1 + a_1 z^{-1} + a_2 z^{-2}} \quad (5.131)$$

Evident, această funcție de transfer diferă de cea a filtrului care include și zerouri, care este

$$H(z) = \frac{b_0 + b_1 z^{-1}}{1 + a_1 z^{-1} + a_2 z^{-2}} \quad (5.132)$$

Conform relației (5.74), dispersia totală de regim permanent a ieșirii datorate zgomotului de rotunjire este

$$\sigma_{0_{zss}}^2 = \frac{\Delta^2}{3} \sum_{\substack{\text{polii din} \\ \text{cercul unitate}}} \text{reziduurile lui } H'(z) H'(z^{-1}) z^{-1} \quad (5.133)$$

cu $H'(z)$ dat de (5.131).

În cazul formei directe I dispersia totală de regim permanent a zgomotului datorat rotunjirii multiplicărilor este

$$\begin{aligned} \sigma_{0zss}^2 &= (M + N + 1) \frac{\Delta^2}{12} \frac{1}{2\pi j} \oint H'(z) H'(z^{-1}) z^{-1} dz = \\ &= (M + N + 1) \frac{\Delta^2}{12} \sum_n |h'[n]|^2 \end{aligned} \quad (5.134)$$

unde $H'(z) = \frac{1}{1 + \sum_{k=1}^N a_k z^{-k}}$ este partea care conține toți polii sistemului.

În continuare, se consideră implementarea canonică (forma directă II) a filtrului descris de (5.126), caz în care erorile de rotunjire pot fi reprezentate ca surse de zgomot poziționate ca în figura 5.17.

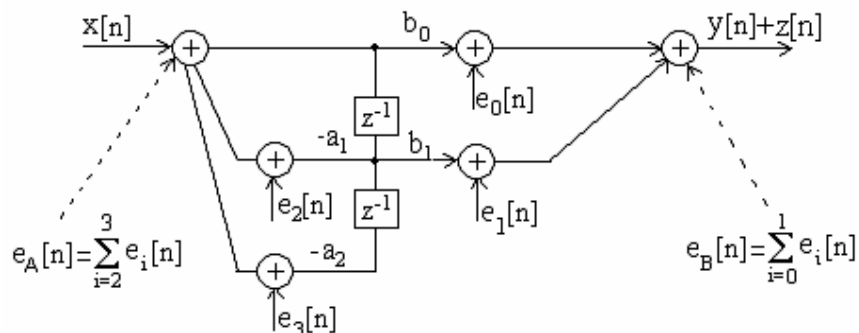


Figura 5.17. Zgomotul de rotunjire al produselor pentru un filtru recursiv implementat în forma canonică

Se observă că semnalul de eroare $e_A[n] = \sum_{i=2}^3 e_i[n]$ cu dispersia $\frac{\Delta^2}{6}$, trece prin tot filtru, în timp ce $e_B[n] = \sum_{i=0}^1 e_i[n]$ cu dispersia $\frac{\Delta^2}{6}$ este un zgomot adunat direct la ieșire. În acest caz dispersia de regim permanent a ieșirii datorată zgomotului de rotunjire a produselor este suma dispersiilor zgomotelor determinate de cele două semnale de eroare $e_A[n]$ și $e_B[n]$.

$$\sigma_{0zss}^2 = \frac{\Delta^2}{6} \left[1 + \sum_{\substack{\text{polii din interiorul} \\ \text{cercului unitate}}} \text{reziduurile lui } H(z) H(z^{-1}) z^{-1} \right] \quad (5.135)$$

cu $H(z)$ dat de (5.132).

Pentru cazul general al formei directe II pentru filtrul IIR, când coeficienții acestuia sunt diferiți de 0 și 1, dispersia de regim permanent a zgomotului de ieșire este

$$\begin{aligned} \sigma_{0zss}^2 &= N \frac{\Delta^2}{12} \frac{1}{2\pi j} \oint_c H(z) H(z^{-1}) z^{-1} dz + (M+1) \frac{\Delta^2}{12} = \\ &= N \frac{\Delta^2}{12} \sum_n |h[n]|^2 + (M+1) \frac{\Delta^2}{12} \end{aligned} \quad (5.136)$$

Fără a considera valori numerice pentru coeficienți, numai din compararea relațiilor (5.134) și (5.136), nu este posibil a decide care dintre aceste forme de implementare produce un zgomot de ieșire mai mic datorat erorii de cuantizare a produselor.

Exemplul 5. 10.

Să se determine dispersia de regim permanent a zgomotului de ieșire, datorat rotunjirii aritmetice, a filtrului cu funcția de sistem

$$H(z) = \frac{b_0 + b_1 z^{-1}}{1 - 2r \cos \theta z^{-1} + r^2 z^{-2}}$$

implementat în

- a) formă directă I
- b) forma directă II

dacă $r=0,9$, $\theta = \pi/4$, $b_0 = 1, 1$, $b_1 = 0, 3$ și pasul de cuantizare Δ .

Soluție. a) Din figura 5.15 și 5.16 rezultă că dispersia de regim permanent a zgomotului de ieșire este

$$\begin{aligned} \sigma_{0zss}^2 &= \frac{\Delta^2}{3} \sum_{\substack{\text{polii din cercul} \\ \text{unitate ai lui } H'(z)}} \text{reziduurile lui } H'(z) H'(z^{-1}) z^{-1} = \\ &= \frac{\Delta^2}{3} \frac{1+r^2}{1-r^2} \frac{1}{r^4 - 2r^2 \cos 2\theta + 1} = 1,92 \Delta^2 \end{aligned}$$

$$\text{cu } H'(z) = \frac{1}{1 - 2r \cos \theta z^{-1} + r^2 z^{-2}}$$

b) Din figura (5.17) rezultă

$$\begin{aligned}\sigma_{0zss}^2 &= \frac{\Delta^2}{6} + \frac{\Delta^2}{6} \sum_{\text{polii lui } H(z)} \text{reziduurile lui } H(z) H(z^{-1}) z^{-1} = \\ &= \frac{\Delta^2}{6} \left[1 + \frac{(b_0^2 + b_1^2)(1 + r^2) - 4b_0 b_1 r \cos \theta}{(r^4 - 2r^2 \cos 2\theta + 1)(1 - r^2)} \right] = 1,07 \Delta^2\end{aligned}$$

Se observă că forma directă II (canonică) produce un zgomot de ieșire mai mic pentru valorile date ale parametrilor decât forma directă I și că valorile b_0 și b_1 nu afectează dispersia zgomotului de ieșire în forma directă I.

Ecuatiile (5.134) și (5.136) arată că structurile în forma directă I și II sunt afectate diferit de cuantizarea produselor în implementarea ecuațiilor cu diferențe corespunzătoare. În general, alte structuri echivalente, cum ar fi cele în cascadă, în paralel, lattice și formele transpuse vor avea dispersii totale ale zgomotului la ieșire diferite de cele din structurile în formă directă. Nu se poate spune care sistem va avea dispersia de zgomot la ieșire cea mai mică, dacă nu se cunosc valorile coeficienților.

Îmbunătățirea performanței de zgomot a sistemelor numerice este posibilă folosind sumatoare și acumulate pe un număr mai mare de biți. Această soluție presupune însă o complicare semnificativă a realizării “hard” a schemei.

5.6. Oscilații cu ciclu-limită în sisteme recursive

În secțiunile anterioare au fost analizate erorile care apar în operațiile aritmetice realizate de un filtru digital. Prezența unuia sau a mai multor cuantizoare în implementarea unui filtru digital, conduce la un dispozitiv neliniar a cărui caracteristică poate fi semnificativ diferită de cea a filtrului ideal. Efectele neliniare datorate aritmeticii cu precizie finită, îngreunează analiza performanțelor unui filtru digital. Pentru a efectua o analiză a efectului cuantizării, s-a adoptat o caracterizare statistică a erorilor de cuantizare, ceea ce a condus în final la un model liniar pentru filtru.

În sistemele recursive, neliniaritatea datorată efectuării operațiilor matematice în aritmetică finită poate cauza oscilații periodice la ieșire, chiar dacă secvența de intrare este zero sau o valoare constantă, nenulă. Astfel de oscilații în sistemele recursive sunt numite *cicluri limită* și pot fi direct atribuite erorii de rotunjire sau trunchiere la multiplicare. Aceste

oscilații pot fi reduse folosind registre pe mai mulți biți. Al doilea tip de oscilații numit *oscilații de depășire* poate apărea când intrarea cuantizorului depășește domeniul dinamic. Aceste oscilații au, de obicei, amplitudine mare și nu pot fi reduse prin creșterea numărului de biți.

5.6.1. Cicluri limită datorate rotunjirii

Fenomenul ciclurilor limită este diferit de comportamentul zgomotului cauzat de cuantizare. Efectele cuantizării se identifică cu zgomotul când nivelul semnalului este mare și foarte variabil, făcând eroarea de cuantizare, la orice moment de timp, aproape independentă de erorile anterioare. Când nivelul semnalului este scăzut, erorile cauzate de cuantizare devin corelate. Ciclurile limită sunt periodice, dar nu neapărat sinusoidale. Ele sunt susceptibile a apărea acolo unde există reacție în filtru; filtrele IIR au întotdeauna mecanisme de reacție în interiorul lor, deci astfel de oscilații pot apărea la ieșirea lor. Spre deosebire de acestea, filtrele FIR nu conțin mecanisme de reacție și, în consecință, ele nu vor prezenta oscilații la ieșire. Acesta este un avantaj al filtrelor FIR față de cele IIR. Tratarea generală a comportării pe cicluri limită a filtrelor digitale este dificilă, motiv pentru care se vor analiza structurile de ordinul 1 și 2.

Pentru a ilustra caracteristica unei oscilații de ciclu limită, se consideră un sistem cu un singur pol, descris de ecuația liniară cu diferențe

$$y[n] = ay[n-1] + x[n] \quad (5.137)$$

în care polul este situat la $z=a$. Sistemul ideal este prezentat în figura 5.18a.

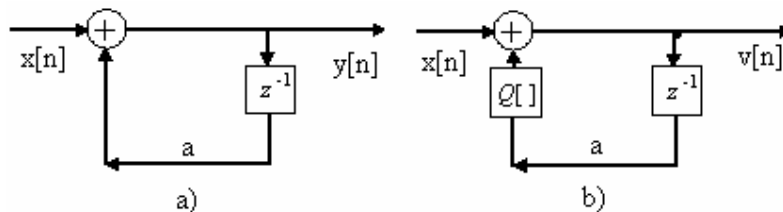


Figura 5.18. a) Sistemul recursiv ideal cu un singur pol b) Sistemul neliniar real

Sistemul real, care este descris de ecuația neliniară cu diferențe

$$v[n] = Q[av[n-1]] + x[n] \quad (5.138)$$

este realizat ca în figura 5.18b.

Se presupune că sistemul real din figura 5.18b este implementat cu o aritmetică în virgulă fixă cu patru biți pentru amplitudine și un bit pentru semn. Cuantizarea care se face după multiplicare este presupusă a rotunji produsul prin adaos. În Tabelul 5.4 se prezintă răspunsul sistemului real pentru patru poziții diferite ale polului $z=a$ și intrarea $x[n]=\beta\delta[n]$, unde $\beta=15/16$, care are reprezentarea binară 0,1111.

Tabel 5.4 Cicluri limită pentru un filtru cu un singur pol

n	$a=0,1000$ = 1/2	$a=1,1000$ = -1/2	$a=0,1100$ = 3/4	$a=1,1000$ = -3/4
0	0,1111 (15/16)	0,1111 (15/16)	0,1011 (11/16)	0,1011 (11/16)
1	0,1000 (7/16)	1,1000 (-7/16)	0,1000 (8/16)	1,1000 (-8/16)
2	0,0100 (3/16)	0,0100 (3/16)	0,0110 (6/16)	0,0110 (6/16)
3	0,0010 (1/16)	1,0010 (-1/16)	0,0101 (5/16)	1,0101 (-5/16)
4	0,0001 (1/16)	0,0001 (1/16)	0,0100 (4/16)	0,0100 (4/16)
5	0,0001 (1/16)	1,0001 (-1/16)	0,0011 (3/16)	1,0011 (-3/16)
6	0,0001 (1/16)	0,0001 (1/16)	0,0010 (2/16)	0,0010 (2/16)
7	0,0001 (1/16)	1,0001 (-1/16)	0,0010 (2/16)	1,0010 (-2/16)
8	0,0001 (1/16)	0,0001 (1/16)	0,0010 (2/16)	0,0010 (2/16)

În mod ideal, răspunsul sistemului ar trebui să scadă exponențial spre zero ($y[n]=a^n \rightarrow 0$ când $n \rightarrow \infty$). În sistemul real, totuși, răspunsul $v[n]$ atinge o stare stabilă periodică la ieșire, cu o perioadă ce depinde de valoarea polului. Când polul este pozitiv, oscilațiile au loc cu perioada $N_p=1$, astfel încât ieșirea atinge o valoare constantă de 1/16 pentru $a=1/2$ și 1/8 pentru $a=3/4$. Acest fenomen este numit ciclu limită cu frecvență zero.

Pe de altă parte, când polul este negativ, secvența de ieșire oscilează între valori pozitive și negative ($\pm 1/16$ pentru $a = -1/2$ și $\pm 1/8$ pentru $a = -3/4$). Prin urmare, perioada este $N_p = 2$. Se obține astfel o oscilație de amplitudine constantă, a cărei pulsație este egală cu π și a cărei amplitudine este $\pm 1/16$ sau $\pm 1/8$.

Aceste cicluri-limită apar ca rezultat al efectului de cuantizare în multiplicări. Când secvența de intrare $x[n]$ devine zero, ieșirea intră într-un ciclu limită după un număr de iterații. Ieșirea rămâne în acest ciclu limită până când este aplicat un alt semnal de intrare, suficient de puternic, pentru a scoate sistemul din ciclu. În mod similar, ciclurile limită cu intrare zero apar din condiții inițiale nenule. Amplitudinea ieșirii

pe perioada ciclului limită este inclusă într-un domeniu de valori care este numit “banda moartă” a filtrului. Frecvența și amplitudinea ciclului limită depind de coeficienți, condiții inițiale, metoda de cuantizare și lungimea cuvântului.

Este interesant de menționat faptul că atunci când răspunsul filtrului cu un pol este în ciclu limită, sistemul neliniar real lucrează ca un sistem liniar echivalent, cu un pol la $z=1$, atunci când polul este pozitiv ($a>0$), și $z = -1$, când polul este negativ ($a<0$). Aceasta înseamnă

$$Q_r[av[n-1]] = \begin{cases} v[n-1], & a > 0 \\ -v[n-1], & a < 0 \end{cases} \quad (5.139)$$

Deoarece produsul $av[n-1]$ este rotunjit, eroarea de cuantizare este limitată de

$$|Q_r[av[n-1]] - av[n-1]| \leq \frac{1}{2} 2^{-b} \quad (5.140)$$

unde b este numărul de biți (exclusiv semnul) utilizat în reprezentarea polului a și a lui $v[n]$. Prin urmare, relațiile (5.139) și (5.140) conduc la

$$|v[n-1]| - |av[n-1]| \leq \frac{1}{2} 2^{-b}$$

și, deci

$$|v[n-1]| \leq \frac{\frac{1}{2} 2^{-b}}{1-|a|} \quad (5.141)$$

Când coeficientul a este pozitiv, răspunsul ciclului limită se numește *de curent continuu* (are amplitudine și semn constante), iar dacă a este negativ comportamentul ciclului limită are amplitudine constantă dar semn alternant.

Expresia din (5.141) definește *zona sau banda moartă* pentru un filtru cu un singur pol. De exemplu, când $b = 4$ și $|a| = 1/2$ banda moartă este cuprinsă în domeniul $(-1/16, 1/16)$ pentru amplitudini, iar pentru $b = 4$ și $|a| = 3/4$, banda moartă crește la $(-1/8, 1/8)$.

Comportarea ciclului limită în cazul unui filtru cu doi poli este mult mai complexă prin faptul că poate apărea o mai mare varietate de oscilații. În acest caz sistemul ideal cu doi poli este descris de ecuația liniară cu diferențe

$$y[n] = a_1 y[n-1] + a_2 y[n-2] + x[n] \quad (5.142)$$

în timp ce sistemul real este descris de ecuația neliniară cu diferențe

$$v[n] = Q_r[a_1 v[n-1]] + Q_r[a_2 v[n-2]] + x[n] \quad (5.143)$$

Când coeficienții filtrului satisfac condiția $a_1^2 < -4a_2$, polii sistemului apar la $z_{1,2} = re^{\pm j\theta}$, unde $a_2 = -r^2$ și $a_1 = 2r\cos\theta$. Ca și în cazul filtrului cu un singur pol, când sistemul este într-un ciclu limită cu intrare zero [49],

$$Q_r[a_2 v[n-2]] = -v[n-2], \quad (5.144)$$

adică sistemul se comportă ca un oscilator cu polii complex-conjugați situați pe cercul unitate ($a_2 = -r^2 = -1$). Rotunjirea produsului $av[n-2]$ implică

$$|Q_r[a_2 v[n-2]] - a_2 v[n-2]| \leq \frac{1}{2} 2^{-b} \quad (5.145)$$

După substituția lui (5.144) în (5.145), se obține

$$|v[n-2] - a_2 v[n-2]| \leq \frac{1}{2} 2^{-b}$$

sau, echivalent

$$|v[n-2]| \leq \frac{\frac{1}{2} 2^{-b}}{1 - |a_2|} \quad (5.146)$$

Expresia din (5.146) definește banda moartă a unui filtru de ordin doi cu poli complex conjugați. Se observă că limitele benzii moarte depind doar de a_2 . Parametrul $a_1 = 2r\cos\theta$ determină doar frecvența oscilațiilor.

Un alt ciclu limită posibil cu intrarea zero, care este numai amintit și care apare ca rezultat al rotunjirii multiplicărilor, corespunde unui sistem echivalent de ordinul doi cu polii la $z = \pm 1$.

Este interesant de menționat cum ciclurile limită descrise anterior au rezultat prin rotunjirea produsului dintre coeficienții filtrului și ieșirile precedente $v[n-1]$ și $v[n-2]$. În locul rotunjirii, se poate alege a trunchia produsul la b biți, caz în care se pot elimina multe din ciclurile limită, dar această soluție nu este foarte agreată, deoarece trunchierea are ca rezultat o deplasare a valorii medii a erorii, excepție făcând cazul când se folosește reprezentarea semn-valoare unde eroarea de trunchiere este simetrică față de zero.

În realizarea în paralel a diverselor sisteme IIR de ordin înalt cu secțiuni de ordinul doi, fiecare secțiune generează propriul ciclu limită, fără interacțiune între secțiunile de filtru de ordin doi. Prin urmare, ieșirea

este o sumă a ciclurilor limită cu intrare zero a secțiunilor individuale. În cazul realizării în cascadă pentru un sistem IIR de ordin înalt, ciclurile limită sunt mult mai greu de analizat. În particular, când prima secțiune de filtru generează un ciclu limită cu intrare zero, acesta este filtrat de secțiunile succesive. Dacă frecvența ciclului limită este apropiată de frecvența de rezonanță a filtrului următor din succesiune, amplitudinea secvenței va fi mărită de caracteristica de rezonanță. În general, trebuie evitate astfel de situații.

5.6.2. Cicluri limită datorate depășirii

Un tip mult mai sever de cicluri limită poate apărea datorită depășirii aritmetice din interiorul filtrelor care folosesc aritmetica în complement față de unu sau în complement față de doi. Aceste cicluri limită sunt cunoscute sub numele de *oscilații de depășire*. O depășire la adunarea a două sau mai multe numere binare apare atunci când suma depășește lungimea disponibilă a cuvântului la implementarea digitală a sistemului.

De exemplu, se consideră secțiunea de filtru de ordin doi prezentată în figura 5.19, în care adunarea se face în aritmetica complementului față de doi.

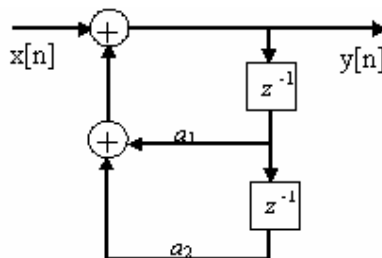


Figura 5.19. Secțiune de filtru de ordinul doi

Ieșirea din filtru se poate scrie

$$y[n] = g[a_1 y[n-1] + a_2 y[n-2] + x[n]] \quad (5.147)$$

unde funcția $g[\cdot]$ reprezintă adunarea în complement față de doi.

Figura 5.20 prezintă caracteristica intrare-ieșire $g[v]$ a sumatorului în complement față de doi.

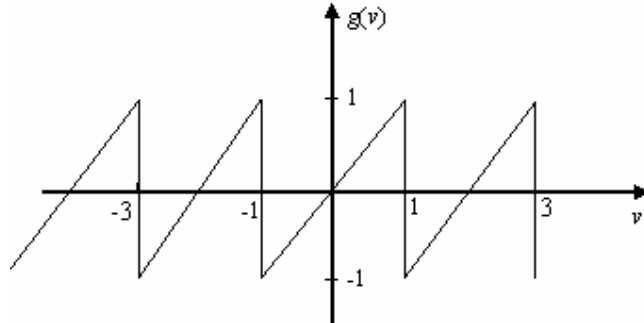


Figura 5.20. Caracteristica funcțională pentru adunarea în complement față de doi a două sau mai multe numere

Domeniul de valori al parametrilor (a_1, a_2) pentru un filtru stabil este precizat de triunghiul de stabilitate [63]. Totuși, aceste condiții nu sunt de ajuns pentru a preveni oscilațiile datorate depășirii din aritmetica în complement față de doi. Condiția necesară și suficientă pentru a nu apărea cicluri limită datorate depășirii, este [49]

$$|a_1| + |a_2| < 1 \quad (5.148)$$

care este o condiție extrem de restrictivă și duce la o constrângere nerezonabilă asupra oricărei secțiuni de filtru de ordin doi.

Un remediu efectiv pentru rezolvarea problemei oscilațiilor provocate de depășire este de a modifica caracteristica sumatorului, ca în figura 5.21, care operează cu saturare numerică. Atunci când este sesizată o depășire (sau o subdepășire), ieșirea sumatorului va avea valoarea maximă de capăt de scară ± 1 . Distorsiunea cauzată de această neliniaritate în sumator este de obicei mică deoarece saturația apare rar. Folosirea unei astfel de neliniarități nu elimină necesitatea scalării semnalelor și a parametrilor sistemului, așa cum va fi descris în paragraful următor.

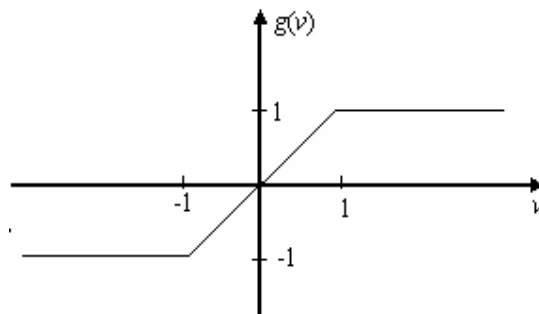


Figura 5.21. Caracteristica funcțională pentru adunare cu saturare la ± 1

Ilustrarea oscilațiilor datorate depășirii se face pe exemplul următor. Se consideră secțiunea de filtru de ordin doi caracterizată de ecuația (5.143) în care adunarea se face în aritmetica complementului față de doi, cu lungimea cuvintelor de 4 biți, incluzând bitul de semn, și se folosește rotunjirea pentru reprezentările în complement față de doi. Se presupune că $a_1 = 3/4 = 0,110$ și $a_2 = -3/4 = 1,010$ și, de asemenea, că $x[n]$ rămâne zero pentru $n \geq 0$.

Se consideră condițiile inițiale $v[-1] = (3/4)_{10} = (0,110)_{2C}$ și $v[-2] = (-3/4)_{10} = (1,010)_{2C}$. Eșantionul de la ieșire la momentul $n=0$ va fi $v[0] = 0,110 \cdot 0,110 + 1,010 \cdot 1,010 = 0,100100 + 0,100100$.

Dacă se rotunjește fiecare produs, rezultă

$$v[0] = 0,101 + 0,101 = 1,010 = -3/4.$$

În mod similar se obține

$$v[1] = 1,011 + 1,011 = 0,110 = 3/4,$$

adică, $v[n]$ va continua să oscileze între $-3/4$ și $3/4$ până ce este aplicat un semnal de intrare care să scoată sistemul din acest ciclu limită. Acesta este un exemplu de oscilații de depășire. Sistemele de ordin mai mare au o comportare mai complexă.

5.7. Scalarea pentru prevenirea depășirii

Saturația aritmetică descrisă în paragraful anterior elimină ciclurile limită datorate depășirii pe de o parte, dar, pe de altă parte, duce la distorsiuni nedorite ale semnalelor, în acest caz nemaifuncționând regula conform căreia, dacă se adună mai multe numere a căror sumă este de modul subunitar, rezultatul este corect, chiar dacă apar depășiri în etapele intermediare de calcul. Pentru a limita aceste distorsiuni neliniare se scalează semnalul de intrare și răspunsul la impuls între intrare și orice nod din sistem, astfel încât să nu se depășească gama dinamică.

Efectul depășirii este mult mai sever pentru un filtru recursiv, decât pentru unul nerecursiv, deoarece erorile sunt filtrate din nou (datorită reacției) ceea ce face ca filtrul să devină inutilizabil în scurt timp. Pentru ambele tipuri de filtre, scalarea este necesară pentru reducerea amplitudinii semnalelor în anumite limite, evitându-se depășirea în condiții normale de lucru. Există mai multe reguli de scalare, care vor fi prezentate în cele ce urmează.

5.7.1. Norme de scalare

5.7.1.1. Scalarea după norma l_1

Se analizează toate nodurile în care ar putea apărea depășiri și fiecare nod din rețea este constrâns să aibă o amplitudine mai mică decât 1, pentru a evita depășirea. Dacă $w_i[n]$ reprezintă valoarea variabilei asociată nodului i iar $h_i[n]$ este răspunsul la impuls de la nodul de intrare, căruia îi este asociată variabila $x[n]$, până la nodul i , atunci se poate scrie

$$|w_i[n]| = \left| \sum_{m=0}^{\infty} x[n-m]h_i[m] \right|. \quad (5.149)$$

Considerând că $x[n-m]$ are valoarea maximă x_{\max} , rezultă

$$|w_i[n]| \leq x_{\max} \sum_{m=0}^{\infty} |h_i[m]|. \quad (5.150)$$

O condiție suficientă ca $|w_i[n]| < 1$ este ca

$$x_{\max} < \frac{1}{\sum_{m=0}^{\infty} |h_i[m]|} \quad (5.151)$$

pentru toate nodurile din rețea. Mărima $l_1 = \|h_i\|_1 = \sum_{m=0}^{\infty} |h_i[m]|$ se numește *norma l_1* a lui h_i . Dacă x_{\max} nu satisface ecuația (5.151), atunci se poate

multiplica $x[n]$ cu factorul de scalare $s_1 < \min_i \left\{ \frac{1}{\|h_i\|_1} \right\}$ la intrarea sistemului, astfel încât $s_1 x_{\max}$ să satisfacă (5.151) pentru toate nodurile din rețea, adică

$$s_1 x_{\max} < \frac{1}{\max_i \left[\sum_{m=0}^{\infty} |h_i[m]| \right]} \quad (5.152)$$

Scalând intrarea pe această cale se garantează că depășirea nu apare niciodată la nici unul din nodurile de rețea. La ieșire se compensează scalarea prin înmulțirea cu $\frac{1}{s_1}$, astfel încât să nu se modifice funcția de transfer a filtrului. Relația (5.152) conduce la o

scalare foarte severă, care se mai numește și *scalare de sumă*. În practică scalarea nu este făcută niciodată așa puternic, pentru că înrăutățește raportul semnal-zgomot, fapt ce va fi arătat ulterior.

5.7.1.2. Scalarea după norma l_∞

Dacă se dispune de cunoștințe suplimentare despre intrare, se poate alege factorul de scalare, s_∞ , mai mare, astfel încât să se garanteze lipsa depășirii. Dacă intrarea este un semnal de bandă îngustă modelat cu $x[n] = x_{\max} \cos(\omega_0 n)$, variabilele de noduri vor fi [39]

$$w_i[n] = |H_i(\omega_0)| x_{\max} \cos[\omega_0 n + \angle H_i(\omega_0)] \quad (5.153)$$

Depășirea este evitată pentru toate semnalele armonice dacă

$$\max_{i, |\omega| \leq \pi} |H_i(\omega)| x_{\max} < 1 \quad (5.154)$$

Mărimea $l_\infty = \|H_i\|_\infty = \max_{|\omega| \leq \pi} |H_i(\omega)|$ se numește norma l_∞ a lui H_i .

Dacă intrarea este scalată prin factorul de scalare $s_\infty < \min_i \left\{ \frac{1}{\|H_i\|_\infty} \right\}$ rezultă

$$s_\infty x_{\max} < \frac{1}{\max_{i, |\omega| \leq \pi} |H_i(\omega)|} \quad (5.155)$$

5.7.1.3. Scalarea după norma l_2

O altă abordare posibilă este de a scala intrarea astfel încât energia fiecărei variabile de nod să fie mai mică sau egală cu energia totală a secvenței de intrare. Se poate obține scalarea corespunzătoare folosind inegalitatea Schwartz Buniacovski și teorema lui Parseval [63].

$$\begin{aligned} |w_i[n]|^2 &= \left| \sum_{k=0}^{\infty} h_i[k] x[n-k] \right|^2 \leq \sum_{k=0}^{\infty} |h_i[k]|^2 \sum_{k=0}^{\infty} |x[n-k]|^2 = \\ &= \sum_{k=0}^{\infty} |h_i[k]|^2 \sum_{k=0}^{\infty} |x[k]|^2 \end{aligned} \quad (5.156)$$

Pentru a asigura condiția de nedepășire a energiei semnalului de intrare de către variabilele de noduri, adică $|w_i[n]|^2 \leq \sum_{n=0}^{\infty} |x[n]|^2$, unde

$\sum_{n=0}^{\infty} |x[n]|^2 = E_x$ este energia semnalului de intrare, se poate multiplica secvența $x[n]$ cu factorul de scalare s_2 , ales astfel încât

$$s_2^2 \leq \frac{1}{\max_i \sum_{n=0}^{\infty} |h_i[n]|^2} = \frac{1}{\max_i \frac{1}{2\pi} \int_{-\pi}^{\pi} |H_i(\omega)|^2 d\omega} \quad (5.157)$$

Mărima $l_2 = \|h_i\|_2 = \left(\sum_{n=0}^{\infty} |h_i[n]|^2 \right)^{1/2}$ se numește *norma* l_2 a lui h_i .

5.7.1.4. Scalarea după norma l_p

Metodele anterioare pot fi generalizate în sensul normei l_p .

Norma l_p a unei transformate Fourier $H(\omega)$ este definită ca [39]

$$l_p = \|H\|_p = \left[\frac{1}{2\pi} \int_{-\pi}^{\pi} |H(\omega)|^p d\omega \right]^{1/p} \quad (5.158)$$

Se poate arăta că, în general, este îndeplinită inegalitatea [26]

$$|w_i[n]| \leq \|X\|_p \|H_i\|_q \quad (5.159)$$

unde p și q sunt întregi astfel încât

$$\frac{1}{p} + \frac{1}{q} = 1. \quad (5.160)$$

Pentru orice secvență $h[n]$ cu transformata Fourier $H(\omega)$ există relația [23]

$$\|H\|_{\infty} \geq \|H\|_p, \text{ oricare ar fi } p \in N^*.$$

Ca urmare, scalarea l_{∞} reduce nivelele de semnal într-o măsură mai mare decât alte scalări de tip l_p . Cele mai folosite scalări sunt l_2 , l_{∞} , precum și scalarea de sumă. Se poate arăta că există relația [23]

$$\sum_{n=0}^{\infty} |h_i[n]|^2 \leq \max_{i,\omega} |H_i(\omega)| \leq \sum_{n=0}^{\infty} |h_i[n]|, \quad (5.161)$$

$$\text{adică } l_2 \leq l_{\infty} \leq l_1.$$

Dintre acestea, cea mai severă este scalarea de sumă, care este și dificil de calculat. Cel mai ușor de evaluat analitic este relația (5.157), deoarece această integrală poate fi calculată folosind teorema reziduurilor a lui Cauchy [1].

Deoarece în implementarea filtrelor recursive intervin mai multe puncte de sumare, ieșirea fiecăruia trebuie scalată pentru a evita depășirea, deci vor fi mai multe răspunsuri la impuls $h_i[n]$ și funcții de sistem corespunzătoare, $H_i(z)$, care fac legătura între intrarea $x[n]$ și semnalele intermediare $w_i[n]$.

5.7.2. Interacțiunea dintre domeniul dinamic și zgomot

Normele de scalare l_2, l_∞, l_1 reprezintă trei moduri de a obține coeficienți de scalare pentru intrarea unui filtru digital. Prin scalarea intrării cu factorul $s_p, p = 1, \infty, 2$, raportul semnal / zgomot de cuantizare la ieșire scade.

În figura 22 a,b, se prezintă un sistem IIR de ordinul doi, implementat în forma directă I și forma directă II, cu intrarea scalată. În figura 22 a, factorul de scalare s-a combinat cu coeficienții b_k , astfel încât sursa de zgomot este aceeași ca în cazul fără scalare, prezentat în figura 5.15. Deoarece acest zgomot este filtrat din nou de partea de filtru care conține polii, puterea zgomotului de ieșire este aceeași pentru sistemul nescalat, reprezentat în figura 5.15 și cel scalat, reprezentat în 5.22a. Pentru sistemul din figura 22a, funcția de sistem este $s_p H(z)$, față de $H(z)$ a sistemului cu intrarea nescalată și, corespunzător, ieșirea este $y'[n] = s_p y[n]$, în loc de $y[n]$. Deoarece zgomotul este injectat după scalare, raportul dintre puterea semnalului și cea a zgomotului în sistemul scalat este de s_p^2 ori raportul semnal/zgomot pentru sistemul nescalat din figura 5.15. Cum $s_p < 1$ atunci când este necesară scalarea, raportul semnal / zgomot la ieșirea filtrului se reduce prin scalare.

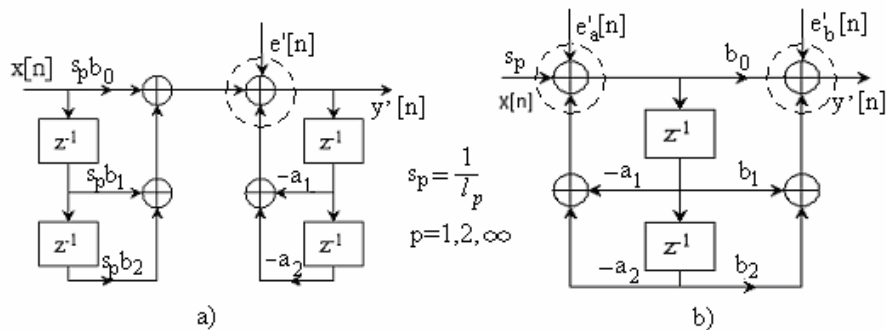


Figura 5.22. Scalarea sistemelor de ordinul doi. a) Forma directă I, b) Forma directă II

În cazul implementării în forma directă II din figura 22b factorul de scalare trebuie determinat astfel încât să se evite depășirea în ambele noduri încercuite. Funcția de sistem a filtrului scalat este $s_p H(z)$. Factorul de scalare $s_p, p = 1, \infty, 2$, contribuie cu o sursă suplimentară de zgomot la $e_a[n]$ a sistemului nescalat reprezentat în figura 5.17. Acest zgomot este filtrat în același mod de sistemul nescalat și de cel scalat. Prin urmare, puterea semnalului se multiplică cu s_p^2 , iar puterea zgomotului de ieșire este dată de relația (5.136), cu N înlocuit cu $(N+1)$, astfel încât raportul semnal/zgomot se reduce și în acest caz, dacă se efectuează scalarea pentru a evita depășirea.

În concluzie, cu cât o regulă de scalare conduce la un factor de scalare mai scăzut, se reduce probabilitatea depășirii, dar se reduce și raportul semnal/zgomot de cuantizare. Acest fapt reprezintă *interacțiunea dintre domeniul dinamic și zgomot*. Din acest motiv prezintă interes găsirea unor structuri caracterizate de zgomot de cuantizare minim în condiții de scalare precizate. Utilizarea unor structuri în formă directă de ordin mare nu conduce la rezultate satisfăcătoare din acest punct de vedere, astfel încât sunt preferate structurile în cascadă sau în paralel, realizate cu secțiuni de ordinul doi.

În continuare sunt date schemele de scalare pentru structurile în cascadă și în paralel.

5.7.3. Scalarea în realizarea în cascadă și în paralel

5.7.3.1. Analiza realizării în cascadă

În figura 5.23 este prezentat un sistem implementat cu K module de ordinul doi, fiecare din acestea implementat în forma canonică, conectate în cascadă.

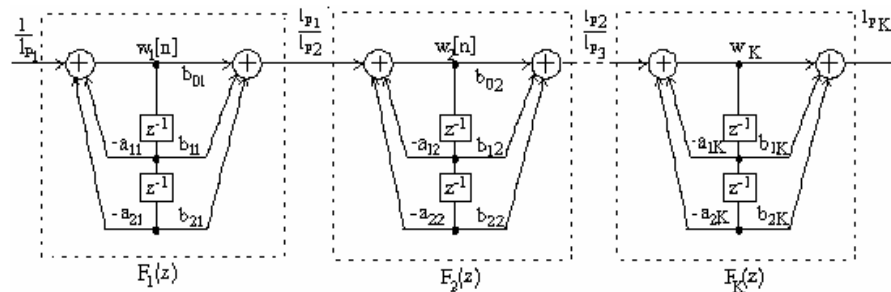


Figura 5.23. Scalarea la realizarea în cascadă a unui filtru cu K celule de ordinul doi

Se notează cu F_k , $k=1, \dots, K$, funcția de sistem a unui modul de ordinul doi.

$$F_k(z) = \frac{b_{0k} + b_{1k}z^{-1} + b_{2k}z^{-2}}{1 + a_{1k}z^{-1} + a_{2k}z^{-2}} \quad (5.162)$$

$l_{pi} = \|H_i(\omega)\|_p$; $i=1, 2, 3, \dots, K$, $p = 1, 2, \infty$, reprezintă norma după care s-a efectuat scalarea, iar $H_i(\omega)$ - funcția de transfer de la intrare la nodul w_i .

$$H_i(z) = \frac{\prod_{k=1}^{i-1} F_k(z)}{1 + a_{1i}z^{-1} + a_{2i}z^{-2}} \quad (5.163)$$

coeficienții $\frac{l_{pi}}{l_{p(i+1)}}$ pot fi încorporați în b_{0i}, b_{1i}, b_{2i} .

Ținând seama de cele prezentate în paragraful precedent, scalarea este propriu-zis necesară numai pentru secțiunile pentru care normele $l_{pi} = \|H_i(\omega)\|_p$ sunt supraunitare. Dacă, însă, $\|H_i(\omega)\|_p \leq 1$, rezultă că nu este necesară scalare pentru celula respectivă, ceea ce ar corespunde unui factor de scalare unitar, fără efect asupra zgomotului de cuantizare. Totuși, dacă se scalează intrarea într-o secțiune de ordinul doi cu un factor supraunitar, care va amplifica semnalul, va crește raportul semnal/zgomot, prin utilizarea eficientă a gamei dinamice a filtrului. Astfel, scalarea poate fi privită nu numai ca un mod de a evita depășirea, ci și de adaptare a nivelului semnalului la gama dinamică a filtrului.

În cazul unui filtru numeric IIR de ordin mare realizat prin conectarea în cascadă a unor structuri de ordinul doi, puterea zgomotului la ieșire depinde de modul în care polii și zerourile sunt împerecheate pentru a forma structuri de ordinul doi și de ordinea secțiunilor în cascadă. Se poate observa că pentru K secțiuni de ordin doi există $K!$ posibilități de a împerechea polii și zerourile și $K!$ posibilități de a ordona secțiunile de ordinul doi rezultate. Rezultă în total $(K!)^2$ sisteme diferite. În plus, se poate alege oricare din formele directe I sau II (sau transpusele lor) pentru implementarea secțiunilor de ordinul doi. Chiar și pentru sisteme de ordin mic problema împerecherii și ordonării nu este simplă, deoarece necesită un volum mare de calcule.

Se definește *factorul (sau câștigul) de vârf* pentru celula k cu relația

$$\rho_k = \frac{\max_{\omega} |H_k(\omega)|}{\left[\frac{1}{2\pi} \int_{-\pi}^{\pi} |H_k(\omega)|^2 d\omega \right]^{\frac{1}{2}}} \quad (5.164)$$

În ciuda dificultății găsirii unei împerecheri și ordonări optime, Jackson a arătat că o grupare optimă minimizează factorii de vârf și a găsit că se pot obține rezultate bune aplicând următoarele reguli simple [23]:

1. Polul care este cel mai apropiat de cercul de rază unitate din planul Z , trebuie împerecheat cu zeroul cel mai apropiat de el;
2. Regula 1 se aplică repetat până ce toți polii și zerourile au fost împerecheate;
3. Secțiunile de ordinul doi rezultate trebuie ordonate în funcție de apropierea polilor de cercul unitate, fie în ordinea crescătoare, fie descrescătoare a apropierii polilor de cercul unitate.

Regulile de împerechere sunt bazate pe observația că subsistemele cu câștig (factor) de vârf foarte mare sunt nedorite pentru că ele pot cauza depășiri și pot amplifica zgomotul de cuantizare. Împerechind un pol ce este apropiat de cercul unitate, cu un zerou adiacent se tinde să se reducă câștigul de vârf al secțiunii.

O motivație pentru regula 3 este aceea că pentru ca spectrul zgomotului de ieșire să nu aibă o alură ascuțită, cu un maxim puternic în apropierea unui pol ce este apropiat de cercul de unitate din planul Z , este de dorit ca acești poli să fie la începutul schemei în cascadă. Pe de altă parte, răspunsul în frecvență la ieșirea unui anumit nod implică produsul răspunsurilor în frecvență ale subsistemelor care preced nodul. Astfel, pentru a evita reducerea excesivă a nivelului de semnal în etajele anterioare ale cascadei ar trebui ca polii ce sunt apropiați de cercul unitate să fie plasați ultimii în cascadă. Se observă că problema ordonării secțiunilor depinde de o varietate de factori, cum ar fi dispersia totală a zgomotului de ieșire și forma spectrului zgomotului de ieșire. Jackson a folosit norme l_p pentru a cuantifica analiza problemei împerecherii și ordonării polilor și zerourilor și a elaborat o serie de reguli empirice pentru obținerea de rezultate satisfăcătoare, fără a evalua toate posibilitățile.

De multe ori, pentru obținerea unui zgomot cât mai mic, celulele se ordonează în sens crescător al factorului de vârf. În figura 5.24 este prezentată ordonarea secțiunilor de ordinul doi în cascadă în ordinea

crescătoare a selectivității, astfel încât celula cea mai selectivă să filtreze zgomotele provenite de la toate filtrele, atenuându-le.

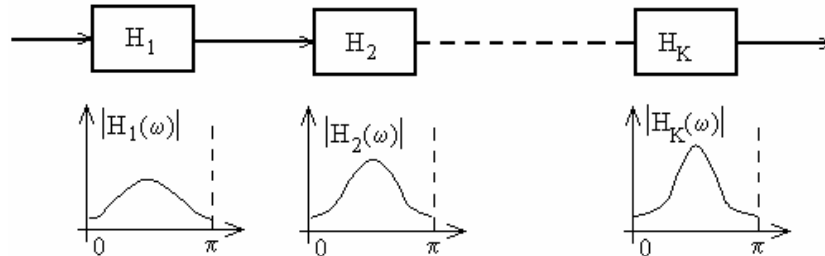


Figura 5. 24. Ordonarea secțiunilor de ordinul doi în cascadă în ordinea crescătoare a selectivității acestora

Următorul exemplu ilustrează punctul de vedere conform căruia ordonarea în cascadă a secțiunilor este importantă în controlarea zgomotului de rotunjire a produselor la ieșirea întregului sistem.

Exemplul 5.11.

Să se determine dispersia zgomotului cauzat de rotunjirea produselor, la ieșirea realizării în cascadă a filtrului causal, cu funcția de sistem

$$H(z) = H_1(z)H_2(z)$$

unde

$$H_1(z) = \frac{1}{1 - \frac{1}{2}z^{-1}}; H_2(z) = \frac{1}{1 - \frac{1}{4}z^{-1}}$$

Soluție. Fie $h[n]$, $h_1[n]$, și $h_2[n]$ răspunsurile la impuls corespunzătoare funcțiilor de transfer $H(z)$, $H_1(z)$ și, respectiv, $H_2(z)$. Acestea sunt:

$$h_1[n] = \left(\frac{1}{2}\right)^n u[n], \quad h_2[n] = \left(\frac{1}{4}\right)^n u[n], \quad h[n] = \left[2\left(\frac{1}{2}\right)^n - \left(\frac{1}{4}\right)^n\right] u[n]$$

Cele două realizări în cascadă sunt prezentate în figura 5.25.

În prima realizare în cascadă, dispersia zgomotului la ieșire este

$$\sigma_{z1}^2 = \sigma_e^2 \left[\sum_{n=0}^{\infty} h^2[n] + \sum_{n=0}^{\infty} h_2^2[n] \right]$$

În a doua realizare în cascadă, dispersia este

$$\sigma_{z_2}^2 = \sigma_e^2 \left[\sum_{n=0}^{\infty} h^2[n] + \sum_{n=0}^{\infty} h_1^2[n] \right]$$

$$\sum_{n=0}^{\infty} h_1^2[n] = \frac{1}{1 - \frac{1}{4}} = \frac{4}{3}; \quad \sum_{n=0}^{\infty} h_2^2[n] = \frac{1}{1 - \frac{1}{16}} = \frac{16}{15}$$

$$\sum_{n=0}^{\infty} h^2[n] = \frac{4}{1 - \frac{1}{4}} - \frac{4}{1 - \frac{1}{8}} + \frac{1}{1 - \frac{1}{16}} = 1,83$$

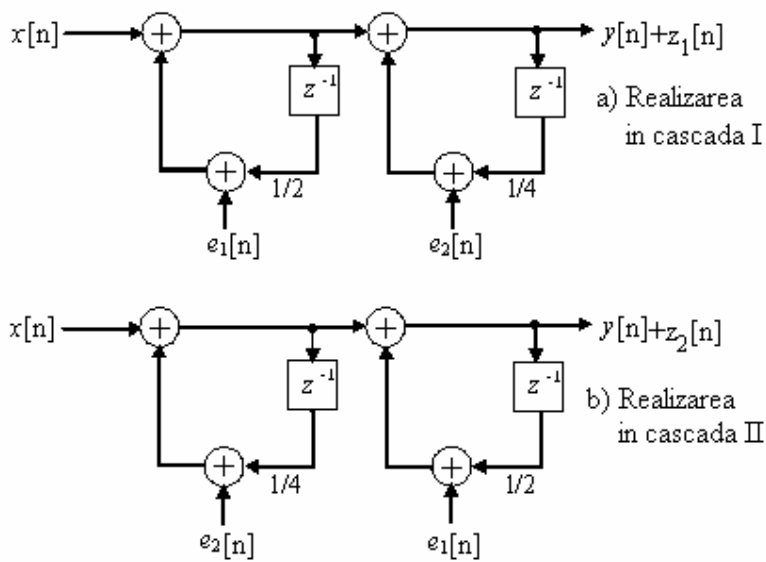


Figura 5. 25. Realizări în cascadă

În consecință,

$$\sigma_{z_1}^2 = 2,90\sigma_e^2$$

$$\sigma_{z_2}^2 = 3,16\sigma_e^2$$

iar raportul dispersiilor zgomotului de ieșire este $\frac{\sigma_{z_2}^2}{\sigma_{z_1}^2} = 1,09$.

Prin urmare, puterea zgomotului în a doua realizare în cascadă este cu 9% mai mare decât în primul caz.

5.7.3.2. Analiza realizării în paralel

În figura 5.26 este prezentat un sistem implementat cu K module de ordinul doi, conectate în paralel.

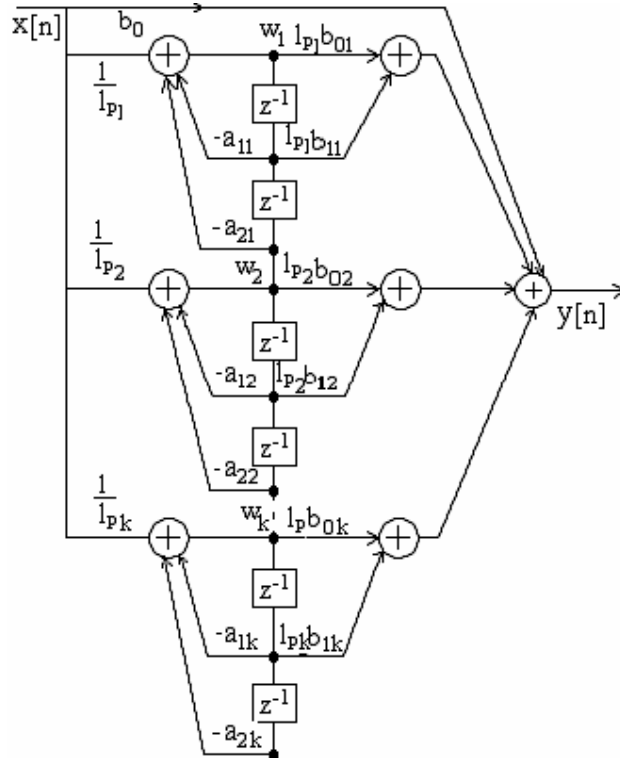


Figura 5.26. Scalarea la realizarea în paralel a unui filtru cu K celule de ordinul 2

$l_{p_i} = \|H_i(\omega)\|_p$; $i=1, 2, 3, \dots, K$, $p = 1, 2, \infty$, reprezintă norma după care s-a efectuat scalarea.

$H_i(\omega)$ - funcția de transfer de la intrare $x[n]$ la nodul w_i . Funcția de sistem corespunzătoare este

$$H_i(z) = \frac{1}{1 + a_{1i}z^{-1} + a_{2i}z^{-2}} ; i = 1, 2, 3, \dots \quad (5.165)$$

Analiza efectelor de cuantizare într-un filtru de ordin doi poate fi direct aplicată la filtrele de ordin superior bazate pe realizări în paralel. În

acest caz, fiecare secțiune de ordinul doi este independentă de celelalte secțiuni și, deci, puterea totală a zgomotului de cuantizare la ieșire este suma puterii zgomotului de cuantizare a fiecărei secțiuni individuale. Tehnicile de împerechere enunțate anterior pot fi aplicate și la formele în paralel unde se poate arăta [23] că puterea de zgomot la ieșire este comparabilă cu cele mai bune împerecheri și ordonări la conectarea în cascadă. Forma în cascadă rămâne totuși cea mai folosită pentru structurile IIR.

Deoarece structurile IIR cu formele directe I și II includ și sistemele FIR în forma directă ca un caz particular, rezultatele și tehnicile de analiză considerate mai sus se aplică la sistemele FIR, dacă se elimină toate referirile la polii funcției de sistem și se elimină căile de reacție în toate grafurile de semnal.

Pentru sistemele FIR cu fază liniară, implementarea se poate face cu aproximativ jumătate din multiplicările sistemului FIR general, ceea ce determină reducerea la jumătate a dispersiei zgomotului la ieșire, dacă produsele sunt cuantizate înainte de adunare.

Rezultatele pentru realizările în cascadă de tip IIR sunt aplicabile și pentru realizările în cascadă de tip FIR, pentru acestea urmărindu-se numai problema ordonării secțiunilor de ordinul doi.

5.7.4. Analiza erorii de cuantizare în cazul scalării intrării

Pentru a obține o imagine mai clară a efectului erorii de cuantizare, se va considera și efectul scalării intrării. Se reia cazul filtrului cu un singur pol din exemplul 5.7 prezentat în figura 5.12. Se presupune că secvența de intrare $\{x[n]\}$ este o secvență de zgomot alb, a cărei amplitudine a fost scalată cu norma l_1 pentru a preveni depășirea la adunare. Atunci

$$|y[n]| \leq x_{\max} \sum_{n=0}^{\infty} h[n]$$

Cum se dorește ca $|y[n]| \leq 1$, rezultă

$$x_{\max} \leq \frac{1}{\sum_{n=0}^{\infty} h[n]} = 1 - |a| \quad (5.166)$$

Dacă se presupune $x[n]$ uniform distribuit în domeniul $(-x_{\max}, x_{\max})$, atunci, dispersia semnalului de intrare este $\sigma_x^2 = (1 - |a|)^2/3$.

Potrivit relației (5.125), puterea zgomotului la ieșirea filtrului este

$$\sigma_z^2 = \frac{\sigma_e^2}{1-a^2}.$$

Puterea semnalului de la ieșirea filtrului este

$$\sigma_y^2 = \sigma_x^2 \sum_{k=0}^{\infty} a^{2k} = \frac{\sigma_x^2}{1-a^2} \quad (5.167)$$

Raportul dintre puterea semnalului de ieșire, σ_y^2 , și puterea erorii de cuantizare, σ_z^2 , este

$$\frac{\sigma_y^2}{\sigma_z^2} = \frac{\sigma_x^2}{\sigma_e^2} = (1-|a|)^2 \cdot 2^{2(b+1)} \quad (5.168)$$

Această expresie pentru raportul semnal/zgomot de la ieșirea filtrului arată prețul plătit ca urmare a scalării intrării, mai ales când polul este apropiat de cercul unitate.

Prin comparație, dacă intrarea nu este scalată și sumatorul are un număr suficient de mare de biți pentru a evita depășirea, amplitudinea semnalului este în intervalul (-1, 1). În acest caz, dispersia semnalului de intrare este $\sigma_x^2 = 1/3$, independentă de poziția polului. Atunci

$$\frac{\sigma_y^2}{\sigma_z^2} = 2^{2(b+1)} \quad (5.169)$$

Diferența dintre rapoartele semnal/zgomot din (5.168) și (5.169) demonstrează necesitatea de a utiliza mai mulți biți la adunare, față de multiplicare. Numărul biților adiționali depinde de poziția polului și trebuie crescut odată cu mutarea polului mai aproape de cercul unitate.

În continuare, se consideră un filtru cu doi poli care, cu precizie infinită, este descris de ecuația liniară cu diferențe

$$y[n] = a_1 y[n-1] + a_2 y[n-2] + x[n] \quad (5.170)$$

unde $a_1 = 2r \cos \theta$ și $a_2 = -r^2$.

Când cele două produse sunt rotunjite, rezultă un sistem care este descris de ecuația neliniară cu diferențe

$$v[n] = Q_r[a_1 v[n-1]] + Q_r[a_2 v[n-2]] + x[n] \quad (5.171)$$

Sistemul este prezentat în schema bloc din figura 5.27.

Fiind două multiplicări, se produc două erori de cuantizare pentru fiecare ieșire.

Prin urmare, trebuie să se introducă două secvențe de zgomot $e_1[n]$ și $e_2[n]$, care corespund ieșirilor cuantizoarelor

$$\begin{aligned} Q_r[a_1 v[n-1]] &= a_1 v[n-1] + e_1[n] \\ Q_r[a_2 v[n-2]] &= a_2 v[n-2] + e_2[n] \end{aligned} \quad (5.172)$$

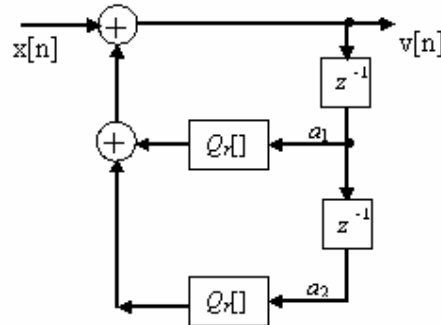


Figura 5.27 Filtru cu doi poli cu cuantizare prin rotunjire a produselor

O diagramă bloc pentru modelul corespunzător este ilustrată în figura 5.28. Se observă că secvențele de eroare $e_1[n]$ și $e_2[n]$ pot fi mutate direct la intrarea filtrului.

Ca și în cazul filtrului de ordinul întâi, ieșirea filtrului de ordin doi poate fi separată în două componente, componenta semnalului dorit și componenta erorii de cuantizare. Prima poate fi descrisă de ecuația cu diferențe

$$y[n] = a_1 y[n-1] + a_2 y[n-2] + x[n] \quad (5.173)$$

în timp ce a doua satisface ecuația cu diferențe

$$z[n] = a_1 z[n-1] + a_2 z[n-2] + e_1[n] + e_2[n] \quad (5.174)$$

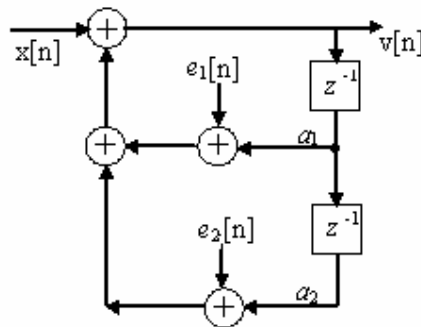


Figura 5.28 Modelul zgomotului aditiv pentru erorile de cuantizare ale unui filtru cu doi poli

Se presupune că secvențele $e_1[n]$ și $e_2[n]$ sunt necorelate.
Răspunsul la impuls al filtrului este [63]

$$h(n) = \frac{r^n}{\sin \theta} \sin(n+1)\theta \cdot u[n] \quad (5.175)$$

Prin urmare,

$$\sum_{n=0}^{\infty} h^2[n] = \frac{1+r^2}{1-r^2} \frac{1}{r^4 + 1 - 2r^2 \cos 2\theta} \quad (5.176)$$

Aplicând (5.122) se obține dispersia erorii de cuantizare la ieșirea filtrului, în forma [47]

$$\sigma_z^2 = \sigma_e^2 \left(\frac{1+r^2}{1-r^2} \frac{1}{r^4 + 1 - 2r^2 \cos 2\theta} \right) \quad (5.177)$$

Dacă semnalul de intrare $x[n]$ este scalat cu norma l_1 ca în (5.151) pentru a evita depășirea, puterea semnalului de ieșire este

$$\sigma_y^2 = \sigma_x^2 \sum_{n=0}^{\infty} h^2[n] \quad (5.178)$$

unde puterea semnalului de intrare $x[n]$ este dată de dispersia

$$\sigma_x^2 = \frac{1}{3 \left[\sum_{n=0}^{\infty} |h[n]| \right]^2} \quad (5.179)$$

În concluzie, raportul semnal/zgomot la ieșirea filtrului cu doi poli este

$$\frac{\sigma_y^2}{\sigma_z^2} = \frac{\sigma_x^2}{\sigma_e^2} = \frac{2^{2(b+1)}}{\left[\sum_{n=0}^{\infty} |h[n]| \right]^2} \quad (5.180)$$

Cu toate că este dificilă evaluarea exactă a numitorului în (5.180), este ușor să determinăm marginile superioară și inferioară ale acestuia. În particular, $|h[n]|$ este mărginită superior

$$|h[n]| \leq \frac{1}{\sin \theta} r^n \quad n \geq 0 \quad (5.181)$$

astfel încât

$$\sum_{n=0}^{\infty} |h[n]| \leq \frac{1}{\sin \theta} \sum_{n=0}^{\infty} r^n = \frac{1}{(1-r) \sin \theta} \quad (5.182)$$

Marginea inferioară se poate obține dacă se observă că

$$|H(\omega)| = \left| \sum_{n=0}^{\infty} h[n] e^{-j\omega n} \right| \leq \sum_{n=0}^{\infty} |h[n]| \quad (5.183)$$

Dar,

$$H(\omega) = \frac{1}{(1 - re^{j\theta} e^{-j\omega})(1 - re^{-j\theta} e^{-j\omega})} \quad (5.184)$$

La $\omega = \theta$, care este frecvența de rezonanță a filtrului, se obține cea mai mare valoare a lui $|H(\omega)|$, deci

$$\sum_{n=0}^{\infty} |h[n]| \geq |H(\theta)| = \frac{1}{(1-r)\sqrt{1+r^2-2r\cos 2\theta}} \quad (5.185)$$

Prin urmare, raportul semnal/zgomot este mărginit superior și inferior conform relației

$$2^{2(b+1)}(1-r)^2 \sin^2 \theta \leq \frac{\sigma_y^2}{\sigma_z^2} \leq 2^{2(b+1)}(1-r)^2(1+r^2-2r\cos 2\theta) \quad (5.186)$$

De exemplu, când $\theta = \pi/2$, expresia din (5.186) devine

$$2^{2(b+1)}(1-r)^2 \leq \frac{\sigma_y^2}{\sigma_z^2} \leq 2^{2(b+1)}(1-r)^2(1+r)^2 \quad (5.187)$$

Termenul dominant în aceste margini este $(1-r)^2$, care poate reduce serios raportul semnal/zgomot odată cu apropierea polilor de cercul unitate. Dacă $\delta = 1-r$ este distanța de la pol la cercul unitate, raportul semnal/zgomot din (5.187) este redus cu δ^2 . Aceste rezultate servesc la întărirea aserțiunii anterioare, referitoare la necesitatea utilizării mai multor biți la adunare decât la multiplicare, ca un mecanism de evitare a erorilor rezultate din operația de scalare.