

Lucrarea 1

ERORI

SCOPUL LUCRĂRII

În prima parte a lucrării se prezintă conceptele fundamentale ale reprezentării numerelor reale, utile în toate problemele de metode numerice. Într-un fel sau altul tehnicile prezentate intervin în etapele de preluare și prelucrare a datelor, chiar dacă nu în mod vizibil.

În cea de-a doua parte a lucrării sunt aduse în discuție grafurile de procedură, ca metodă de calcul a marginii erorii relative pentru diferite expresii și aplicarea acestei metode în câteva situații practice concrete.

1. PREZENTARE TEORETICĂ

1.1.ERORI ABSOLUTE ȘI ERORI RELATIVE

Deoarece valoarea măsurată sau calculată a unei mărimi nu reprezintă valoarea adevărată sau exactă a mărimii, o numim *valoare aproximativă*.

Definiția 1.1. Numim eroare diferența dintre valoarea adevărată (exactă) și valoarea aproximativă. Matematic definiția poate fi exprimată prin formula:

$$e_x = x - \bar{x}, \quad (1.1)$$

unde:

- e_x reprezintă eroarea,
- x valoarea adevărată,
- \bar{x} valoarea aproximativă.

Deoarece unii autori definesc eroarea ca diferența:

$$e_x = \bar{x} - x, \quad (1.2)$$

pentru a reuni cele două definiții ale erorii într-una singură s-a trecut la definirea *erorii absolute*, astfel:

$$e_x = |x - \bar{x}|$$

Definiția 1.2. Numim *eroare relativă* raportul dintre eroarea absolută și valoarea aproximativă absolută:

$$\varepsilon_x = |e_x/\bar{x}| \quad (1.3)$$

2. REPREZENTAREA GENERALĂ A NUMERELOR REALE

2.1.FORMA GENERALĂ A UNUI NUMĂR REAL. ENTITĂȚILE CONSTITUTIVE ALE ACESTUIA.

Un număr real se reprezintă sub forma următoare:

$$\text{nrReal} = \text{mantisa} * \text{baza}^{\text{exponent}} \quad (1.4)$$

Mantisa se mai numește și *fracție*, și este un număr real.

Exponentul este un număr întreg, deci poate fi sau nu interpretat cu semn.

Există și o reprezentare a numerelor reale de forma:

$$\text{nrReal} = \pm d_0.d_1d_2d_3\dots d_{p-1} * b^e \quad (1.5)$$

unde:

- *mantisa*, constă din pozițiile d_i , în număr de p poziții;
- *baza* este notată prin b ;
- *exponentul* este notat prin e .

Cele două notații sunt echivalente. Cea de-a doua prezintă în detaliu construcția mantisei, pe când prima notație este mai directă, și arată entitățile principale ce alcătuiesc un număr real.

2.2.CONVENȚII DE REPREZENTARE

O cifră în reprezentarea în bază 2 se numește *bit*, și va fi notată în text prin litera **b** minusculă. Entitatea formată din 8 biți consecutivi se numește *octet* (*byte* în limba engleză) și va fi referită în text prin litera **B**, majusculă.

În reprezentările din memorie ale celor două entități (în speță *mantisa* și *exponentul*), fiecare poziție binară va fi reprezentată folosind simbolurile:

- **F**, de la fracție, pentru fiecare dintre biții mantisei;
- **E**, de la exponent, pentru fiecare dintre biții exponentului.
- **S**, de la semn, pentru bitul de semn.

3. DIMENSIUNILE TIPURILOR REALE FUNDAMENTALE

3.1. SPAȚIUL DE MEMORIE

Orice tip de variabilă are sens într-un limbaj de programare dacă există în memorie. Deci dacă are asociat un *spațiu de reprezentare*. Spațiul de reprezentare, din punct de vedere *binar*, este format din *biți*. Care se grupează în *octeți*.

Spațiul minim adresabil este *octetul*. Lățimile în biți/octeți pentru fiecare tip de dată fundamental C sunt stabilite prin standard (v. standardul ISO/IEC 9899:1999).

În discuția de față este important cui atribuim / alocăm acești biți, referindu-ne la mantisă, bază și exponent.

Din punct de vedere hardware baza de reprezentare este 2 (*sistemul de numerație binar*). Celelalte baze de numerație uzual folosite în domeniul informatic sunt:

- octal (baza 8) ;
- hexazecimal (baza 16).

La ecran (și deci nu în memorie) o afișare făcută pentru un utilizator obișnuit (utilizatorii diferiți de cei specializați, precum programatori sau depanatori) va fi în sistemul de numerație zecimal (baza 10). Pentru a reuși astfel de afișări trebuie parcurse etape de *conversie*, care aduc numărul binar în echivalentul său dorit (zecimal/octal/hexazecimal).

În concluzie, pentru baza de numerație nu se va rezerva spațiu de memorie: aceasta este 2 și rămâne neschimbată (nu se poate schimba baza de numerație pe parcursul calculelor, în funcție de situații convenabile, pentru că dacă s-ar recurge la așa ceva, caracterul unitar al prelucrărilor s-ar pierde, generându-se confuzii). În plus, nu există dispozitive hardware în logică multiplă (deci care să folosească diferite baze de numerație).

Rămân de analizat mantisa și exponentul, pentru care evident există un număr de biți specifici, în funcție de tipul de dată real folosit.

3.2.ECHIVALENȚA ÎNTRE DENUMIRILE TIPURILOR REALE PREZENTE ÎN STANDARD ȘI CELE ADOPTATE DE LIMBAJELE C/C++

Standardul numește trei tipuri de date reale:

- *simplă precizie*;
- *dublă precizie*;
- *precizie extinsă*.

Limbajul C/C++ are un echivalent pentru fiecare. Aceste echivalențe sunt, respectiv:

- float
- double
- long double

Spațiul de memorie stabilit prin standardul ANSI C (ISO/IEC 9899:1999), diferențiat pentru fiecare dintre entitățile mantisă și exponent este următorul:

- float $\Rightarrow 4B = 32b$ repartizați astfel:
 - o 23b pentru mantisă;
 - o 8b exponentul, care pe poziția MSb conține semnul;
 - o 1b pentru semnul mantisei, deci semnul numărului real ca atare.
- double $\Rightarrow 8B = 64b$, repartizați astfel:
 - o 52b pentru mantisă;
 - o 11b pentru exponent;
 - o 1b pentru semnul mantisei, deci semnul numărului real.
- long double $\Rightarrow 10B = 80b$ (denumit și *ten bytes*)
 - o 64b pentru mantisă;
 - o 15b pentru exponent;
 - o 1b pentru semnul mantisei, deci semnul numărului real.

Folosind notațiile simbolice anunțate în secțiunea *Convenții de reprezentare*, vom detalia așezarea în memorie a celor 32 de biți ai tipului float.

Astfel, pentru acest tip de dată, și prin generalizare și pentru celelalte două, în memorie există următoarea așezare a biților:

$$S \overset{31}{|} \overset{30}{E} \overset{23}{E} \overset{22}{E} \overset{21}{E} \overset{20}{E} \overset{19}{E} \overset{18}{E} \overset{17}{E} \overset{16}{E} \overset{15}{E} \overset{14}{E} \overset{13}{E} \overset{12}{E} \overset{11}{E} \overset{10}{E} \overset{9}{E} \overset{8}{E} \overset{7}{E} \overset{6}{E} \overset{5}{E} \overset{4}{E} \overset{3}{E} \overset{2}{E} \overset{1}{E} \overset{0}{E}$$

Am marcat spațierea dintre fiecare grup de biți printr-o bară verticală. Deasupra pozițiilor importante am marcat indexul (puterea bazei) asociat acelei poziții. Plecând din dreapta, indexul biților crește. Primul bit din dreapta are indexul zero, și este denumit bitul cel mai puțin semnificativ (*Least Significant Bit*), iar primul bit din stânga cum privim așezarea de ansamblu a biților, poartă denumirea de bit cel mai semnificativ (*Most Significant Bit*).

Bitul LSB are utilitate de exemplu în determinarea parității numerelor reprezentate în baza 2. Bitul MSB este convențional asociat semnului, anume:

- 0, pentru numerele pozitive;
- 1, pentru numerele negative.

3.3.INTERPRETAREA GAMEI ASOCIATE FIECĂRUI TIP DE DATĂ REAL

La prima vedere, pentru tipurile reale `float` și `double` se respectă dubla precizie: tipul `double` are de două ori mai mulți biți decât `float`. Privind mai atent însă, acest lucru este îndeplinit și oarecum depășit, în sensul următor: la tipul `double`, pentru mantisă, dubla precizie ar fi cerut un număr de 46b, dacă interpretarea se face în sens strict. Însă, fiind vorba de baza 2, se știe că orice bit luat în considerare în reprezentare, conduce automat la dublarea gamei. Diferența de biți între cele două cantități este de:

$$52b - 23b = 29b$$

Aceasta este considerabilă dacă avem în vedere afirmația că fiecare bit dublează gama. Cu 29b spațiul câștigat în reprezentarea numerelor reale este:

$$2^{29} = 2^9 * 2^{20} = 256 \text{ Mega,}$$

deci 256 milioane de combinații suplimentare.

Și pentru exponent situația este mulțumitoare, pentru că nu avem strict 1b în plus (care ar duce efectiv la dublarea gamei), ci avem 3b suplimentari, așadar o operație de dublare a gamei repetată de 3 ori.

3.4.SPECIFICATORII DE FORMAT

Pentru fiecare dintre tipurile de date reale, limbajul C rezervă următoarele combinații de litere:

- `float: %f`
- `double: %lf`
- `long double: %Lf`

Deși sunt tipuri de date compatibile, recomandarea în practică este să se folosească specificatorii asociați fiecăruia dintre tipuri. În caz contrar, funcțiile de citire/scriere a datelor fiind instruite într-un anumit sens, datele de intrare sosesc în alt sens, și de aici rezultatele calculelor pot fi neașteptate, sau mai grav, nedeterminate.

4. VALORI LIMITĂ ÎN REPREZENTAREA NUMERELOR REALE

4.1. REPREZENTĂRILE PSEUDO-NUMERELOR: *zero, infinit și NaN*

După cum am anunțat în secțiunea *Convenții de reprezentare* vom folosi mai jos acele notații simbolice, pentru fiecare dintre pozițiile binare ale mantisei și exponentului, în reprezentarea unor numere speciale pentru mulțimea \mathbf{R} : zero, infinit și NaN (*Not a Number*).

Sunt date mai jos reprezentările binare ale acestor situații speciale, cu observația că s-a presupus tipul real simplă precizie (float în C), pentru a oferi expunerii un caracter comprehensiv. În cele din urmă nu spațiul este esențial, ci modul de repartitie al biților pentru fiecare entitate componentă a numărului real.

$$0 = \begin{cases} \begin{matrix} 31 & 30 & & 23 & 22 & & & & & & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \end{matrix} \\ \begin{matrix} 31 & 30 & & 23 & 22 & & & & & & 0 \\ 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \end{matrix} \end{cases}$$

$$\infty = \begin{cases} \begin{matrix} 31 & 30 & & 23 & 22 & & & & & & 0 \\ 0 & 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 & 0 \end{matrix} \\ \begin{matrix} 31 & 30 & & 23 & 22 & & & & & & 0 \\ 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 & 0 \end{matrix} \end{cases}$$

$$NaN = \begin{cases} \begin{matrix} 31 & 30 & & 23 & 22 & & & & & & 0 \\ 0 & 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 & 0 \end{matrix} \\ \dots\dots\dots \\ \begin{matrix} 31 & 30 & & 23 & 22 & & & & & & 0 \\ 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 \end{matrix} \end{cases}$$

Se pot face câteva *observații*, pe baza reprezentărilor de mai sus:

- toate cele trei reprezentări sunt obligatoriu *duale* (pozitive și negative);
- în valoare absolută, NaN îl depășește pe $+\infty$ (cum am spune 'la dreapta' lui $+\infty$), iar -NaN este ca valoare sub cea a lui $-\infty$ ('la stânga' lui $-\infty$);
- când exponentul are toate pozițiile pe 1 binar, atunci indiferent de valoarea mantisei se obține o reprezentare rezervată: fie *infinit*, fie *NaN*; aceasta demonstrează că există o gamă întreagă pe care standardul o rezervă acestor reprezentări (datorată valorilor mantisei).

Relativ la ultima observație, merită să reprezentăm exponentul pentru cele două situații, și să-l reprezentăm în baza 10. Valorile pe care le vom descoperi dau limitele superioară și inferioară ale exponentului numerelor reale, și sunt importante în găsirea valorilor limită inferioară și superioară.

În paragraful următor detaliem acest lucru.

4.1.1. INTERPRETAREA EXPONENTULUI AVÂND TOATE POZIȚIILE PE VALOAREA 1 BINAR

Reprezentarea pe 8b cu toate valorile 1 binar poate avea o semnificație dublă în baza 10, în funcție de cum considerăm util bitul de semn (MSb-ul reprezentării binare).

Fără semn, avem utile în reprezentare toate pozițiile, deci 8b. Aceasta conduce către valoarea zecimală 255. Bitul LSb este pe 1, deci numărul binar este impar.

Cu semn, numărul de biți utili este 7. Semnul este dat de MSb, și cum el este 1 numărul este negativ. Pentru cele 7 poziții, cum LSb este 1 rezultă că numărul obținut trebuie să fie impar. Obținem: -127.

Aceste valori trebuie reținute, având în vedere discuția de la reprezentarea normalizată. Acolo *polarizarea* se alege astfel încât să aducă în zero cel mai mic exponent negativ. Vedeți secțiunea *Reprezentarea normalizată* pentru mai multe detalii.

4.2. GAMA REZERVATĂ PSEUDO-NUMERELOR

Se poate observa că atunci când exponentul are pentru toți biții valoarea 1, indiferent de valorile biților mantisei, compilatorul va trata acea reprezentare drept o excepție: sau infinit, sau NaN.

Aceasta arată că, pentru cazul în care exponentul are toți biții 1, pentru tipul `float`, cu cele 23 de poziții ale mantisei se obțin 2^{23} combinații care nu pot fi utilizate în program. Apare un 'decupaj' din gama totală pentru numerele reale, decupaj care se referă exclusiv la situații imposibile de folosit în practică.

Pornind de la această situație, gama de reprezentare rămasă nu prezintă dezavantaj în sensul limitării posibilităților de reprezentare.

4.3. EXEMPLE DE SITUAȚII PRACTICE ÎN CARE APARE FIECARE DINTRE REPREZENTĂRILE SPECIALE

Zero este un număr uzual. Putem spune că poate apărea în mod normal într-o situație practică, fără să constituie o situație excepțională. Poate fi folosit fără restricții. Trebuie însă avute în vedere sunt situațiile în care acest număr apare la numitorul unei fracții. Dacă aceste cazuri nu sunt tratate corespunzător, rezultatul calculelor nu mai poate fi prezis cu acuratețe, dat fiind contextul variabil în care o asemenea situație poate apărea.

Infinit este deja o reprezentare ce nu poate fi folosită în mod normal într-un caz practic. Dacă am încerca să stabilim manual biții unei variabile, astfel încât ei să corespundă celor prezentați, în momentul în care încercăm să

folosim în mod obișnuit acea variabilă, rezultatul va fi semnalat de către compilator.

NaN este o reprezentare pur convențională a situațiilor speciale. Este într-adevăr un pseudo-număr. De exemplu, situația unei nedeterminări de tip $0/0$ este modelată în standard prin NaN. Sau alte situații similare: infinit la zero, zero la zero etc.

De exemplu, se poate forța afișarea la ecran a acestei reprezentări din partea compilatorului, prin simularea unei limite care generează o nedeterminare:

$$\lim_{x \rightarrow 0} \left(\frac{\sin(x)}{x} \right) = NaN$$

sau prin calculul unei fracții ai cărui numitor și numărător sunt ambii stabiliți la zero.

Pentru detalii consultați desfășurarea lucrării, secțiunea *Probleme de laborator*, exercițiul 6.

5. REPREZENTAREA NUMERELOR REALE ÎN FORMA NORMALIZATĂ

5.1. DENUMIREA DE NOTAȚIE ȘTIINȚIFICĂ

Numerele reale se pot reprezenta în două variante:

- normalizat;
- nenormalizat.

Utilitatea majoră o are *reprezentarea normalizată*. În această notație numerele foarte mari ca și cele foarte mici nu-și pierd semnificația și valoarea, și pot fi interpretate ca ordin de mărime. De exemplu, pentru a face distincția dintre cele două moduri de reprezentare anunțate, să alegem numărul următor:

$0.0000000001|_{10}$

După cum se poate ușor observa, acest număr este foarte mic, în speță foarte apropiat de 0. Prima cifră diferită de zero a mantisei apare pe poziția a 10^a . Aceasta înseamnă că, pentru a putea fi folosit ca o entitate diferită de zero, rezoluția de reprezentare trebuie să fie de cel puțin 10 zecimale. Ca o consecință a acestui fapt, ar trebui ca și aparatele folosite în măsurarea unei astfel de mărimi să permită rezoluții foarte bune.

Cum în mod obișnuit aplicațiile folosesc de obicei maximum 6 zecimale exacte, acest număr va fi considerat zero (urmăre operației de rotunjire, în care se păstrează doar zecimalele utile).

Aici intervine *reprezentarea normalizată*, în care numărul de mai sus nu-și va pierde valoarea, putând fi interpretat ca un număr foarte mic, dar nu zero. Scrierea sa echivalentă folosind reprezentarea normalizată este:

$$1.2 * 10^{-10}$$

dacă considerăm că este reprezentat în *baza 10*. În acest moment, numărul nu mai este zero, ba chiar mai permite încă 5 poziții zecimale, dacă am considera o reprezentare cu 6 zecimale exacte.

5.2. UTILITATEA REPREZENTĂRII NORMALIZATE

Avantajul reprezentării normalizate este acela de a permite comparația numerelor foarte mari sau foarte mici, fără a pierde din vedere pozițiile utile (cele diferite de zero) care pot apărea, pentru numerele foarte mari după multe poziții de zero în dreapta cifrei utile, iar pentru numerele foarte mici după multe poziții de zero în stânga cifrei utile.

Denumirea echivalentă a reprezentării normalizate în limbajul matematic obișnuit este *notație științifică*, sau *notație în puteri ale lui 10*. Prima denumire este cea uzuală.

Putem considera două exemple sugestive de numere, care fără apel la reprezentarea normalizată (notația științifică) nu ar avea sens în calcule.

Primul exemplu este *sarcina electronului*: $1.60217646 \times 10^{-19}$ coulombs. Este binecunoscută această reprezentare din problemele de fizică, și nu numai.

Al doilea exemplu este *numărul lui Avogadro*: $NA = 6.0221415 * 10^{23}$. Acest număr este foarte utilizat în problemele de chimie.

După cum se observă, reprezentările au sens și sunt natural interpretate prin intermediul notației științifice. În concluzie, pentru o reprezentare nenormalizată, aceste numere ar fi, presupunând un număr de zecimale impus, fie zero (pentru sarcina electronului) fie o valoare limită superioară conform gamei de reprezentare aleasă, pentru numărul lui Avogadro.

5.3. REGULILE DE REPREZENTARE IMPUSE DE FORMA NORMALIZATĂ

Pentru a putea folosi reprezentarea normalizată, trebuie să respectăm următoarele reguli, referitoare la mantisă și la exponent. Aceste reguli nu țin de baza de reprezentare.

Regulile prezentate mai jos trebuie respectate în ordinea menționată, deoarece schimbări ale valorilor mantisei influențează exponentul, și nu vice-versa.

Așadar:

- 1) *Regula pentru mantisă* - cere ca aceasta să fie situată în intervalul: [1, baza)

- 2) *Regula pentru exponent* - cere ca exponentul să fie reprezentat fără semn **în memorie**. Pentru a putea fi reprezentat fără semn, o valoare de exponent întâlnită în mod obișnuit în calcule va trebui convertită într-una fără semn folosind ceea ce se numește *polarizare (bias)* în standardul IEEE754.

Polarizarea este un număr special ales pentru fiecare formă de reprezentare reală (simplă, dublă și precizie extinsă) astfel încât să aducă în 0 valoarea cea mai mică cu semn a exponentului în reprezentarea din acea bază.

Valoarea *polarizării* se calculează cu relația:

$$\text{bias} = 2^{\text{nrBitsExponent}-1} - 1 \quad (1.6)$$

Rezultă următoarele valori, pentru fiecare tip real posibil:

- *simplă precizie (float)*: $\text{bias} = 2^{8-1} - 1 = 127$
- *dublă precizie (double)*: $\text{bias} = 2^{11-1} - 1 = 1023$
- *precizie extinsă (long double)*: $\text{bias} = 2^{15-1} - 1 = 16383$

De exemplu, numărul 10 în baza 10 se scrie în reprezentarea normalizată drept:

$$+1.0 * 10^1$$

Observații:

- 1) Există două reprezentări distincte presupuse de standard, și respectate și în limbajul C:
 - a. *reprezentarea internă* a unui număr, fiind vorba despre o reprezentare **în memorie**;
 - b. *reprezentarea externă* a unui număr real, aici fiind vorba despre o reprezentare **la ecran**.
 Cele două reprezentări sunt echivalente, dar se fac în mod diferit.
- 2) Polarizarea este un *număr impar*.

Vom da în următoarele două paragrafe două exemple de calcul care sunt foarte utile în înțelegerea mecanismelor care stau la baza funcționării unei reprezentări în virgulă mobilă:

- conversia unui număr real din baza 10 în baza 2;
- conversia unui număr real din baza 2 în baza 10;

5.3.1. REPREZENTAREA UNUI NUMĂR REAL BINAR ÎN FORMĂ NORMALIZATĂ PLECÂND DIN BAZA 10

Plecăm de la numărul zecimal: $-300.65 |_{10}$

Acesta prezintă și parte întreagă și zecimală pentru mantisă. Partea zecimală sau cea întreagă putea lipsi.

Trebuie să convertim cele două componente ale mantisei, succesiv.

Începem cu **conversia părții întregi**.

După cum vă reamintiți, conversia unui număr întreg din baza 10 în baza 2 se face prin împărțiri succesive la 2. Pentru fiecare pas se reține restul împărțirii, iar noul deîmpărțit se consideră câtul determinat la pasul curent. Resturile vor fi în final scrise în ordine inversă, astfel încât primul rest găsit devine LSB al reprezentării binare a numărului. Dacă numărul de resturi nu este o putere a lui 2, se vor face completări cu poziții de 0 înspre MSb pentru atâtea poziții lipsă până se ajunge la un număr par, putere a lui 2 (4, 8, 16 etc).

Așadar, pentru $300|_{10}$ avem succesiunea:

```

300   |2
rest 0 150|2
      rest 0 75|2
            rest 1 37|2
                  rest 1 18|2
                        rest 0 9|2
                              rest 1 4|2
                                    rest 0 2|2
                                          rest 0 1

```

După cum spuneam, primul rest devine LSB, deci este prima poziție din dreapta, în reprezentarea binară către care ne îndreptăm. Ultimul rest devine primul bit din reprezentarea binară căutată. Avem:

1 0010 1100

deci o reprezentare pe 9 poziții binare. În mod normal, la stânga primului bit trebuie să completăm cu 0 până găsim un număr de poziții putere a lui 2. Pentru situația concretă prezentată, următorul număr de poziții putere a lui 2, ce poate reprezenta dimensiunea unui tip de dată, este 16 (pentru int), deci ar trebui să completăm cu poziții de 0, astfel:

0000 0001 0010 1100

unde biții zero scriși înclinat reprezintă adăugările noastre.

Aceasta este operația completă de conversie a unui număr întreg din bază 10 în bază 2.

Pentru situația de față însă **nu vom face completările cu 0**, pentru că aceasta generează biți suplimentari, care nu rezultă din conversia efectivă, și acest lucru are efecte nedorite în 'economia' biților pentru reprezentările reale (vă reamintiți că fiecare număr real are un număr fix, impus, de poziții binare alocate mantisei). Deci, pentru lucrarea de față reprezentarea binară considerată nu va conține și adăugările de poziții 0. Așadar, reprezentarea binară finală pentru partea întreagă a mantisei este:

1 0010 1100 |₂

Avem până acum 9 biți din cei 23b ai mantisei (reamintim că exemplele considerate sunt pentru reprezentări simplă precizie, adică float în limbajul C).

Urmează **conversia părții fracționare**. Pentru aceasta reamintim că algoritmul de conversie este următorul: numărul fracționar inițial se înmulțește cu 2. Rezultă un nou număr fracționar. Se reține partea întreagă, iar pentru continuarea calculelor se preia partea fracționară, care se înmulțește din nou cu 2. Rezultă un nou număr fracționar, pentru care se reține partea întreagă, iar partea fracționară intră în operația de înmulțire cu 2 ș.a.m.d.

Avem, succesiv:

$0.65 \times 2 = 1.3$, rețin 1; continui cu 0.3
 $0.3 \times 2 = 0.6$, rețin 0, continui cu 0.6
 $0.6 \times 2 = 1.2$, rețin 1, continui cu 0.2
 $0.2 \times 2 = 0.4$, rețin 0, continui cu 0.4
 $0.4 \times 2 = 0.8$, rețin 0, continui cu 0.8
 $0.8 \times 2 = 1.6$, rețin 1, continui cu 0.6
 ș.a.m.d

Continuăm înmulțirile fie până când întâlnim o parte fracționară zero, fie până numărul de poziții reținute ne este suficient în cazul concret de lucru.

Pentru situația de față, se observă că nu se poate ca în urma unei înmulțiri cu 2 partea fracționară să ajungă 0. Astfel, vom continua înmulțirile până contorizăm un număr de biți dorit (valorile reținute se observă că sunt doar de 0 și de 1, deci pot fi considerate biți). Cum partea întreagă a mantisei a avut nevoie de 9b, vom continua înmulțirile până când 'colecționăm' 14b pentru partea fracționară.

Grupând cele două reprezentări de până acum (în total 23b) și separându-le prin punctul flotant, avem:

100101100.10100110011001

După cum s-a arătat în secțiunea 6.3, mantisa trebuie să aparțină intervalului $[1,2)$ dacă dorim o reprezentare normalizată. Acest lucru este aici posibil doar dacă partea întreagă a mantisei este 1. Pentru aceasta vom muta punctul mobil la stânga până imediat în dreapta bitului MSb. Noua reprezentare este:

1.0010110010100110011001

Am deplasat punctul mobil cu 8 poziții la stânga, deci va trebui să înmulțim numărul cu 2^8 , pentru ca nimic să nu se schimbe. Deci exponentul este pentru acest caz egal cu 8. Avem:

1.0010110010100110011001 * 2^8

Pentru că am aflat care este valoarea exponentului, și ținând cont că acum facem o reprezentare în memorie, va trebui să polarizăm exponentul. Ținând cont că reprezentăm pe gama float, polarizarea este 127. Plecând de la aceste observații, în memorie va fi stocată pentru exponent valoarea polarizată:

8 + 127 = 135

Căutăm acum reprezentarea binară a lui 135 și aceasta va fi cea memorată.

Pentru 135 găsim următoarea reprezentare (pe baza algoritmului folosit anterior):

$$1000\ 0111|_2 = 87|_{16} = (8 \cdot 16 + 7 \cdot 16^0)|_{10} = 128 + 7|_{10} = 135|_{10}$$

Am marcat pentru fiecare pas baza de numerație în care s-au făcut reprezentările de fiecare dată. Exponentul este aici este impar și se vede că LSb este 1 binar (în acest fel am verificat suplimentar rezultatul obținut).

Numărul în bază 10 a fost anunțat ca -300.65 , adică un număr negativ. Pentru aceasta va trebui ca MSb al reprezentării binare din memorie să fie 1.

Avem toate datele pentru a putea da reprezentarea binară finală a numărului zecimal propus, pentru gama de reprezentare float. Aceasta este:

$$\begin{array}{cccccccccccccccc} & 31 & 30 & & & 23 & 22 & & & & & & & & & 0 \\ 1 & 1 & 00001111 & & & 1 & 0010110010100110011001 & & & & & & & & & \end{array}$$

Ca și în prezentarea teoretică am marcat și indicii pozițiilor binare importante a le numărului.

Observații:

- 1) Fiind vorba despre o discuție în *sistem binar*, regula pentru mantisă pentru reprezentarea normalizată impune acesteia să aparțină intervalului $[1, 2)$. Acest lucru este posibil dacă și numai dacă partea întregă este $1|_2$. Orice biți am avea pentru partea fracționară, chiar și toți de $1|_2$, tot nu se reușește atingerea unității. Aceasta pentru că în partea dreaptă a punctului zecimal sunt puterile negative ale lui 2, deci o sumă de puteri negative ale lui 2, adică a unor numere *subunitare*, care nu poate atinge valoarea $1|_{10}$.
- 2) Pentru că se lucrează în formă normalizată, deci mantisa aparține obligatoriu intervalului $[1, 2)$ se poate ca primul bit din stânga al mantisei, cel de index 22, să nu mai fie memorat. Acesta se numește **hidden bit** și este presupus implicit. Fără el mantisa nu ar mai aparține intervalului indicat. Deci se poate câștiga 1b în reprezentarea numărului, ceea ce poate fi important atunci când avem un număr inexact (reprezentat pe un număr infinit de poziții într-o anumită bază).
- 3) Ca un corolar al observației 2), pot exista numere care să aibă o reprezentare exactă într-o bază de numerație și inexactă într-o alta. Este și cazul de față. Din această cauză, orice bit câștigat în reprezentarea poate fi important.

5.3.2. REPREZENTAREA UNUI NUMĂR REAL ZECIMAL ÎN FORMĂ NORMALIZATĂ PLECÂND DIN BAZA 2

Presupunem următoarea reprezentare binară ca punct de plecare în conversia reciprocă, din bază 2 în bază 10.

```
1 10000111 0010110000000000000000 |2
```

Dacă presupunem reprezentarea normalizată, înseamnă că avem deja presupus bitul ascuns. Prin urmare, ce există în memorie ca biți ai mantisei reprezintă *partea fracționară* a mantisei:

```
0010110000000000000000
```

ce va trebui prefixată cu valoarea aceluși bit ascuns, pentru a genera reprezentarea completă a mantisei.

```
1.0010110000000000000000
```

5.4. ECHIVALENȚA ÎN LIMBAJUL C A REPREZENTĂRII NORMALIZATE

După cum ne-am propus, prezentarea se face în paralel între standard și limbajul C. Spuneam că limbajul respectă standardul. Astfel, pentru acest moment al prezentării vom explica cum se poate forța o afișare în notație normalizată (format științific) a valorilor reale folosite în calcule.

Pentru C, funcțiile de citire/afișare de la consolă (tandemul tastatură + ecran) permit și interpretarea numerelor reale în reprezentare normalizată. Vom denumi această reprezentare *notație științifică*.

Așadar, în limbaj există posibilitatea afișării datelor la ecran sub forma normalizată folosind unul dintre specificatorii de format **%e** sau **%E**. Diferența între aceștia este doar de formă, și nu de fond. Pentru prima reprezentare, la ecran va apărea litera e minusculă, iar în cea de-a doua litera E majusculă.

Exemplu:

- 1) Pentru numărul 1234.567, notația științifică se obține în C astfel:

```
float varR = 1234.567;
printf("\n Notatia stiintifica: %e", varR);
Rezultatul este: 1.234567e+3
```

- 2) Pentru numărul 123456.789 se poate afișa formatul științific și în varianta cu %E:

```
float varR = 123456.789;
printf("\n Notatia stiintifica: %E", varR);
Rezultatul este: 1.23456789E+5
```

Se poate vedea că rezultatele sunt similare, doar simbolul afișat este diferit.

6. CIFRE SEMNIFICATIVE

6.1. DEFINIȚIA UNEI CIFRE SEMNIFICATIVE

Acea cifră care nu poate lipsi (nu poate fi ignorată) din reprezentarea unui număr real poartă numele de *cifră semnificativă*. Problema este mai delicată mai ales cu cifrele aflate pe ultimele poziții, privind de la stânga la dreapta. Și în mod special pentru cifrele de zero situate în aceste poziții extreme.

6.2. CIFRELE DE ZERO SEMNIFICATIVE ÎN REPREZENTARE NORMALIZATĂ ȘI NENORMALIZATĂ

Dacă numărul este în reprezentare ne-normalizată, deci digitul d_0 este zero, atunci se poate spune că acea poziție poate lipsi, pentru că există o reprezentare în care acea cifră nu mai apare - anume *reprezentarea normalizată* - v. secțiunea *Reprezentarea numerelor în forma normalizată*.

În mod normal, se poate ignora o cifră doar dacă este într-una dintre extremitățile numărului. Și trebuie ținut obligatoriu seama de regulile implicate în această decizie.

6.2.1. CIFRA SITUATĂ LA ÎNCEPUTUL NUMĂRULUI

Astfel, în reprezentarea nenormalizată, mantisa începe sigur cu o poziție zero, în stânga punctului zecimal. Așa cum am spus în paragraful anterior, în această situație acea cifră poate fi ignorată, pentru că există o reprezentare echivalentă ce nu cuprinde și acea poziție. Informația pe care acea cifră o oferă este asupra modului în care trebuie făcută translația punctului zecimal, deci se reflectă în valoarea exponentului. Este jocul care trebuie avut în vedere de fiecare dată: translația punctului zecimal influențează mărimea numărului real, deci are influență asupra exponentului.

6.2.2. CIFRA SITUATĂ LA SFÂRȘITUL NUMĂRULUI

O cifră *zero* situată pe prima poziție din dreapta numărului se poate ignora, pentru că nu are vecin în dreapta sa.

Această operație însă nu se face fără efecte. Ignorarea unei cifre în situația aceasta are două consecințe:

- pierderea în reprezentare a unui ordin de mărime (discutăm acest aspect în contextul bazelor de numerație *poziționale*);
- ignorarea eventualelor operații anterioare, care au determinat forma curentă a numărului, și care este posibil să fi presupus deja operații de trunchiere/rotunjire. Dacă se alege înlăturarea ultimei cifre se ignoră în

mod direct aceste operații. În cele din urmă acele presupuse operații se pot repeta, dar de această dată pentru un număr mai redus de zecimale, pentru că prin îndepărtarea cifrei finale, trebuie avute în vedere regulile tandem pentru trunchiere/rotunjire - v. secțiunea *Rotunjiri obișnuite*.

Așadar, în situațiile practice în care o astfel de situație apare, nu trebuie să înlăturăm cifra situată pe ultima poziție (în citirea numărului de la stânga la dreapta, ultima poziție este prima din dreapta).

De exemplu, să considerăm situația numărului următor:

2.560

pentru care dorim înlăturarea ultimei cifre, acel 0 care nu mai are vecin în dreapta sa.

Care ar fi consecințele?

Numărul acesta este preluat ca atare de către noi, fără a-i cunoaște istoria. Acest număr ar fi putut fi în prealabil supus unor operații de rotunjire/trunchiere, în urma cărora a căpătat valoarea de față. Astfel, altcineva ar fi putut deja aplica rotunjirea la 3 zecimale, pornind însă de la numărul:

3.5597₁₀

Conform regulilor de rotunjire/trunchiere, dacă dorim 3 zecimale exacte pentru numărul 2.5597 obținem chiar 2.560. În acest caz, ignorarea din partea noastră a cifrei 0 din final ar da peste cap o operație anterioară, pentru care putem presupune că s-a avut în vedere un motiv oarecare.

De asemenea trebuie reținut din exemplul prezentat că operația de trunchiere are rol important în reprezentarea datelor. Vom vedea mai concret acest lucru când discutăm despre rotunjirile obișnuite - v. secțiunea *Rotunjiri obișnuite*.

7. ROTUNJIRI OBIȘNUITE

Se fac prin comparația față de 5 a cifrei ce urmează să fie înlăturată. În această operație apar și câteva situații la limită, în care foarte importantă este cantitatea ce trebuie trunchiată.

Primul pas într-o operație de rotunjire este *trunchierea*. Astfel, dacă numărul de zecimale din reprezentarea dată este mai mare decât cel avut în vedere în rotunjire, atunci toate cifrele din acea reprezentare începând cu cifra din dreapta celei care va trebui comparată cu 5 vor fi ignorate. Această cifră ce urmează să fie comparată cu 5 este absolut necesară rotunjirii. Odată comparația încheiată această cifră va fi la rândul său ignorată.

Tocmai această trunchiere trebuie avută în vedere în rotunjire. Nu putem ignora pur și simplu zecimalele suplimentare. Trunchierea are o influență concretă și foarte importantă asupra valorii finale a numărului de rotunjit.

Pentru a aplica rotunjirea vom studia câteva exemple concrete, care acoperă toate situațiile reale.

Numerele pe care dorim să le rotunjim sunt:

4.3244

4.3236

4.3275

4.3285

Ultimele două constituie cazuri limită. Vom vedea imediat despre ce este vorba.

Am marcat în **aldin** cifra ce urmează să fie comparată cu 5. Orice alte zecimale situate în dreapta sa vor fi trunchiate. Putem avea o parte trunchiată diferită sau egală cu 0. Pentru al doilea caz rotunjirea ține cont de o *convenție*.

Să începem cu primul număr. Trunchiem așadar toate zecimalele situate în dreapta cifrei marcată în aldin (îngroșat). Iar această cifră o comparăm cu 5. Fiind mai mică, avem de-a face cu o *rotunjire prin lipsă*, așa cum ne va arăta și eroarea de rotunjire. În urma acestui tip de rotunjire lășăm neschimbată valoarea digitului din stânga celui cu care facem comparația cu 5. De fapt înlăturăm o anumită valoare astfel încât ultima cifră a numărului să devină 0. Deoarece dorim evident rotunjirea vom ignora fără probleme acea cifră (vedeți și paragraful **6 - Cifre semnificative** pentru mai multe detalii privind renunțarea la o cifră 0 situată pe ultima poziție a unui număr real).

4.3244 -> 4.324, $err = |nrRotunjit - nrInitial| = |0.0004| -$

Am pus un semn minus în dreapta rezultatului calculat în valoare absolută pentru a reține aspectul menționat: avem de-a face cu o rotunjire *prin lipsă*. Valoarea înlăturată este dată de diferența de mai sus.

Pentru al doilea număr. Trunchiem așadar toate zecimalele situate în dreapta cifrei marcată în aldin. Comparăm cifra rămasă cu 5. În acest caz este mai mare decât 5 și vom aplica *rotunjirea prin adaos*. Valoarea digitului anterior celui intrat în operația de comparare cu 5 va crește cu o unitate. Aceasta provine din faptul ca de fapt adăugăm o anumită valoare astfel încât ultima cifră a numărului să devină 0. Cum ceea ce dorim este rotunjirea vom putea ignora fără probleme acea cifră (vedeți și paragraful **6 - Cifre semnificative** pentru mai multe detalii privind renunțarea la o cifră 0 situată pe ultima poziție a unui număr real).

4.3236 -> 4.324, $err = |nrRotunjit - nrInitial| = |0.0004| +$

Am pus un semn plus în dreapta rezultatului calculat în valoare absolută pentru a reține aspectul menționat: avem de-a face cu o rotunjire *prin adaos*. Valoarea adăugată este dată de diferența de mai sus.

Pentru ultimele două situații există o nuanță suplimentară ce trebuie avută în vedere: valoarea părții trunchiate. Anunțăm acest lucru mai sus.

Așadar, sunt două situații: valoare diferită de zero SAU valoare nulă.

Pentru *prima situație*, indiferent de paritatea cifrei anterioare celei cu care facem comparația, rotunjirea se face prin adaos.

În schimb, dacă partea trunchiată este zero, atunci ținem cont de convenția care spune că dacă paritatea cifrei anterioare celei care ajută în comparația cu 5 este impară rotunjirea se va face prin adaos. Iar în cazul unei parități pare rotunjirea se face prin lipsă.

Considerând parte trunchiată diferită de zero rezultatele finale ale rotunjirii sunt:

```
4.3275 -> 4.328, err = |0.0005|+
4.3285 -> 4.329, err=|0.0005|+
```

Această convenție mai poartă numele și de *convenția digitului impar*, marcând prin această denumire situația în care se aplică rotunjirea prin adaos.

Observații:

- 1) Aplicarea convenției digitului impar are fundamentare în experimentele de laborator. Justificarea sa este legată de eroarea de rotunjire pe ansamblu. Anume, dacă observăm că printre simbolurile reprezentării zecimale avem atâtea cifre pare câte impare aceasta înseamnă că per ansamblu, dacă ținem cont de această convenție, eroarea calculelor ce implică rotunjirile se compensează și tinde la 0. Acest lucru avantajează calculele intense. Pentru acele situații nu avem studiată o statistică (sau o distribuție oarecare) prin care să putem justifica probabilitatea de apariție a simbolurilor impare în defavoarea celor pare sau reciproc. Doar eventualele erori de calibrare ale aparatelor de măsură folosite pot determina deviații de această natură.
- 2) Pentru primele două exemple de rotunjire se observă că rezultatul final este același, deși eroarea este diferită. Așadar nu putem spune dinainte ce fel de rotunjire s-a aplicat dacă avem un număr dat.
- 3) Compilatorul de C nu respectă *convenția digitului impar*. Pentru cazurile limită amintite mai sus acesta va ține cont de paritatea-imparitatea digitului anterior cifrei comparate cu 5, doar ca tehnica utilizată pare oarecum complementara convenției prezentate, în sensul ca va rotunji *prin adaos* dacă digitul este par și *prin lipsă* pentru digit impar. Se poate verifica acest lucru rulând următorul scurt program:

```
/*Rotunjirea obisnuita - cazuri limita*/
#include<stdio.h>

int main()
{
    float nr1=4.3275, nr2=4.3285;

    /*impun reprezentare cu 3 zecimale*/
    printf("Rotunjire la limita -> nr1: %.3f, nr2: %.3f\n",nr1,nr2);

    scanf("%f", &nr1);
    return 0;
}
```

- 4) Numărul implicit de zecimale folosit în reprezentările numerelor reale în ANSI C este 6. Puteți observa acest lucru dacă se renunță la limitările de reprezentare din programul anterior (construcțiile %A.Bf)
- 5) Exemplele alese anterior pentru situațiile limită în contextul operației de rotunjire, justifică faptul că eroarea maximă de rotunjire este $5 \cdot 10^{-t}$, unde t reprezintă numărul de zecimale pe care-l are un anumit număr ulterior trunchierii și anterior operației de rotunjire. Avem aici și legătura cu paragraful următor, în care sunt studiate grafurile de procedură. Acolo, operația de majorare în calculul erorilor expresiilor este determinată unilateral chiar de valoarea $5 \cdot 10^{-t}$.

8. GRAFURI DE PROCEDURĂ

Aceste grafuri sunt *binare*, adică într-un nod intră cel mult două mărimi și iese o singură mărime. Pentru fiecare operație aritmetică marca grafului se calculează diferit. Marca grafului pentru operația de adunare este raportul între valoarea din nodul care se însumează și suma valorilor din nodurile care se însumează, procedeul fiind prezentat prin graful din fig.1.1 - **Operația de adunare**.

Se dau valorile erorilor relative ε_x și ε_y ale mărimilor care se însumează și eroarea de rotunjire ε_r .

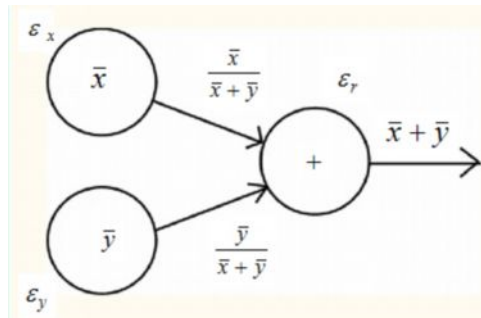


Fig.1.1. - Graful operației de adunare

Eroarea la ieșirea grafului este:

$$\varepsilon_{x+y} = \frac{\bar{x}}{\bar{x} + \bar{y}} \varepsilon_x + \frac{\bar{y}}{\bar{x} + \bar{y}} \varepsilon_y + \varepsilon_r \quad (1.7)$$

Eroarea relativă propagată prin orice operație, se calculează ca suma produselor erorilor relative ale valorilor care intră în operație cu marca grafurilor corespunzătoare, la care se mai însumează eroarea de rotunjire

datorată operației de adunare. Acest mod de calcul trebuie efectuat pe întregul graf până la ieșire, obținând în final eroarea relativă a grafului. Grafurile pot fi mai simple sau mai complexe, funcție de expresia pe care o implementează. În fiecare nod al grafului se realizează o operație aritmetică ce are doi termeni sau facto operație aritmetică care are doi termeni sau factori. Eroarea în nodul grafului se calculează după modul expus anterior și nodul poate deveni la rândul său un termen sau un factor al următoarei operații.

În calculul erorii se ține seama de valoarea obținută în fiecare nod prin operația corespunzătoare expresiei a cărei valoare a erorii relative o calculăm, pentru a determina marca ramurii grafului. Pentru operațiile de înmulțire, împărțire și ridicare la putere marca grafului este constantă indiferent de valoarea obținută în nodul grafului.

Marca grafului pentru descăzut la **operația de scădere** se calculează ca raportul între valoarea descăzutului și diferența valorilor descăzutului și scăzătorului, iar pentru scăzător se calculează ca raportul dintre valoarea scăzătorului și diferența valorilor descăzutului și scăzătorului, luat cu semnul minus. Eroarea la ieșirea grafului este dată în formula (1.8), iar graficul este prezentat în fig. 1.2.

$$\varepsilon_{x-y} = \frac{\bar{x}}{\bar{x}-\bar{y}} \varepsilon_x - \frac{\bar{y}}{\bar{x}-\bar{y}} \varepsilon_y + \varepsilon_r \quad (1.8)$$

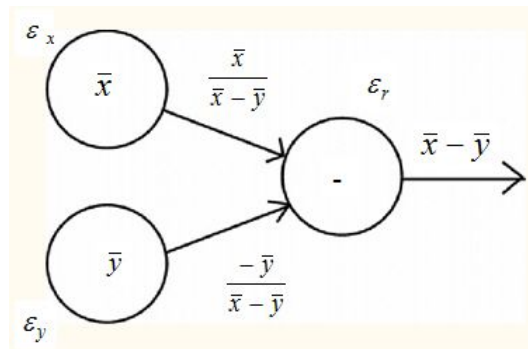


Fig.1.2. - Graficul operației de scădere

Marca grafului pentru operația de înmulțire este totdeauna constantă egală cu unitatea atât pentru înmulțitor cât și pentru deînmulțit. Graficul pentru această operație este prezentat în figura 1.3. Eroarea la ieșirea grafului este dată de expresia următoare:

$$\varepsilon_{xy} = \varepsilon_x + \varepsilon_y + \varepsilon_r \quad (1.9)$$

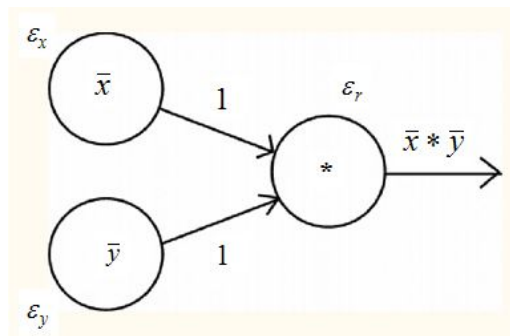


Fig.1.3. - Graful operației de înmulțire

Marca grafului pentru operația de **împărțire** este totdeauna constantă, având valoarea 1 pentru deîmpărțit și -1 pentru împărțitor. În figura 1.4 este prezentat graful pentru operația aritmetică de împărțire. Eroarea la ieșirea grafului este:

$$\varepsilon_{x/y} = \varepsilon_x - \varepsilon_y + \varepsilon_r \tag{1.10}$$

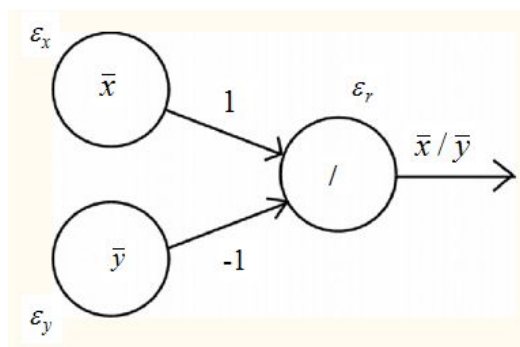


Fig.1.4. - Graful operației de împărțire

Operația de extragere a rădăcinii pătrate are marca grafului constantă egală cu $1/2$ și este prezentată în figura 1.5.

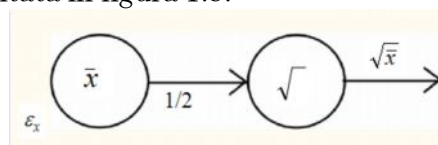


Fig.1.5. - Graful operației de extragere de rădăcină pătrată

Eroarea mărimii la ieșirea grafului este:

$$\varepsilon_{\sqrt{x}} = (1/2)\varepsilon_x + \varepsilon_r \tag{1.11}$$

Pentru ușurința calculului mărcilor grafului se ține seama de modul de calcul prezentat anterior pentru fiecare operație aritmetică. Cunoșcând erorile relative ale mărimilor care intră într-o expresie în care avem operațiile de bază ale aritmeticii, putem să determinăm eroarea finală a expresiei cu ajutorul grafurilor de procedură.

9. CONCLUZII

1. Utilizarea în programe a unui număr de zecimale trebuie să fie dictată prioritar de modul implicit de reprezentare din C (6 zecimale exacte);
2. C nu este un limbaj științific (ci unul de uz general) și din acest motiv nu ține cont de *convenția digitului impar*. Dar totuși se comportă diferit în cele două situații limită, cea în care digitul din stânga celui cu care facem comparația este par sau impar. Tehnica folosită de acesta pare oarecum opusă celei prezentate și s-ar putea numi *convenția digitului par* (dacă avem digit par atunci rotunjirea se face prin adaos);
3. Tipul `float` ar trebui să fie suficient pentru toate programele acestui laborator. Tipul `double` este tipul exclusiv al funcțiilor matematice de bibliotecă (biblioteca matematică C este construită exclusiv pe `double`). Orice argument `float` al acestor funcții va fi automat promovat la `double` de către compilator. În sfârșit, tipul de dată `long double` este utilizat în coprocesoarele matematice. Știind dimensiunea acestui tip de date înțelegem performanța obligatorie pe care coprocesoarele trebuie să o ofere.
4. Există o echivalență strânsă între impunerile standardului și felul în care ANSI C implementează aceste cerințe;
5. Reprezentările din program trebuie să aibă în vedere situația practică din acel moment. Astfel, trebuie folosiți specificatorii de format tipici fiecărei situații;
6. Nu trebuie ignorată o cifră de zero situată ultima într-o reprezentare reală oarecare. Mai general, cifra finală din reprezentarea unui număr nu se va ignora;
7. Utilizarea doar a specificatorilor de format asociați fiecărui tip de date (conform cerințelor de limbaj) și evitarea interpretării acestor tipuri de către compilator (chiar și pentru tipurile compatibile, cum sunt cele reale de exemplu, folosirea în tandem a specificatorilor diferiți generează erori de citire/interpretare/afișare a datelor).

10. DESFĂȘURAREA LUCRĂRII

10.1 . ÎNTREBĂRI

1. Unde sunt utile grafurile de procedură?

2. Cum se modelează o operație oarecare prin grafurile de procedură (de exemplu operația de adunare)? Dar o operație compusă?
3. De ce în grafurile de procedură nu este necesar să cunoaștem ordinul de mărime al erorilor operanzilor?
4. Deduceți ponderile (valorile asociate arcelor) pentru operația de scădere.
5. Care este numărul de biți rezervați exponentului pentru reprezentarea reală în simplă precizie (float în C)?
6. Care sunt cele două condiții impuse reprezentării normalizate?
7. Explicați ce semnificație are polarizarea în reprezentarea normalizată a numerelor reale.
8. Rotunjiți obișnuit numerele:
 - a. 0.30510
 - b. 0.36500
 presupunând că se cer: două și apoi trei zecimale exacte, pentru fiecare dintre numere.
9. Aduceți la forma normalizată numerele:
 - a. 809.31
 - b. 0.05925 × 10²
 - c. 72.090 × 10⁻²
 - d. 0.6115 × 10¹⁰⁰
10. De unde rezultă și ce semnificație are cantitatea 5×10^{-t} ?

10.2. PROBLEME DE LABORATOR

1. Să se determine o limită superioară a erorii relative pentru expresia:

$$E(x) = (ax^2 + bx + c) / (dx + f)$$
 în punctul x_0 știind că a, b, c, d, f, x_0 sunt pozitive și au respectiv erorile relative $\varepsilon_a, \varepsilon_b, \varepsilon_c, \varepsilon_d, \varepsilon_f, \varepsilon_{x_0}$.
2. Fie a un număr pozitiv rotunjit în mod obișnuit, iar 3 nu are eroare. Se dau expresiile:

$$u = a + a + a, \text{ și}$$

$$v = 3a$$
 Să se arate, cu ajutorul grafurilor de procedură, că marginea erorii relative a lui u este mai mare ca marginea erorii relative a lui v .
3. Se consideră expresiile $u = a^2 + 2a + 1$ și $v = (a + 1)^2$ unde a este un număr pozitiv rotunjit în mod obișnuit și presupunem că 1 și 2 nu au erori. Să se arate că pentru a foarte mare marginile erorii relative a lui u și v devin aproximativ egale iar pentru a foarte mic marginea erorii relative a lui u devine aproximativ de trei ori mai mică decât marginea

erorii relative a lui v . Numărul maxim de cifre al mantisei admise de calculator este t . Să se calculeze marginile erorii relative ale expresiilor date și pentru cazul când a este exact.

4. Fie a, b și x trei numere pozitive și exacte. Să se traseze grafurile de procedură și să se calculeze marginile erorilor relative de rotunjire pentru expresiile:

$$u = ax + bx^2, \text{ și}$$

$$v = x(a + bx).$$

Să se compare cele două margini ale erorilor, obținute. Să se calculeze marginile erorilor relative și pentru cazul când cele trei numere au erori de rotunjire.

5. Se consideră un circuit R, L, C alimentat de la o sursă U. Știind că frecvența sursei are variația ε_f , iar R, L, C au variațiile $\varepsilon_R, \varepsilon_L, \varepsilon_C$, să se traseze graful impedanței circuitului și să se determine o margine a erorii relative a ei.

6. *Exemplificări ale situațiilor speciale pentru numerele reale*

Ca exemplu s-a construit sursa următoare, care cuprinde exemple practice pentru toate situațiile în care apar numerele speciale (*pseudonumere*) la care ne-am referit.

```
// Exemple de situatii in care apar numerele reale speciale
#include<iostream>
#include<math.h>
using namespace std;

int main()
{
    union sitSpeciala {long int nrI; float nrR;} situatie;
    float numerator = 0;
    float denominator = 0;
    float underSqrt = -1;

    // zero
    situatie.nrI = 0x00000000; // zero cu plus
    cout << hex << "\n Calcul eronat... zero 1... " << situatie.nrR;

    situatie.nrI = 0x80000000; // zero cu minus
    cout << "\n Calcul eronat... zero 2... " << situatie.nrR;

    // infinit - cu plus si cu minus
    cout << "\n Calcul eronat... infinit 1... " << 1/denominator;
    cout << "\n Calcul eronat... infinit 2... " << 1/-denominator;

    // NaN - calcul eronat - nedeterminari
    cout << "\n Calcul eronat... nedeterminare 1... " <<
```



```

                                numerator/denominator;
    cout << "\n Calcul eronat... nedeterminare 2... " <<
                                sqrt(underSqrt);

// final
    cin >> numerator;
    return 0;
}

```

7. Folosirea simultană în program a tipurilor reale și întregi

Vom alege un scurt program, ce presupune utilizarea ambelor tipuri de date (reale și întregi) pentru a rezolva problema propusă.

Orice programator va trebui să folosească cu atenție tipurile mixte, pentru că este vorba despre spații de memorie în principiu diferite, dar mai ales de două mulțimi de numere, în speță \mathbf{Z} și \mathbf{R} , care nu au o echivalență imediată între valorile lor.

Va trebui utilizat și mecanismul *conversiilor explicite*, pentru a garanta păstrarea informației, și a evita pierderea acesteia - datorată de exemplu micșorării spațiului de memorie prin trecerea de la real la întreg.

Pentru programul următor, compilatorul avertizează programatorul asupra unor conversii între tipuri de date incompatibile. Aceste avertismente se referă la situațiile în care nu apare operatorul de conversie explicită. Pentru celelalte situații totul este în regulă, pentru că odată ce întâlnește conversia explicită compilatorul presupune că programatorul știe ce face. Avertismentele țin și de modul în care este configurat compilatorul, având în vedere că se pot impune opțiuni în linia de comandă în momentul lansării sale în execuție.

```

// Tipuri mixte de date
#include<iostream>
using namespace std;

int main()
{
    int nr ;                // implicit cu semn
    unsigned int nr1;      // fara semn
    float rez = -2.1;

// folosirea mixta a tipurilor de date - fara conversie implicita
// Numarul intreg este cu semn
    nr = rez;
    cout << "\n Fara conversie implicita -> numarul intreg este "
         << nr;

// folosirea mixta a tipurilor de date - cu conversie implicita
// Numarul intreg este cu semn
    nr = (int)rez;
    cout << "\n Cu conversie implicita -> numarul intreg este "
         << nr;

```

```

// =====
// folosirea mixta a tipurilor de date - cu conversie implicita
// Numarul intreg este fara semn
nr1 = rez;
cout << "\n Fara conversie implicita -> numarul intreg este "
      << nr1;

// folosirea mixta a tipurilor de date - cu conversie implicita
// Numarul intreg este fara semn
nr1 = (int)rez;
cout << "\n Cu conversie implicita -> numarul intreg este "
      << nr1;

// final program
cin >> nr;
return 0;
}

```

Compilerul va avertiza programatorul prin mesaje:

Line 14 [Warning] assignment to `int' from `float'

Line 25 [Warning] assignment to `unsigned int' from `float'

Este vorba despre liniile 14 și 25 din program, în care nu apar operatorii de conversie explicită, și se impune unui număr real să fie salvat într-un număr întreg. Astfel de situații pot ușor să fie trecute cu vederea, ceea ce are consecințe importante în momentul rulării programului.

Pentru situația în care numărului real i se impune să fie reținut în spațiul asociat unui număr întreg scurt (2B), fără a se face apel la conversia explicită, s-ar putea da următoarea interpretare operației: se pleacă de la bitul 0, se numără 16b (adică 2B) și ceea ce rezultă este salvat în numărul întreg. Însă nu este vorba despre o trunchiere fizică a biților mantisei! Rezultatul oferit de compiler în absența conversiei explicite este într-adevăr eronat, dar pentru situația în care numărul întreg este fără semn. Pentru celălalt caz, deși rezultatele sunt similare și cu și fără operatorul de conversie explicită, este indicat ca de fiecare dată să se folosească operația de conversie, pentru că nu se știe ce generație de compiler este folosit.

10.3. TEME DE CASĂ

1. Presupunem că x și y sunt două numere pozitive rotunjite simetric. Să se traseze grafurile de procedură ale expresiilor $u = (x + x + x + x) \cdot y$ și $v = 4xy$ și să se arate că marginea erorii relative a lui u este mai mare decât marginea erorii relative a lui v .
2. Se consideră expresiile $u = a^3 + 3a^2 + 3a + 1$ și $v = (a + 1)^3$ unde a este un număr pozitiv rotunjit în mod obișnuit și presupunem că 1 și 2 nu au erori. Să se calculeze marginile erorii relative a lui u și v pentru a foarte mic și pentru a foarte mare (infini). Comparați marginile

erorilor relative. Numărul de cifre semnificative a lui a este t . Toate erorile care intervin în calcul sunt mai mici decât eroarea relativă de rotunjire simetrică.

3. Se consideră circuitul din figura 1.6 în care $R_1 = 4.7k\Omega$ și toleranța $\varepsilon_1 = 10\%$, $R_2 = 6.8k\Omega$ și toleranța $\varepsilon_2 = 5\%$, $R_3 = 7.5k\Omega$ cu toleranța $\varepsilon_3 = 10\%$ și $R_4 = 10k\Omega$ cu toleranța $\varepsilon_4 = 10\%$. Să se calculeze erorile relative a rezistențelor. Să se deducă formula de calcul a rezistenței echivalente a circuitului între punctele A și B. Scrieți toate valorile în virgulă mobilă și considerând că toate erorile sunt mai mici decât eroarea relativă a rotunjirii simetrice determinați eroarea relativă maximă a rezistenței echivalente. Considerați numărul de cifre semnificative a valorilor $t = 2$.

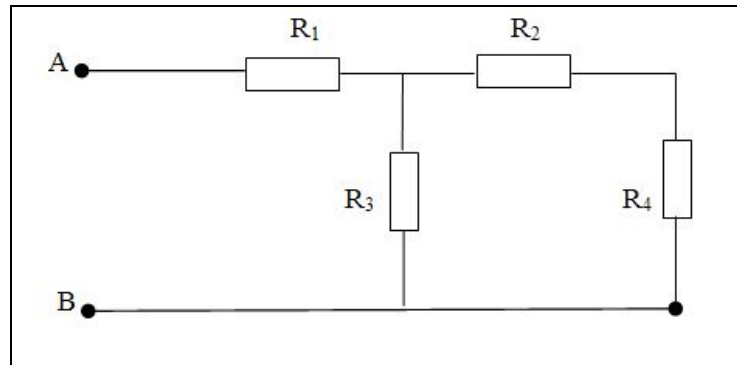


Fig. 1.6 - Circuitul pentru problema 8

4. Se consideră amplificatorul operațional din figura 1.7. Considerând rezistențele $R_1 = 100k\Omega$ cu toleranța de $\varepsilon_1 = 10\%$ și $R_2 = 10M\Omega$ cu toleranța de $\varepsilon_2 = 10\%$ să se determine eroarea maximă a amplificării. Se calculează erorile relative a rezistențelor, se fac calculele în virgulă mobilă și considerăm numărul de cifre semnificative $t = 2$.

$$A = -\frac{R_2}{R_1}$$

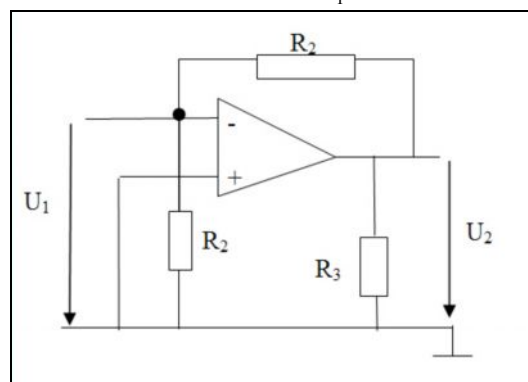


Fig. 1.7 - Circuitul pentru problema 9

BIBLIOGRAFIE

1. David Goldberg - "*What Every Computer Scientist Should Know About Floating-Point Arithmetic*", ACM Computing Surveys, Vol 23, No 1, March 1991, pp 6-48.
2. John R. Hauser - "*Handling Floating-Point Exceptions in Numeric Programs*", ACM Transactions on Programming Languages and Systems, Vol. 18, No. 2, March 1996, pp. 139-174.
3. Ioan Rusu - "*Metode Numerice în electronică - aplicații în limbaj C*", Ed. Tehnică, București, 1997
4. "*IEEE floating point representations*", material disponibil la adresa:
<http://www.math.frii.edu/~cstone/courses/fundamentals/IEEE-reals.html>
5. ***Net* - documentații din diferite surse publice.