

Αν  $\beta \geq 2$  είναι η "βάση" ενός συστήματος για την παράσταση αριθμών και  $\alpha_k = 0, 1, \dots, \beta-1$  τα αντίστοιχα "ψηφία" τότε το σύμβολο  $(\alpha_N \dots \alpha_0 \cdot \alpha_{-1} \alpha_{-2} \dots)_\beta$  παριστάει τον αριθμό  $\sum_{k=-\infty}^N \alpha_k \beta^k$ . Η σειρά ονομάζεται υποδιαστολή.

Αν  $\beta = 2, 3, 8, 10, 16, \dots$  λέμε για δυαδικό, τριαδικό, οκταδικό, δεκαδικό, δεκαεξάδικό σύστημα αριθμών. Περιγραφή αναπαράστασης (finite representation) ως  $\exists k_{min}: \alpha_k = 0, \forall k < k_{min}$

π.χ.  $(100110.11)_2 = 1 \cdot 2^5 + 0 \cdot 2^4 + 0 \cdot 2^3 + 1 \cdot 2^2 + 1 \cdot 2^1 + 0 \cdot 2^0 + 1 \cdot 2^{-1} + 1 \cdot 2^{-2} = 38.75_{10} = 38.75$

π.χ.  $53473_8 = 5 \cdot 8^4 + 3 \cdot 8^3 + 4 \cdot 8^2 + 7 \cdot 8^1 + 3 \cdot 8^0 = 22331_{10}$

ή ακόμα

$$53473_8 = 3 + 8 \cdot (7 + 8 \cdot (4 + 8 \cdot (3 + 8 \cdot 5))) = 3 + 8 \cdot 7 + 4 \cdot 8^2 + 3 \cdot 8^3 + 5 \cdot 8^4$$

= κάνουμε τις πράξεις ευκολότερα από "μέσα" προς τα "έξω"

Γενικά για τον χύεραο  $(\alpha_N \dots \alpha_0)_\beta = \alpha_0 + \beta \cdot (\alpha_1 + \beta \cdot (\alpha_2 + \beta \cdot (\dots + \beta \cdot (\alpha_{N-1} + \beta \cdot \alpha_N) \dots))_{10}$  (σχήμα Horner). Η πράξη που ακολουθεί σε μεταβλητή  $y$  με  $\alpha_i + \beta y$  γίνεται flop, επομένως κατά το σχήμα Horner χρειάζονται  $N$  flop.

π.χ.  $307.17_8 = 3 \cdot 8^2 + 0 \cdot 8^1 + 7 \cdot 8^0 + 1 \cdot 8^{-1} + 7 \cdot 8^{-2} = 199.234375_{10}$

π.χ.  $0.11_2 = 1 \cdot 2^{-1} + 1 \cdot 2^{-2} = \frac{1}{2} + \frac{1}{4} = 0.75_{10}$

Η παράσταση ενός αριθμού είναι περιοδική αν  $\alpha_k = 0$  ο αριθμός είναι ρητός, π.χ.  $4.1299\dots = 4 + 0.12 + 9 \cdot 10^{-3} + 9 \cdot 10^{-4} + \dots = 4.12 + 9 \cdot \sum_{k=0}^{\infty} 10^{-k} - 9 \cdot 10^{-0} - 9 \cdot 10^{-1} - 9 \cdot 10^{-2} = 4.12 - 0.9 - 0.09 + 9 \cdot \frac{1}{1 - \frac{1}{10}} - 9 = 3.13 + 1 = 4.13$

π.χ.  $F2B_{16} = F \cdot 16^2 + 2 \cdot 16^1 + B \cdot 16^0 = 15 \cdot 16^2 + 2 \cdot 16 + 11 = 3883_{10}$   
 10, 1, 2, 3, 4, 5, 6, 7, 8, 9, A, B, C, D, E, F

Μετατροπή ακέραιου από το δεκαδικό σύστημα στο σύστημα με βάση  $\beta$

π.χ.  $x = 369_{10} = (\dots \alpha_2 \alpha_1 \alpha_0)_8 = \alpha_0 + 8 \cdot (\alpha_1 + 8 \cdot (\alpha_2 + \dots))$

$\Rightarrow \alpha_0 = x \bmod 8$ ,  $x \div 8 = 369 \div 8 = 46 + 1$   $x = 369 = 8 \cdot 46 + 1 \Rightarrow \alpha_0 = 1$ ,  $y_0 = 46$   
 $\Rightarrow y_0 = \alpha_1 + 8 \cdot (\alpha_2 + \dots)$ ,  $\alpha_1 = y_0 \bmod 8$ ,  $y_0 = 46 = 6 + 8 \cdot 5 \Rightarrow \alpha_1 = 6$ ,  $y_1 = 5$   
 $\Rightarrow y_1 = \alpha_2 + 8 \cdot (\alpha_3 + \dots)$ ,  $\alpha_2 = y_1 \bmod 8$ ,  $y_1 = 5 = 5 + 8 \cdot 0 \Rightarrow \alpha_2 = 5$ ,  $y_2 = 0$   
 $\Rightarrow 369_{10} = 561_8$

Γενικά  $\alpha_{k+1} = y_k \bmod \beta$ ,  $k = -1, 0, \dots, N-1$ ,  $y_k = \begin{cases} x, & k = -1 \\ \frac{y_{k-1} - \alpha_k}{\beta}, & k = 0, 1, \dots, N \end{cases}$   
 είναι ο αλγόριθμος

π.χ.  $x = 111_{10}$ ,  $k = -1: y_{-1} = x$ ,  $\alpha_0 = y_{-1} \bmod 2 = 111 \bmod 2 = 1$   
 $k = 0: y_0 = \frac{y_{-1} - \alpha_0}{2} = \frac{111 - 1}{2} = 55$ ,  $\alpha_1 = 55 \bmod 2 = 1$  |  $k = 4: y_4 = \frac{y_3 - \alpha_4}{2} = \frac{6 - 0}{2} = 3$   
 $k = 1: y_1 = \frac{y_0 - \alpha_1}{2} = \frac{55 - 1}{2} = 27$ ,  $\alpha_2 = 27 \bmod 2 = 1$  |  $\alpha_5 = 3 \bmod 2 = 1$   
 $k = 2: y_2 = \frac{y_1 - \alpha_2}{2} = \frac{27 - 1}{2} = 13$ ,  $\alpha_3 = 13 \bmod 2 = 1$  |  $k = 5: y_5 = \frac{y_4 - \alpha_5}{2} = \frac{3 - 1}{2} = 1$   
 $k = 3: y_3 = \frac{y_2 - \alpha_3}{2} = \frac{13 - 1}{2} = 6$ ,  $\alpha_4 = 6 \bmod 2 = 0$  |  $\alpha_6 = 1 \bmod 2 = 1$

$k$	$y_k$	$\beta$	$y_{k+1}$	$\alpha_{k+1}$
-1	111	: 2	= 55	+ 1
0	55	: 2	= 27	+ 1
1	27	: 2	= 13	+ 1
2	13	: 2	= 6	+ 1
3	6	: 2	= 3	+ 0
4	3	: 2	= 1	+ 1
5	1	: 2	= 0	+ 1

$\Rightarrow 111_{10} = 110111_2$

ANX.  $3883_{10} \rightarrow (\dots)_{16}$

$k$	$y_k$	$\beta$	$y_{k+1}$	$\alpha_{k+1}$
-1	3883	: 16	= 242	+ 11 = B
0	242	: 16	= 15	+ 2
1	15	: 16	= 0	+ 15 = F

$\Rightarrow (F2B)_{16}$

Ένας αέρας με βάση  $\beta_1$  παραμένει αέρας αν εκφραστεί σε οποιοδήποτε άλλο σύστημα με βάση  $\beta_2$ .

-3-

Μετατροπή κλασματικού από το δεκαδικό στο σύστημα με βάση  $\beta$ .

$$\text{π.χ. } x = 0.372_{10} = (\cdot \alpha_{-1} \alpha_{-2} \alpha_{-3} \dots)_2 \Rightarrow \alpha_{-1} 2^{-1} + \alpha_{-2} 2^{-2} + \alpha_{-3} 2^{-3} + \dots$$

$$\Rightarrow 2x = 0.744 = \alpha_{-1} 2^0 + \alpha_{-2} 2^{-1} + \alpha_{-3} 2^{-2} + \dots = (\alpha_{-1} \alpha_{-2} \alpha_{-3} \dots)_2$$

$$\Rightarrow \alpha_{-1} = 0, \quad y_1 = 0.744 = (\cdot \alpha_{-2} \alpha_{-3} \dots)_2 = \alpha_{-2} 2^{-1} + \alpha_{-3} 2^{-2} + \dots$$

$$\Rightarrow 2y_1 = 1.488 = \alpha_{-2} 2^0 + \alpha_{-3} 2^{-1} + \dots = (\alpha_{-2} \alpha_{-3} \dots)_2$$

$$\Rightarrow \alpha_{-2} = 1, \quad y_2 = 0.488 = (\cdot \alpha_{-3} \alpha_{-4} \dots)_2 = \alpha_{-3} 2^{-1} + \alpha_{-4} 2^{-2} + \dots$$

$$\Rightarrow 2y_2 = 0.976 = \alpha_{-3} 2^0 + \alpha_{-4} 2^{-1} + \dots = (\alpha_{-3} \alpha_{-4} \dots)_2$$

$$\Rightarrow \alpha_{-3} = 0, \quad y_3 = 0.976 = (\cdot \alpha_{-4} \alpha_{-5} \dots)_2 = \alpha_{-4} 2^{-1} + \alpha_{-5} 2^{-2} + \dots$$

$$\Rightarrow 2y_3 = 1.952 = \alpha_{-4} 2^0 + \alpha_{-5} 2^{-1} + \dots = (\alpha_{-4} \alpha_{-5} \dots)_2$$

$$\Rightarrow \alpha_{-4} = 1, \quad y_4 = 0.952 = (\cdot \alpha_{-5} \alpha_{-6} \dots)_2$$

$$\text{Άρα } \{0.372\}_{10} = (\cdot 0101 \dots)_2$$

$$\text{Απόδειξη } \alpha_{-k-1} = [\beta y_{-k}], \quad k=0,1,\dots \quad y_{-k} = \begin{cases} x & k=0 \\ \beta y_{-k+1} - \alpha_{-k} & k=1,2,\dots \end{cases}$$

$$\text{π.χ. } k=0: y_0 = x = 0.372, \quad \alpha_{-1} = 0$$

$$k=1: y_{-1} = 2y_0 - \alpha_{-1} = 2 \cdot 0.372 = 0.744, \quad \alpha_{-2} = 1$$

$$k=2: y_{-2} = 2y_{-1} - \alpha_{-2} = 2 \cdot 0.744 - 1 = 0.488, \quad \alpha_{-3} = 0$$

$$k=3: y_{-3} = 2y_{-2} - \alpha_{-3} = 2 \cdot 0.488 = 0.976, \quad \alpha_{-4} = 1$$

$n$	$k$	$y_{-k}$	$\beta$	$y_{-k-1}$	$\alpha_{-k-1}$
	0	0.372	$\times 2$	= 0.744	+ 0
	1	0.744	$\times 2$	= 0.488	+ 1
	2	0.488	$\times 2$	= 0.976	+ 0
	3	0.976	$\times 2$	= 0.952	+ 1

n.x.  $x = 0.59375_{10} \rightarrow (\dots)_2$

k	$y_{-k}$	$\beta$	$y_{-k-1}$	$\alpha_{-k-1}$
0	$0.59375 \times 2$		$= 0.1875$	$+ 1$
1	$0.1875 \times 2$		$= 0.375$	$+ 0$
2	$0.375 \times 2$		$= 0.75$	$+ 0$
3	$0.75 \times 2$		$= 0.5$	$+ 1$
4	$0.5 \times 2$		$= 0$	$+ 1$

$\Rightarrow (0.10011)_2$

n.x.  $x = 0.5625_{10} \rightarrow (\dots)_2$

k	$y_{-k}$	$\beta$	$y_{-k-1}$	$\alpha_{-k-1}$
0	$0.5625 \times 2$		$= 0.125$	$+ 1$
1	$0.125 \times 2$		$= 0.25$	$+ 0$
2	$0.25 \times 2$		$= 0.5$	$+ 0$
3	$0.5 \times 2$		$= 0$	$+ 1$

$\Rightarrow (0.1001)_2$

n.x.  $x = 0.890625_{10} \rightarrow (\dots)_8$

k	$y_{-k}$	$\beta$	$y_{-k-1}$	$\alpha_{-k-1}$
0	$0.890625 \times 8$		$= 0.125$	$+ 7$
1	$0.125 \times 8$		$= 0$	$+ 1$

$\Rightarrow (0.71)_8$

Αν το ηλίθιο του φυσικού του x είναι ανεπαρκές σε κάποιο άσφα, τότε προφανώς ο x είναι πρώτο. Αν λοιπόν σε κάποιο άλλο άσφα το ηλίθιο του φυσικού του είναι άπειρο, αναγκαστικά θα είναι περιόδικο.

n.x.  $0.1_{10} = \frac{1}{10} = \frac{3}{2} \cdot \frac{1}{15} = \frac{3}{2} \left( \frac{16}{15} - 1 \right) = \frac{3}{2} \left( \sum_{n=0}^{\infty} \frac{1}{16^n} - 1 \right) = \frac{3}{2} \sum_{n=1}^{\infty} \frac{1}{2^{4n}} =$

$$= \sum_{n=1}^{\infty} \frac{1}{2^{4n}} + \sum_{n=2}^{\infty} \frac{1}{2^{4n+1}} = 2^{-4} + 2^{-5} + 2^{-8} + 2^{-9} + 2^{-12} + 2^{-13} + \dots =$$

$= (0.000110011001100\dots)_2$

Προφανώς, κάθε άρρητο x έχει άπειρα ηλίθια σε κάθε άσφα. (γιατί αν σε κάποιο άσφα ήταν περικόδιο ή περιπερικόδιο θα ήταν πρώτο)

Αν πάρουμε <sup>η προσέγγιση (πρόσθεσμο)</sup>  $x = (0.000110011)_2$

$$= 2^{-4} + 2^{-5} + 2^{-8} + 2^{-9} = (0.099609375)_{10}$$

$$\text{Αν πάρουμε } x = (0.0001100110011)_2 = 2^{-4} + 2^{-5} + 2^{-8} + 2^{-9} + 2^{-12} + 2^{-13}$$

$$= (0.0999908447265625)_{10}$$

Όσο πιο πολύ απλά, τόσο περισσότερο υπολογιστέ με τον (πρώτο για αναρίθμητα αριθμητικά ψηφία, πρώτο...)

ii)  $\frac{1}{2}$  андроз

k	$y_{-k}$	$\beta$	$y_{-k-1}$	$\alpha_{-k-1}$
0	0.1	$\times 2$	= 0.2	+ 0
1	0.2	$\times 2$	= 0.4	+ 0
2	0.4	$\times 2$	= 0.8	+ 0
3	0.8	$\times 2$	= 0.6	+ 1
4	0.6	$\times 2$	= 0.2	+ 1
5	0.2	$\times 2$	= 0.4	+ 0
6	0.4	$\times 2$	= 0.8	+ 0
7	0.8	$\times 2$	= 0.6	+ 1
8	0.6	$\times 2$	= 0.2	+ 1

$\rightarrow (0.000110011\dots)_2$   
не периодично

iii)  $\frac{1}{3} = 0.3_{10} \rightarrow (\dots)_2$

k	$y_k$	$\beta$	$y_{-k-1}$	$\alpha_{-k-1}$
0	0.3	$\times 2$	= 0.6	+ 0
1	0.6	$\times 2$	= 0.2	+ 1
2	0.2	$\times 2$	= 0.4	+ 0
3	0.4	$\times 2$	= 0.8	+ 0
4	0.8	$\times 2$	= 0.6	+ 1
5	0.6	$\times 2$	= 0.2	+ 1
6	0.2	$\times 2$	= 0.4	+ 0
7	0.4	$\times 2$	= 0.8	+ 0
8	0.8	$\times 2$	= 0.6	+ 1
9	0.6	$\times 2$	= 0.2	+ 1

$\rightarrow (0.0100110011\dots)_2$   
не периодично

Σημαντικά ψηφία (significant digits) λέγονται όλα τα ψηφία του αριθμού εκτός των μηδενικών που βρίσκονται στο αρχή του αριθμού

π.χ.  $x = 0.0007374$ , έχει 4 σημαντικά ψηφία με πρώτο σημαντικό το 7

π.χ.  $x = 0.0997$  έχει 3 σημαντικά ψηφία 997 με πρώτο σημαντικό το 9

π.χ.  $x = 0.099700$  έχει 5 σημαντικά ψηφία 99700 με πρώτο το 9

π.χ.  $x = 410.7$  έχει 4 σημαντικά 4107 με πρώτο το 4

π.χ.  $x = 5.70$  έχει 3 σημαντικά 570 με πρώτο το 5

π.χ.  $x = 0.0079$  έχει 2 σημαντικά 79 με πρώτο το 7

Κανονικοποιημένη αναπαράσταση (μορφή) (normalized floating point representation)  $x = \pm (0.d_1d_2\dots) \cdot \beta^E$ ,  $d_1 = 1, 2, \dots, \beta-1$ ,  $d_i = 0, 1, \dots, \beta-1 (i=2, 3, \dots)$ .

Προφανώς  $\frac{1}{\beta} \leq 0.d_1d_2\dots \leq 1$  ( $0.d_1d_2\dots = d_1\beta^{-1} + d_2\beta^{-2} + \dots = \frac{d_1}{\beta} + \frac{d_2}{\beta^2} + \dots > \frac{d_1}{\beta} \geq \frac{1}{\beta}$ ) ( $E \in \mathbb{Z}$ )

Το  $E$  λέγεται εκθέτης (exponent) ή χαρακτηριστικό (characteristic)

Ο  $0.d_1d_2\dots$  λέγεται κλάσμα ή κλάσμα ή κλασματικό μέρος (fraction, mantissa) του  $x$ .

π.χ.  $-0.00598_{10} = -0.598_{10} \times 10^{-2}$ ,  $(FBA3.DE2)_{16} = 0.FBA3DE2_{16} \times 16^4$

$111.001_2 = 0.111001_2 \times 2^3$

Η μετατροπή μιας μορφής κλασματικού αναπαράστασης από το ένα σύστημα στο άλλο γίνεται με μετα-

σύνολο αριθμών μηχανής  $M = \{ \pm 0.d_1d_2\dots d_t \cdot \beta^E$

τροπή του κλασματικού μέρους και συνθέτων  $\beta^E = 10^{E_{10}} \Leftrightarrow E \log_{10} \beta = E_{10}$ .

Σύνολο αριθμών μηχανής  $M = \{ \pm 0.d_1d_2\dots d_t \cdot \beta^E \}$ ,  $L \leq E \leq U$  ή  $\{0\}$

$t$ : ακρίβεια ή μήκος

Το ελάχιστο θετικό στοιχείο των  $M$  είναι  $0.10\dots 0 \cdot \beta^L$  ενώ το μέγιστο στοιχείο των  $M$  είναι  $(0.\beta-1, \beta-1, \dots, \beta-1) \cdot \beta^U$ .

π.χ.  $\beta=2, t=4, L=-1, U=2 \Rightarrow M = \{ 0.1000 \cdot 2^{-1} = \frac{1}{4}, \dots, 0.1111 \cdot 2^2 = 3.75 \} \cup \{0\}$

[π.χ. αν  $t=3$  και  $x = 0.000598$  και η μηχανή μπορεί να αποθηκεύσει 3 δεκαδικά ψηφία  $\Rightarrow$

Το  $M$  δεν είναι σίγουρα η προκύπτουσα και του πο/αφού-

π.χ. το  $M$  δεν είναι κλειστό αφού  $(0.1 \cdot \beta^L) (0.1 \cdot \beta^L) \notin M$

π.χ. για  $\beta=10, t=5, 1 = 0.1 \times 10^1 \in M, 10^{-5} = 0.1 \times 10^{-4} \in M$ , αλλά  $1 + 10^{-5} = 1.00001 \notin M$   
(αφού έχει 6 σημαντικά ψηφία)

Υπερχείλιση (overflow) ή αδυναμία παράσταση ενός αριθμού  $x$  με  $|x| > \max M$

Υπενχείλιση (underflow) ή αδυναμία παράσταση ενός αριθμού  $x$  με  $0 < |x| < \min M = 0.1 \times \beta^L$

⇒  $\lambda$  (αριθμός)  $\lambda = 0,000 \text{ dm}$  και ένα σημαντικό ψηφίο. Μτ αν αναφέρεται  
 για τη μήκη υποδιόδοτης  $\lambda = 0,598 \times 10^{-3} \text{ m}$  σημαίνει για 5,98 μm επί μ<sup>3</sup>,  
 από η απόδοσης και απόδοσης.

⇒ Η αναφορά είναι, προσεκτικώς, επιβεβαιώνω ιδίως σε Μ.



t=3

0.100 0.101 0.102

0.12437...  
x' x''  
0.124 0.125

Αν  $x = q\beta^E$ ,  $q = 0.d_1d_2\dots d_t d_{t+1}\dots$ , τότε οι διαδοχικοί  $x', x'' \in M$  με  $x' \leq x < x''$  είναι  $x' = 0.d_1d_2\dots d_t \cdot \beta^E$ ,  $x'' = (0.d_1d_2\dots d_t + \beta^{-t})\beta^E$  και ισχύει  $|x' - x''| = \beta^{E-t}$ .

Η προσέγγιση κάποιων  $x$  από ένα "κοντινό" και ακριβές μηχανής αριθμητική  $fl(x)$ .

Υπάρχουν 2 είδη ζήτησης προσέγγισης. Ο ένας είναι η αποκοπή, δηλαδή

αποκόπουμε τα  $d_{t+1}\dots$  και άρα  $fl(x) = x'$ , άρα  $|fl(x) - x| \leq |x' - x''| = \beta^{E-t}$

Ο άλλος είναι η στρογγυλοποίηση, όπου  $|fl(x) - x| \leq |x - y|$ ,  $\forall y \in M$ , άρα π.χ. αν  $\beta = 10$

και  $d_{t+1} \geq 5$  τότε  $fl(x) = x'' = (0.d_1d_2\dots d_t + 10^{-t})10^E$ , ενώ αν  $d_{t+1} \leq 4$  τότε

$fl(x) = x' = (0.d_1d_2\dots d_t) \cdot 10^E$ , ενώ αν  $d_{t+1} = 5$  με  $d_i = 0, i \geq t+2$  τότε

$fl(x) = x'$  ή  $x''$ . Ισχύει  $|fl(x) - x| \leq \frac{1}{2}|x' - x''| = \frac{1}{2}\beta^{E-t}$  [όσο μικρότερη  $t$ , τόσο χειρότερη προσέγγιση]

(Το σφάλμα στρογγυλοποίησης σε  $t$  δεκαδικά ψηφία είναι  $\frac{1}{2}\beta^{-t}$ )

π.χ. στρογγυλοποίηση του  $x = 2.6457513$  σε 6, 5, 4, 3, 2, 1 δεκαδικά ψηφία

Για  $t=6 \rightsquigarrow 2.645751$ , Για  $t=5 \rightsquigarrow 2.64575$ , Για  $t=4 \rightsquigarrow 2.6458$ ,

Για  $t=3 \rightsquigarrow 2.646$ , Για  $t=2 \rightsquigarrow 2.65$ , Για  $t=1 \rightsquigarrow 2.6$

(βέβαια σε μορφή κινήσεως δεκαδικών  $0.26457513 \times 10^1$  ο αριθμός των ψηφίων προσέγγισης διαφέρει)

π.χ. στρογγυλοποίηση του  $x = 3.1415926$  σε 6, 5, 4, 3, 2, 1 δεκαδικά ψηφία

Για  $t=6 \rightsquigarrow 3.141593$ , Για  $t=5 \rightsquigarrow 3.14159$ , Για  $t=4 \rightsquigarrow 3.1416$

Για  $t=3 \rightsquigarrow 3.142$ , Για  $t=2 \rightsquigarrow 3.14$ , Για  $t=1 \rightsquigarrow 3.1$

Ισχύει  $\frac{|fl(x) - x|}{x} \leq u = \begin{cases} \beta^{1-t} & \text{για αποκοπή} \\ \frac{1}{2}\beta^{1-t} & \text{για στρογγυλοποίηση} \end{cases}$

Διότι (αν  $x > 0$ ) αν  $\theta = \begin{cases} 1 & \text{για αποκοπή} \\ \frac{1}{2} & \text{για στρογγυλοποίηση} \end{cases}$  τότε  $|fl(x) - x| \leq \theta |x' - x''| \Rightarrow$

$\frac{|fl(x) - x|}{x} \leq \theta \frac{|x' - x''|}{|x|} = \theta \frac{\beta^{E-t}}{q\beta^E} \leq \theta \beta^{-t} \beta = \theta \beta^{1-t} = u$  (Διότι π.χ.  $q = 0.145 > 0.1 = 10^{-1} = \beta^{-1}$ )

Ισχύει  $fl(x) = x(1 + \epsilon)$ ,  $\epsilon \in \mathbb{R}$ ,  $|\epsilon| \leq u$  (αφού αν  $\epsilon = \frac{fl(x) - x}{x} \Rightarrow |\epsilon| = \frac{|fl(x) - x|}{x} \leq u$ )

Σφάλμα (error)  $\epsilon = x^* - x$ ,  $r = -\epsilon = x - x^*$  διόρθωση (correction)

Απόλυτο σφάλμα  $|\epsilon| = |x^* - x|$

Επιπλέον σφάλμα (relative error)  $\delta = \frac{\epsilon}{x} = \frac{x^* - x}{x} \approx \frac{x^* - x}{x^*}$  διότι  $\frac{x^* - x}{x^*} = \frac{x^* - x}{x^*} \frac{x^*}{x} = \frac{x^* - x}{x^*} (1 + \frac{\epsilon}{x}) \approx \frac{x^* - x}{x^*}$

Απόλυτο σχετικό σφάλμα  $|\delta| = \left| \frac{\epsilon}{x} \right| = \left| \frac{x^* - x}{x} \right| \approx \left| \frac{x^* - x}{x^*} \right|$  (absolute relative error)

Παράδειγμα πολλαπλών H.Y.

= Συστήματα 8200 αντιστοιχούν σε 0s 1  
 1 bit (0s) 1  
 □  
 για να προσέχουμε  
 σε x  
 31 bit  
 □□□...□  
 για το 9  
 bit, t=31

1 bit  
 □  
 για να προσέχουμε  
 σε x

7 bit  
 □□...□  
 για να προσέχουμε  
 σε x

Συνολικά 40 bits για να αντιστοιχούν σε 20 χαρακτήρες x  
 =  $5 \times (8 \text{ bits}) = 5 \text{ bytes}$   
 byte = χαρακτήρας

π.λ. αν  $x=4$ ,  $x^*=3.96$  είναι  $\varepsilon = x^* - x = 3.96 - 4 = -0.04$ ,  $r = -\varepsilon = 0.04$ ,  
 $|\varepsilon| = |-0.04| = 0.04$ ,  $\delta = \frac{\varepsilon}{x} = \frac{-0.04}{4} = -0.01$ ,  $|\delta| = |-0.01| = 0.01$   
 ή 1%

π.λ.  $x = 100\mu$ ,  $x^* = 101\mu$ ,  $|\varepsilon| = |x^* - x| = 1\mu$ ,  $\delta = \frac{\varepsilon}{x} = \frac{1}{100} = 0.01$   
 $y = 10000\mu$ ,  $y^* = 9999\mu$ ,  $|\tilde{\varepsilon}| = |y^* - y| = 1\mu$ ,  $\tilde{\delta} = \frac{\tilde{\varepsilon}}{y} = -\frac{1}{10000} = -0.0001$

Άρα το σχετικό σφάλμα είναι πιο μικρότερο.

Αν όμως ήταν  $x=0$ , τότε  $\delta \rightarrow \infty$ , άρα το  $\varepsilon$  πιο σημαντικό.

π.λ.  $x = 0.0000001$ ,  $x^* = 0.0000005 \rightarrow |\varepsilon| = 0.0000004$ ,  $\delta = \frac{\varepsilon}{x} = 4$   
 $y = 10000001$ ,  $y^* = 10000005 \rightarrow |\tilde{\varepsilon}| = 4$ ,  $\tilde{\delta} = \frac{\tilde{\varepsilon}}{y} = 0.0000004$   
 (πιο σημαντικό) (πιο σημαντικό)

Σφάλμα πρόσθεσης Το απόλυτο σφάλμα των πρόσθετων δύο αριθμών είναι μικρότερο ή ίσο με το άθροισμα των απόλυτων σφαλμάτων των αριθμών αυτών, δηλ.

$$|\varepsilon| \leq |\varepsilon_1| + |\varepsilon_2|$$

Πράγματι, αν  $\varepsilon_1 = x_1^* - x_1$ ,  $\varepsilon_2 = x_2^* - x_2$  τότε

$$\varepsilon = (x_1^* + x_2^*) - (x_1 + x_2) = (x_1^* - x_1) + (x_2^* - x_2) = \varepsilon_1 + \varepsilon_2 \Rightarrow |\varepsilon| = |\varepsilon_1 + \varepsilon_2| \leq |\varepsilon_1| + |\varepsilon_2|$$

Επίσης  $\tilde{\varepsilon} = (x_1^* - x_2^*) - (x_1 - x_2) = (x_1^* - x_1) - (x_2^* - x_2) = \varepsilon_1 - \varepsilon_2 \Rightarrow |\tilde{\varepsilon}| = |\varepsilon_1 - \varepsilon_2| \leq |\varepsilon_1| + |\varepsilon_2|$

Σφάλμα γινόμενου Το απόλυτο σχετικό σφάλμα των γινόμενων δύο αριθμών είναι μικρότερο ή ίσο από το άθροισμα των απόλυτων σχετικών σφαλμάτων των αριθμών αυτών, δηλ.  $|\delta| \leq |\delta_1| + |\delta_2|$

Πράγματι, αν  $\varepsilon_1 = x_1^* - x_1$ ,  $\varepsilon_2 = x_2^* - x_2$ , τότε  $\delta_1 = \frac{\varepsilon_1}{x_1}$ ,  $\delta_2 = \frac{\varepsilon_2}{x_2}$ , τότε

$$\varepsilon = x_1^* x_2^* - x_1 x_2 = (x_1 + \varepsilon_1)(x_2 + \varepsilon_2) - x_1 x_2 = \varepsilon_1 x_2 + \varepsilon_2 x_1 + \varepsilon_1 \varepsilon_2 \approx \varepsilon_1 x_2 + \varepsilon_2 x_1$$

$$\delta = \frac{\varepsilon}{x_1 x_2} = \frac{\varepsilon_1 x_2 + \varepsilon_2 x_1}{x_1 x_2} = \frac{\varepsilon_1}{x_1} + \frac{\varepsilon_2}{x_2} = \delta_1 + \delta_2 \Rightarrow |\delta| = |\delta_1 + \delta_2| \leq |\delta_1| + |\delta_2|$$

Επίσης  $\tilde{\varepsilon} = \frac{x_1^*}{x_2^*} - \frac{x_1}{x_2} = \frac{x_1^* x_2 - x_2^* x_1}{x_2^* x_2} = \frac{(x_1 + \varepsilon_1) x_2 - (x_2 + \varepsilon_2) x_1}{x_2^* x_2} = \frac{\varepsilon_1 x_2 - \varepsilon_2 x_1}{x_2^* x_2}$

$$\tilde{\delta} = \frac{\tilde{\varepsilon}}{\frac{x_1}{x_2}} = \tilde{\varepsilon} \frac{x_2}{x_1} = \frac{\varepsilon_1 x_2 - \varepsilon_2 x_1}{x_2^* x_2} \frac{x_2}{x_1} = \frac{\varepsilon_1 x_2 - \varepsilon_2 x_1}{x_2^* x_1} \approx \frac{\varepsilon_1 x_2 - \varepsilon_2 x_1}{x_2 x_1} = \frac{\varepsilon_1}{x_1} - \frac{\varepsilon_2}{x_2} = \delta_1 - \delta_2$$

$$\Rightarrow |\delta| = |\delta_1 - \delta_2| \leq |\delta_1| + |\delta_2|$$

π.χ.  $x = 3.1x_1 + 2.9x_2$ , όπου  $x_1, x_2$  αποθηκευμένα σε 2 δεκαδικά ψηφία

$$\varepsilon = x^* - x = (3.1x_1^* + 2.9x_2^*) - (3.1x_1 + 2.9x_2) =$$

$$= 3.1(x_1^* - x_1) + 2.9(x_2^* - x_2) = 3.1\varepsilon_1 + 2.9\varepsilon_2$$

$$|\varepsilon_1| \leq \frac{1}{2} 10^{-2}, \quad |\varepsilon_2| \leq \frac{1}{2} 10^{-2}$$

$$|\varepsilon| = |3.1\varepsilon_1 + 2.9\varepsilon_2| \leq 3.1|\varepsilon_1| + 2.9|\varepsilon_2| = 3.1 \times \frac{1}{2} 10^{-2} + 2.9 \times \frac{1}{2} 10^{-2} = 0.003$$

π.χ.  $a = 0.731$ ,  $b = 9.12^{0.912 \times 10^3}$  αποθηκευμένα σε 3 οκταδικά ψηφία τότε

$$|\varepsilon_1| \leq \frac{1}{2} 10^{-3}, \quad |\varepsilon_2| \leq \frac{1}{2} 10^{-2} \quad \left(\frac{1}{2} \beta^{E-b} = \frac{1}{2} 10^{2+3}\right)$$

Το σφάλμα  $\varepsilon$  του  $a+b$  έχει  $|\varepsilon| \leq |\varepsilon_1| + |\varepsilon_2| \leq \frac{1}{2} 10^{-3} + \frac{1}{2} 10^{-2} = 0.0055$

Το ανάλογο σχήμα σφάλμα  $|\delta|$  του  $\frac{a}{b}$  είναι

$$|\delta| \leq |\delta_1| + |\delta_2| = \frac{|\varepsilon_1|}{|a|} + \frac{|\varepsilon_2|}{|b|} \leq \frac{1}{|a|} \frac{1}{2} 10^{-3} + \frac{1}{|b|} \frac{1}{2} 10^{-2} = \frac{1}{0.731} \frac{10^{-3}}{2} + \frac{1}{9.12} \frac{10^{-2}}{2}$$

$$= 0.0012321 \leq 0.0015$$

Πράξη των υπολοίπων  $x * y \rightarrow z = fl(fl(x) * fl(y))$   
 $+, -, \dots$

π.χ.  $\beta=10, t=5, U=-L=10$ ,  $fl(\cdot)$  με αποψλοδωρισμό.

$x = 5891.26, y = 0.0773414$

$fl(x) = 0.58913 \times 10^4, fl(y) = 0.77341 \times 10^{-1} = 0.000077341 \times 10^4$

$fl(x) + fl(y) = 0.5891377341 \times 10^4 = 5891.377341$

$z = fl(fl(x) + fl(y)) = 0.58914 \times 10^4 = 5891.4$

$x + y = 5891.3373414, fl(x + y) = 0.58913 \times 10^4, fl(x) + fl(y) = \dots = 5891.3$

π.χ.  $a=1, b=3 \times 10^{-5}, c=3 \times 10^{-5}$

$a + (b + c) \rightarrow z = fl(fl(a) + fl(fl(b) + fl(c)))$

$fl(a) = 1, fl(b) = 0.3 \times 10^{-4} = fl(c) = 0.00003$

$fl(a) + fl(fl(b) + fl(c)) = 1.00006 = 0.100006 \times 10^1$

$z = 0.10001 \times 10 = 1.0001$

$(a + b) + c \rightarrow \tilde{z} = fl(fl(fl(a) + fl(b)) + fl(c))$

$fl(a) + fl(b) = 1.00003, fl(fl(a) + fl(b)) = 0.100003 \times 10^1 = 1.00003$   
 $= 0.100003 \times 10^1$

$fl(fl(a) + fl(b)) + fl(c) = 0.100003 \times 10^1 + 0.00003 \times 10^1 = 0.100003 \times 10^1$

$\tilde{z} = fl(fl(fl(a) + fl(b)) + fl(c)) = 0.100003 \times 10^1 = 1.00003$

$\Rightarrow z \neq \tilde{z}$

π.χ.  $1 + 4 \times 10^{-5} = 1.00004 = 0.100004 \times 10^1$

$fl(1 + 4 \times 10^{-5}) = 0.10000 \times 10^1 = 1.0000 = 1$

όμοια  $fl(1+x) = 1, \forall 0 < x < 5 \times 10^{-5}$  (γιατί  $\forall 0 < x < \frac{1-\beta^{1-t}}{2}$ )  
 $fl(1+x) = 1$

ήπει αυτή τα x θα ήταν το πρώτο του μέρους -

Ο μικρότερος αριθμός επίθετα μετρώμε είναι  $0.1 \times 10^{-10} = 10^{-11}$ , ήπει αυτό  
 κληρονομία από το σφαιρικό μετρώμε της μετρώμε.

Καταστροφική ακύρωση σημαντικών ψηφίων

Η αφαίρεση δύο περίπου ίσων προσεγγιστικών αριθμών οδηγεί σε μείωση των ακριβειών του αποτελέσματος.

π.χ.  $e^{-5.5} = \sum_{n=0}^{\infty} \frac{(-5.5)^n}{n!} \approx \sum_{n=0}^{24} \frac{(-5.5)^n}{n!} = 0.0026363$  προς σημαντικά ψηφία, ενώ  
 $e^{-5.5} = 0.0040868$  , δηλαδή κανένα σημαντικό ψηφίο δεν είναι σωστό

Μπορούμε να χρησιμοποιήσουμε ακριβέστερα αριθμοδύναμη σημαντικών ψηφίων που κοστίζει σε μνήμη και υπολογιστικό χρόνο.

Εναλλακτικά στο παράδειγμα μπορούμε  $e^{-a} = \frac{1}{e^a}$  όπου δεν υπάρχει καταστροφική ακύρωση.

π.χ.  $\sqrt{708} - \sqrt{707} = 26.61 - 26.59 = 0.02$  4 σημαντικά ψηφία  
[Η πραγματική τιμή  $\sqrt{708} - \sqrt{707} = 0.018797791...$  και η ακρίβεια του από το 0.02 είναι 0.001202909... , δηλαδή 2 ψηφία των κλιμακωμένων σφάλλων]

Τα σφάλματα απορροδοούνται ως αριθμούς  $\sqrt{708}, \sqrt{707}$  είναι  $\frac{1}{2} \times 10^{-2}$  , άρα το άθροιστο σφάλμα με διαφοράς τους είναι  $\frac{1}{2} \times 10^{-2} + \frac{1}{2} \times 10^{-2} = 10^{-2} = 0.01$   ~~$\frac{1}{2} \times 10^{-2}$~~

Άλλως  $\sqrt{708} - \sqrt{707} \stackrel{\text{6 σημαντικά}}{=} 26.6083 - 26.5895 = 0.0188$  (με σφάλλων από την αριθμοδύναμη της 0.0000002309 με βάση τον κλιμακωτισμό)  
 $\stackrel{\text{8 ψηφ.}}{=} 0.0187977$  (με σφάλλων με βάση τον δίστομο κλιμακωτισμό)

Άλλως  $\sqrt{708} - \sqrt{707} = \frac{(\sqrt{708} - \sqrt{707})(\sqrt{708} + \sqrt{707})}{\sqrt{708} + \sqrt{707}} = \frac{708 - 707}{\sqrt{708} + \sqrt{707}} \stackrel{\text{4 ψηφ.}}{=} \frac{1}{26.61 + 26.59}$   
 $= \frac{1}{53.20} = 0.01880$  (πολύ ακριβέστερα με βάση 4 ψηφ. ψηφία) με βάση τον κλιμακωτισμό

π.χ.  $x_{1,2} = \frac{-b \pm \sqrt{b^2 - 4ac}}{2a}$  ,  $b^2 \gg 4ac$

$x_1 = \frac{(-b + \sqrt{b^2 - 4ac})(b + \sqrt{b^2 - 4ac})}{2a(b + \sqrt{b^2 - 4ac})} = \frac{-4ac}{2a(b + \sqrt{b^2 - 4ac})} = \frac{-2c}{b + \sqrt{b^2 - 4ac}}$

π.χ.  $\beta=10, t=10, \sqrt{7892} - \sqrt{7891} = 0.8883692926 \times 10^2 - 0.8883130079 \times 10^2$   
 $= 0.5628470000 \times 10^{-2}$  σημαντική ανώτατη ακρίβεια

$\sqrt{7892} - \sqrt{7891} = \frac{1}{\sqrt{7892} + \sqrt{7891}} = 0.5628468294 \times 10^{-2}$

π.χ.  $x - \sin x$  ,  $|x|$  μικρό (αφαίρεση σχεδόν ίσων αριθμών)

Αλλά  $\sin x \approx x - \frac{x^3}{3!} + \varepsilon(x)$  ,  $|\varepsilon(x)| \leq \frac{|x|^5}{120}$

$\rightarrow x - \sin x \approx \frac{x^3}{6}$

$\sin x = \sum_{n=0}^{\infty} \frac{(-1)^n}{(2n+1)!} x^{2n+1}$   
 $= x - \frac{x^3}{3!} + \frac{x^5}{120}$   
 $n=0 \quad n=1 \quad n=2$

Σφάλμα πολλαπλασιασμού Αν  $z = fl(fl(x) fl(y))$  τότε  $\left| \frac{z - xy}{xy} \right| \leq 3u$

Πράγματι, έστω  $fl(x) = x(1 + \epsilon_1)$ ,  $fl(y) = y(1 + \epsilon_2)$ ,  $|\epsilon_i| \leq u$

$z = fl(fl(x) fl(y)) = xy(1 + \epsilon_1)(1 + \epsilon_2)(1 + \epsilon_3)$

Γοητεύω  $\prod_{i=1}^m (1 + \epsilon_i) = (1 + \epsilon)^m$ , όπου  $|\epsilon_i| \leq u < 1$ ,  $|\epsilon| \leq u$  διότι αν  $\lambda = \prod_{i=1}^m (1 + \epsilon_i)$  τότε  $(1 - u)^m = (1 - u) \dots (1 - u) \leq (1 + \epsilon_1) \dots (1 + \epsilon_m) \leq (1 + u) \dots (1 + u) = (1 + u)^m \Rightarrow (1 - u)^m \leq \lambda \leq (1 + u)^m$  • Το δεξί μέρος είναι λάθος.

αφού για τα ορισμένα ορίσματα  $(1+x)^m$  στο  $[-u, u]$  διότι αν  $\exists \epsilon \in [-u, u]: (1 + \epsilon)^m = 1$

Αρα  $z = xy(1 + \epsilon)^3$ ,  $|\epsilon| \leq u \Rightarrow \left| \frac{z - xy}{xy} \right| = |1 - (1 + \epsilon)^3| = |3\epsilon + 3\epsilon^2 + \epsilon^3| \leq 3|\epsilon| + 3|\epsilon|^2 + |\epsilon|^3 < 3|\epsilon| + 4|\epsilon|^2 \leq 3u + 4u^2$

Σφάλμα ορισμού Αν  $z = fl\left(\frac{fl(x)}{fl(y)}\right)$  τότε  $\left| \frac{z - x/y}{x/y} \right| \leq 3u$

Πράγματι,  $fl(x) = x(1 + \epsilon_1)$ ,  $fl(y) = y(1 + \epsilon_2)$  τότε

$z = fl\left(\frac{x(1 + \epsilon_1)}{y(1 + \epsilon_2)}\right) = \frac{x}{y} \frac{1 + \epsilon_1}{1 + \epsilon_2} (1 + \epsilon_3)$ ,  $|\epsilon_i| \leq u$

$= \frac{x}{y} (1 + \epsilon_1)(1 + \epsilon_3) \left(1 - \frac{\epsilon_2}{1 + \epsilon_2}\right) = \frac{x}{y} (1 + \epsilon)^2 (1 + \delta)$ ,  $|\epsilon| \leq u$

$\delta = -\frac{\epsilon_2}{1 + \epsilon_2}$

$\Rightarrow \left| \frac{z - x/y}{x/y} \right| = |1 - (1 + \epsilon)^2 (1 + \delta)| = |2\epsilon + \delta + \epsilon^2 + 2\epsilon\delta + \delta\epsilon^2|$

$\leq 2|\epsilon| + |\delta| + |\epsilon|^2 + 2|\epsilon||\delta| + |\delta||\epsilon|^2$

$\leq 2u + \frac{u}{1-u} + u^2 + \dots \leq 3u$

$|\delta| = \frac{|\epsilon_2|}{|1 + \epsilon_2|} \leq \frac{|\epsilon_2|}{1 - |\epsilon_2|}$

$\leq \frac{u}{1-u}$

Σφάλμα πρόσθεσης-αφαίρεσης Αν  $z = fl(fl(x) + fl(y))$  τότε  $\left| \frac{z - (x+y)}{x+y} \right| \leq 2u \frac{|x| + |y|}{|x+y|}$

Πράγματι, αν  $fl(x) = x(1 + \epsilon_1)$ ,  $fl(y) = y(1 + \epsilon_2)$ , τότε

$z = fl(x(1 + \epsilon_1) + y(1 + \epsilon_2)) = (x(1 + \epsilon_1) + y(1 + \epsilon_2))(1 + \epsilon_3) = (x + y + \epsilon x + \delta y)(1 + \delta)$ ,  $|\epsilon|, |\delta| \leq 2u$

$\Rightarrow z = x + y + 2(\epsilon x + \delta y) + \epsilon^2 x + \delta^2 y \approx (x + y) + 2(\epsilon x + \delta y)$

$\Rightarrow \left| \frac{z - (x+y)}{x+y} \right| \approx 2 \left| \frac{\epsilon x + \delta y}{x+y} \right| \leq 2 \frac{|\epsilon||x| + |\delta||y|}{|x+y|} \leq 2u \frac{|x| + |y|}{|x+y|}$

Αν  $x, y$  ομόσημοι, τότε  $|x| + |y| = |x+y| \Rightarrow \left| \frac{z - (x+y)}{x+y} \right| \leq 2u$

Αν  $x, y$  ετερόσημοι π.χ.  $x=4, y=-3$  τότε  $\frac{|x| + |y|}{|x+y|} = \frac{7}{1} = 7$  οπότε μεγάλως το σφάλμα αν  $x$  είναι σφίχτα. Αν  $y = -x$ , τότε το σφάλμα τριπλασιάζεται.

(Μετασχηματίζοντας εύκολα τον σφίχτα αριθμό αν αφαίρεσε δύο αριθμούς κοντά μεταξύ τους)

