

## Εισαγωγή στην Αριθμητική Ανάλυση

1. Αριθμητική κινητής υποδιαστολής  
Σφάλματα στρογγύλευσης

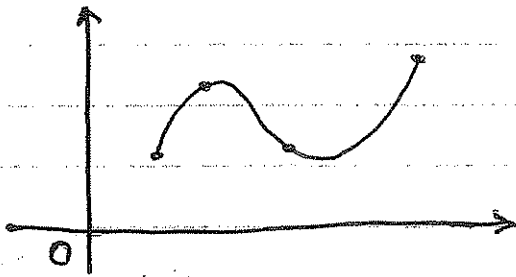
2. Μη γραμμικές εξισώσεις  
 $f: \mathbb{R} \rightarrow \mathbb{R}$

Να βρούμε ρίζες της  $f$ , δηλαδή  $x$   
τέτοιο ώστε  $f(x) = 0$

3. Γραμμικά συστήματα  
 $A \in \mathbb{R}^{n,n}$ ,  $b \in \mathbb{R}^n$

Ζητείται  $x \in \mathbb{R}^n$  τέτοιο ώστε  $Ax = b$

4. Παρεμβολή



5. Αριθμητική Ολοκλήρωση

$$\int_a^b f(x) dx = F(b) - F(a)$$

$$F' = f$$

$$\int_a^b f(x) dx \approx \sum_{i=1}^n w_i f(x_i)$$

2

# 1. Αριθμητική κινήσης υποδιαστολής Σφάλματα στρογγύλευσης

★ 1<sup>ο</sup> Παράδειγμα

$$I_n = \int_0^1 x^n e^{x-1} dx, n \in \mathbb{N}$$

•  $I_n > 0$

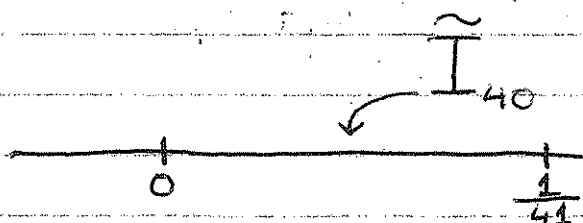
•  $I_{n+1} < I_n$ , άρα η ακολουθία είναι γνήσια φθίνουσα

•  $I_n \rightarrow 0, n \rightarrow +\infty$ , άρα μηδενική ακολουθία

•  $0 < I_n < \frac{1}{n+1}$

$$\begin{cases} I_1 = \frac{1}{e} \\ I_n = 1 - n I_{n-1}, n \geq 2 \end{cases} \quad (\text{ασταθής})$$

$$\begin{cases} I_{n-1} = \frac{1 - I_n}{n}, n = 40, 39, \dots \\ \tilde{I}_{40} = 0 \text{ ή } \frac{1}{2 \cdot 41} \end{cases} \quad (\text{ευσταθής})$$



★ 2<sup>ο</sup> Παράδειγμα

Προσέγγιση του  $\pi$  με τη μέθοδο του Αρχιμήδη

$$\begin{cases} y_n = 2^n \sin \frac{\pi}{2^n}, & n \in \mathbb{N} \\ y_1 = 2 \end{cases}$$

$(y_n)_{n \in \mathbb{N}}$  γνήσια αύξουσα

$$y_n \rightarrow \pi, \quad n \rightarrow +\infty$$

$$\begin{cases} y_{n+1} = 2^{n+1} \sqrt{\frac{1}{2} (1 - \sqrt{1 - (2^{-n} y_n)^2})}, & n = 1, 2, \dots \\ y_1 = 2 \end{cases} \quad (\text{ασταθής})$$

$$\begin{cases} y_{n+1} = \sqrt{\frac{2}{1 + \sqrt{1 - (2^{-n} y_n)^2}}} y_n, & n = 1, 2, \dots \\ y_1 = 2 \end{cases} \quad (\text{ευσταθής})$$

Ποιότητα αριθμητικής μεθόδου

1) Απαιτούμενη κνίξη.

2) Απαιτούμενο κόστος.

3) Ακρίβεια.

4

★ Ανάλυση των δύο προηγούμενων παραδειγμάτων

1<sup>ο</sup> Παράδειγμα

$$I_n = \int_0^1 x^n e^{x-1} dx, n \in \mathbb{N}$$

•  $I_n > 0$ , προφανές

•  $x^{n+1} < x^n$  στο  $(0, 1) \Rightarrow \int_0^1 x^{n+1} e^{x-1} dx < \int_0^1 x^n e^{x-1} dx$

δηλαδή  $I_{n+1} < I_n$

•  $x-1 \leq 0 \Rightarrow e^{x-1} \leq 1 \Rightarrow I_n = \int_0^1 x^n e^{x-1} dx < \int_0^1 x^n \cdot 1 dx = \frac{1}{n+1}$

$$\Rightarrow 0 < I_n < \frac{1}{n+1}$$

•  $I_n = \int_0^1 x^n e^{x-1} dx = \int_0^1 x^n (e^{x-1})' dx =$

$$= x^n e^{x-1} \Big|_{x=0}^{x=1} - \int_0^1 (x^n)' e^{x-1} dx =$$

$$= 1 - n \int_0^1 x^{n-1} e^{x-1} dx = 1 - n I_{n-1}$$

•  $n=1 \Rightarrow I_1 = 1 - \int_0^1 e^{x-1} dx =$

$$= 1 - e^{x-1} \Big|_{x=0}^{x=1} = 1 - (e^0 - e^{-1}) =$$

$$= \frac{1}{e}$$

5

$$\tilde{I}_1$$

$$\tilde{I}_n = 1 - n \tilde{I}_{n-2}, \quad n \geq 2$$

$$\text{Θ έστωμε } \varepsilon_n := I_n - \tilde{I}_n$$

Έχουμε:

$$I_n - \tilde{I}_n = -n I_{n-1} + n \tilde{I}_{n-1} = -n(I_{n-1} - \tilde{I}_{n-1})$$

$$\Rightarrow \boxed{\varepsilon_n = -n \varepsilon_{n-1}}$$

$$\varepsilon_2 = -2 \varepsilon_1$$

$$\varepsilon_3 = (-3) \varepsilon_2 = (-3)(-2) \varepsilon_1 = (-1)^2 3! \varepsilon_1$$

$$\dots$$
$$\boxed{\varepsilon_n = (-1)^{n-1} n! \varepsilon_1}$$

Επαγωγή:  $n=2$

$$\underline{n \rightarrow n+1}: \varepsilon_{n+1} = -(n+1) \varepsilon_n = -(n+1) (-1)^{n-1} n! \varepsilon_1 =$$

$$= (-1)^n (n+1)! \varepsilon_1$$

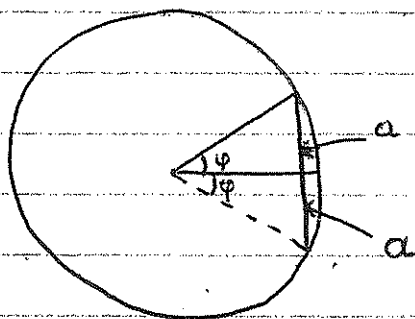
6

2<sup>ο</sup> Παράδειγμα (υπολογισμός του  $\pi$  με την μέθοδο του Αρχιμήδη)

•  $y_n = 2^n \sin \frac{\pi}{2^n}$ ,  $n=1, 2, \dots$

•  $y_1 = 2 \sin \frac{\pi}{2} = 2$

• Για  $n \geq 2$ :



$$\varphi = \frac{\pi}{2^n}, \quad a = \sin \frac{\pi}{2^n} \Leftrightarrow 2\varphi = \frac{2\pi}{2^n}, \quad 2a = 2 \sin \frac{\pi}{2^n}$$

2a: πλευρά του κανονικού πολυγώνου με  $2^n$  πλευρές, εγγεγραμμένη στον μοναδιαίο κύκλο.

Περίμετρος:  $2^n(2a) = 2^{n+1} \sin \frac{\pi}{2^n} = 2y_n$

• Η ακολουθία  $(y_n)_{n \in \mathbb{N}}$  είναι γνησίως αύξουσα και συχλίνει στο  $\pi$ :

$$y_n = 2^n \sin \frac{\pi}{2^n}, \quad \sin(2x) = 2 \sin x \cos x$$

$$\Rightarrow y_n = 2^n \cdot 2 \sin \frac{\pi}{2^{n+1}} \cos \frac{\pi}{2^{n+1}} = y_{n+1} \cos \frac{\pi}{2^{n+1}}$$

$$0 < \cos \frac{\pi}{2^{n+1}} < 1 \Rightarrow y_n < y_{n+1}$$

• Γνωρίζουμε ότι:

$$\lim_{x \rightarrow 0} \frac{\sin x}{x} = 1 \implies \lim_{x \rightarrow 0} \frac{\sin(ax)}{ax} = 1, a \in \mathbb{R}, a \neq 0$$

$$\implies \lim_{x \rightarrow 0} \frac{\sin(ax)}{x} = a$$

Άρα, για  $y_n = 2^n \sin \frac{\pi}{2^n} = \frac{\sin \frac{\pi}{2^n}}{\frac{1}{2^n}}$ , βάσει της

Παραπάνω σχέσης, όπου  $x = \frac{1}{2^n}$  και  $a = \pi$ , έχουμε:

$$\lim_{n \rightarrow +\infty} y_n = \pi$$

• Γνωρίζουμε ότι:

$$\cos(2x) = 1 - 2\sin^2(x)$$

$$y_{n+1} = 2^{n+1} \sin \frac{\pi}{2^{n+1}} \iff$$

$$y_{n+1} = 2^{n+1} \sqrt{\frac{1 - \cos \frac{\pi}{2^n}}{2}} \iff$$

$$y_{n+1} = 2^{n+1} \sqrt{\frac{1}{2} (1 - \sqrt{1 - \sin^2 \frac{\pi}{2^n}})} \iff$$

$$y_{n+1} = 2^{n+1} \sqrt{\frac{1}{2} (1 - \sqrt{1 - (2^{-n} y_n)^2})}$$

8

Άρα, καταλήγουμε:

$$\begin{cases} y_1 = 2 \\ y_{n+1} = 2^{n+1} \sqrt{\frac{1}{2} (1 - \sqrt{1 - (2^{-n} y_n)^2})} \end{cases}, \quad n = 1, 2, \dots$$

Ο αλγόριθμος είναι ασταθής λόγω αφαίρεσης σχεδόν ίσων αριθμών.

Γνωρίζουμε ότι:

$$\sqrt{x} - \sqrt{y} = \frac{(\sqrt{x} - \sqrt{y})(\sqrt{x} + \sqrt{y})}{\sqrt{x} + \sqrt{y}} = \frac{x - y}{\sqrt{x} + \sqrt{y}}$$

$$1 - \sqrt{1 - (2^{-n} y_n)^2} = \frac{1 - (1 - (2^{-n} y_n)^2)}{1 + \sqrt{1 - (2^{-n} y_n)^2}} = \frac{(2^{-n} y_n)^2}{1 + \sqrt{1 - (2^{-n} y_n)^2}}$$

$$y_{n+1} = 2^{n+1} \sqrt{\frac{1}{2} \left( \frac{(2^{-n} y_n)^2}{1 + \sqrt{1 - (2^{-n} y_n)^2}} \right)} \iff$$

$$y_{n+1} = 2^{n+1} \sqrt{\frac{1}{2(1 + \sqrt{1 - (2^{-n} y_n)^2})}} 2^{-n} y_n \iff$$

$$y_{n+1} = \sqrt{\frac{2}{1 + \sqrt{1 - (2^{-n} y_n)^2}}} y_n$$

Άρα, καταλήγουμε:

$$\begin{cases} y_1 = 2 \\ y_{n+1} = \sqrt{\frac{2}{1 + \sqrt{1 - (2^{-n} y_n)^2}}} y_n \end{cases}, \quad n = 1, 2, \dots$$

Ο αλγόριθμος είναι ευσταθής.

Με αυτόν τον αλγόριθμο αποφεύγουμε την αφαίρεση σχεδόν ίσων αριθμών.



## Παράσταση αριθμών ως προς οποιαδήποτε βάση

### Δεκαδικό σύστημα

Βάση: 10, Ψηφία: 0, 1, 2, ..., 9

Παράδειγμα:  $(3.14159)_{10} = 3 \cdot 10^0 + 1 \cdot 10^{-1} + 4 \cdot 10^{-2} +$   
 $+ 1 \cdot 10^{-3} + 5 \cdot 10^{-4} + 9 \cdot 10^{-5}$

Γενικά:  $\pm (a_N a_{N-1} \dots a_1 a_0 a_{-1} a_{-2} \dots)_{10} =$   
 $= \pm (a_N \cdot 10^N + a_{N-1} \cdot 10^{N-1} + \dots + a_1 \cdot 10^1 + a_0 \cdot 10^0 + a_{-1} \cdot 10^{-1} + a_{-2} \cdot 10^{-2} + \dots)$

Ακέραιο μέρος:  $a_N a_{N-1} \dots a_0$

$$p(x) = a_N x^N + a_{N-1} x^{N-1} + \dots + a_1 x + a_0$$

Το ακέραιο μέρος είναι η τιμή του  $p$  στο  $x=10$ ,  
 $p(10)$

Κλασματικό μέρος:  $a_{-1} a_{-2} \dots$

Το κλασματικό μέρος είναι η τιμή της δυναμοσειράς

$$\sum_{k=1}^{+\infty} a_{-k} x^k, \text{ για } x = \frac{1}{10}$$

Η σειρά μπορεί να έχει άπειρο ή πεπερασμένο  
 πλήθος όρων.

10

Μοναδικότητα της Παράστασης;

Ισχυρισμός:  $4.130 = 4.12\bar{9} = 4.12999\dots$

$$3 \cdot 10^{-2} = 2 \cdot 10^{-2} + 9 \cdot 10^{-3} + 9 \cdot 10^{-4} + \dots \Leftrightarrow$$

$$3 \cdot 10^{-2} = 2 \cdot 10^{-2} + 9 \cdot (10^{-3} + 10^{-4} + \dots) \Leftrightarrow$$

$$3 \cdot 10^{-2} = 2 \cdot 10^{-2} + 9 \cdot 10^{-3} \cdot (1 + 10^{-1} + 10^{-2} + \dots)$$

Γνωρίζουμε ότι:

$$\sum_{k=0}^{+\infty} \omega^k = \frac{1}{1-\omega}, \quad |\omega| < 1$$

Άρα, έχουμε:

$$3 \cdot 10^{-2} = 2 \cdot 10^{-2} + 9 \cdot 10^{-3} \cdot \left( \frac{1}{1-10^{-1}} \right) \Leftrightarrow$$

$$3 \cdot 10^{-2} = 2 \cdot 10^{-2} + 9 \cdot 10^{-3} \cdot \left( \frac{1}{9} \cdot 10 \right) \Leftrightarrow$$

$$3 \cdot 10^{-2} = 2 \cdot 10^{-2} + 10^{-2}$$

Για να έχουμε μοναδικότητα αυτό δεν το επιτρέπουμε.

Αν υποθέσουμε ότι:

$$\forall k_0 \in \mathbb{N}, \exists k \geq k_0 \text{ τέτοιο ώστε } a_k \neq 9,$$

τότε η παράσταση είναι μοναδική.

• Σύστημα με βάση  $\theta$ ,  $\theta \geq 2, \theta \in \mathbb{N}$ .

Βάση:  $\theta$ , Ψηφία:  $0, 1, \dots, \theta-1$

Παράδειγμα:  $(100110.11)_2 = 1 \cdot 2^5 + 1 \cdot 2^2 +$   
 $+ 1 \cdot 2^1 + 1 \cdot 2^{-1} + 1 \cdot 2^{-2} = (38.75)_{10}$

Γενικά:  $\pm (a_N a_{N-1} \dots a_1 a_0 a_{-1} a_{-2} \dots)_\theta =$   
 $= \pm (a_N \theta^N + a_{N-1} \theta^{N-1} + \dots + a_1 \theta^1 + a_0 \theta^0 + a_{-1} \theta^{-1} + a_{-2} \theta^{-2} + \dots)$

i) Μετατροπή από σύστημα με βάση  $\theta$  στο δεκαδικό.

α) Ακέραιων αριθμών

Παράδειγμα:  $(53473)_8 = 5 \cdot 8^4 + 3 \cdot 8^3 + 4 \cdot 8^2 +$   
 $+ 7 \cdot 8^1 + 3 \cdot 8^0 \Leftrightarrow$

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

Κάνουμε τις πράξεις από "μέσα" προς τα "έξω".

Σχήμα του Horner:

$$p(x) = a_N x^N + a_{N-1} x^{N-1} + \dots + a_1 x + a_0 \Leftrightarrow$$

$$p(x) = a_0 + x(a_1 + x(a_2 + \dots + x(a_{N-1} + a_N \cdot x) \dots))$$

12

## Διαδικασία υπολογισμού

$$\begin{array}{l} y \leftarrow a_N \\ \text{για } i = N-1, \dots, 0: \\ \quad y \leftarrow a_i + x \cdot y \\ \text{Τότε } y = p(x) \end{array}$$

$$y \leftarrow a_i + x \cdot y, \text{ flop (FLoating point OPeration)}$$

Το σχήμα του Horner απαιτεί  $N$  flop.

## β) Κλασματικών αριθμών

$$\text{Παράδειγμα: } (.11)_2 = 1 \cdot 2^{-1} + 1 \cdot 2^{-2} = (0.75)_{10}$$

★ ii) Μετατροπή από το Δεκαδικό σε ένα σύστημα με βάση 8

α) Ακεραίων αριθμών

Βασίζεται στον αλγόριθμο της διαίρεσης!

Παράδειγμα:  $(369)_{10} \rightsquigarrow$  στο οκταδικό σύστημα

$$(369)_{10} = (\dots a_2 a_1 a_0)_8 \Leftrightarrow$$

$$(369)_{10} = \underbrace{a_0}_{\text{πηλίκιο}} + 8 \underbrace{(a_1 + 8(a_2 + \dots))}_{\text{ακέραιος}}$$

Συμπέρασμα: Το  $a_0$  είναι το υπόλοιπο της διαίρεσης  $369:8$  και το  $a_1 + 8(a_2 + \dots)$  είναι το πηλίκιο.

$$\begin{array}{r} 369 \overline{)8} \\ 49 \overline{)46} \\ 1 \end{array}$$

Άρα,  $a_0 = 1$  και  
 $a_1 + 8(a_2 + \dots) = 46$

$$\begin{array}{r} 46 \overline{)8} \\ 6 \overline{)5} \end{array}$$

Άρα,  $a_1 = 6$  και  
 $a_2 + 8(a_3 + \dots) = 5$

Επομένως:  $a_2 = 5, a_3 + \dots = 0$

Συμπέρασμα:  $(369)_{10} = (561)_8$

Επαλήθευση:  $(561)_8 = 5 \cdot 8^2 + 6 \cdot 8 + 1 = 369$

### β) Κλασματικών αριθμών

$0 < x < 1$  στο δεκαδικό σύστημα

$$x = (.a_{-1}a_{-2}\dots)_\theta = a_{-1}\theta^{-1} + a_{-2}\theta^{-2} + \dots$$

$$\Rightarrow \theta x = \underbrace{a_{-1}} + a_{-2}\theta^{-1} + \dots$$

Άρα, το  $a_{-1}$  είναι το ακέραιο μέρος του  $\theta x$

Παράδειγμα:  $x = (.372)_{10} \rightsquigarrow$  δυαδικό

$$(.372)_{10} = (.a_{-1}a_{-2}\dots)_2$$

Έχουμε:

$$2x = 0.744, \text{ άρα } a_{-1} = 0, \gamma_1 := 0.744$$

$$2\gamma_1 = 1.488, \text{ άρα } a_{-2} = 1, \gamma_2 := 0.488$$

$$2\gamma_2 = 0.976, \text{ άρα } a_{-3} = 0, \gamma_3 := 0.976$$

$$2\gamma_3 = 1.952, \text{ άρα } a_{-4} = 1, \gamma_4 := 0.952$$

...

Συμπέρασμα:  $(.372)_{10} = (.0101\dots)_2$

★ Παράδειγμα:  $\sum_{n=1}^{+\infty} \left( \frac{1}{2^{4n}} + \frac{1}{2^{4n+1}} \right) = \frac{1}{10}$

$$\sum_{n=1}^{+\infty} \left( \frac{1}{2^{4n}} + \frac{1}{2^{4n+1}} \right) =$$

$$= \sum_{n=1}^{+\infty} \left( \frac{1}{2^{4n}} + \frac{1}{2 \cdot 2^{4n}} \right) =$$

$$= \frac{3}{2} \sum_{n=1}^{+\infty} \frac{1}{2^{4n}} =$$

$$= \frac{3}{2} \sum_{n=1}^{+\infty} \left( \frac{1}{2^4} \right)^n =$$

$$= \frac{3}{2} \frac{\frac{1}{2^4}}{1 - \frac{1}{2^4}} =$$

$$= \frac{3}{2} \frac{\frac{1}{16}}{\frac{15}{16}} =$$

$$= \frac{1}{10}$$

Γνωρίζουμε ότι:  $\sum_{n=k}^{+\infty} \omega^n = \frac{\omega^k}{1-\omega}$ ,  $|\omega| < 1$

Άρα, έχουμε:

$$(0.1)_{10} = \sum_{n=1}^{+\infty} \left( \frac{1}{2^{4n}} + \frac{1}{2^{4n+1}} \right) \Leftrightarrow$$

$$(0.1)_{10} = 2^{-4} + 2^{-5} + 2^{-8} + 2^{-9} + 2^{-12} + 2^{-13} + \dots \Leftrightarrow$$

$$(0.1)_{10} = (0.000\underline{1100}110011\dots)_2 \Leftrightarrow$$

$$(0.1)_{10} = (0.000\underline{1100})_2$$

Στο δεκαδικό σύστημα έχουμε μόνο έναν όρο, ενώ στο δυαδικό σύστημα έχουμε άπειρο πλήθος όρων.

(16)

## Αριθμοί μηχανής

Έστω  $x \in \mathbb{R}$ ,  $x \neq 0$

Σε ένα σύστημα με βάση  $\theta$ , ο  $x$  μπορεί να γραφεί ως:

$$(*) \quad x = \pm \cdot \overset{\textcircled{1}}{d_1} d_2 d_3 \dots \cdot \theta^e, \quad \text{με } d_1 \neq 0,$$

όπου  $d_i$  ψηφία ως προς τη βάση  $\theta$   
και  $e$  κατάλληλος ακέραιος.

Η μορφή  $(*)$  λέγεται (κανονική) μορφή κινητής υποδιαστολής.

Το σύνολο των αριθμών μηχανής  $M = M(\theta, t, L, U)$  χαρακτηρίζεται από τις παραμέτρους:

- $\theta$  = βάση του αριθμητικού συστήματος
- $t$  = ακρίβεια = πλήθος των ψηφίων του κλάσματος των αριθμών
- $L$  = κάτω φράγμα του εκθέτη  $e$ .
- $U$  = άνω φράγμα του εκθέτη  $e$ .

δηλαδή  $L \leq e \leq U$

$L, U$  ακέραιοι, συνήθως  $L \approx -U$



Κάθε  $x \in M$ ,  $x \neq 0$ , είναι της μορφής:

$$(+)\ x = \pm d_1 d_2 \dots d_t \cdot \theta^e,$$

με  $d_1 \neq 0$  και  $L \leq e \leq U$

Το  $M$  αποτελείται από όλους τους αριθμούς της μορφής (+) και το μηδέν.

Το  $M$  είναι πεπερασμένο σύνολο.

Μέγιστο κατ' απόλυτη τιμή στοιχείο του  $M$ :

$$d_i = \theta - 1, \quad i = 1, \dots, t \quad \text{και} \quad e = U$$

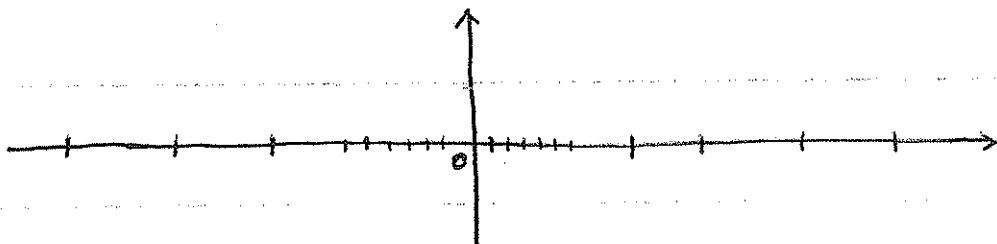
Ελάχιστο κατ' απόλυτη τιμή μη-μηδενικό στοιχείο του  $M$ :

$$d_1 = 1, \quad d_i = 0, \quad i = 2, \dots, t \quad \text{και} \quad e = L$$

Η απόσταση μεταξύ δύο διαδοχικών στοιχείων του  $M$  δεν είναι σταθερή.

$$\tilde{x} = x + \theta^{-t} \cdot \theta^e \iff$$

$$\tilde{x} - x = \theta^{e-t}$$



18

Το  $M$  δεν είναι κλειστό ως προς την  
πολλαπλασιασμό, δηλαδή:

$$x, x^* \in M \not\Rightarrow x \cdot x^* \in M$$

Παράδειγμα:  $(\underbrace{100 \dots 0 \cdot 10^L}_{\in M} \times \underbrace{100 \dots 0 \cdot 10^L}_{\in M}) \notin M$

Το  $M$  δεν είναι κλειστό ως προς την  
πρόσθεση, δηλαδή:

$$x, x^* \in M \not\Rightarrow x + x^* \in M$$

Παράδειγμα:  $\theta = 10, t = 5, 1, 10^{-5} \in M$

$$1 + 10^{-5} = 1.00001 = 0.\underbrace{100001}_{\substack{\uparrow \\ \text{6 ψηφία}}} \cdot 10^1 \notin M$$

Μας ενδιαφέρει το  $M$  να είναι όσο πιο πυκνό  
και όσο πιο επύ γίνεται, δηλαδή μεγάλο  $t$   
και μεγάλο διάστημα  $[L, U]$ .

## Προέγχιση πραγματικών αριθμών με αριθμους μηχανής.

i)  $|x| >$  μέγιστο στοιχείο του  $M$ .

Υπερχείλιση (overflow).

Οι υπολογισμοί σταματούν.

ii)  $0 < |x| < 100 \dots 0 \cdot \theta^L$

Υπεκχείλιση.

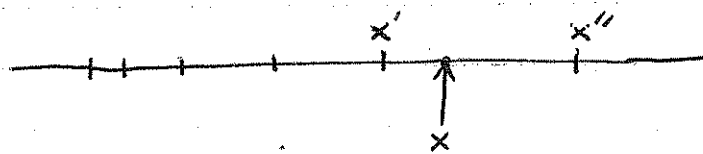
Συνήθως, ο  $x$  προσεγγίζεται με το μηδέν και συνεχίζονται οι υπολογισμοί.

iii)  $100 \dots 0 \cdot \theta^L \leq |x| \leq d_1 d_2 \dots d_t \cdot \theta^u$

με  $d_i = \theta - 1, i = 1, \dots, t$

Ο  $x$  προσεγγίζεται με ένα στοιχείο  $fl(x) \in M$

⊕ Συνήθως ισχύει:  $|x - fl(x)| \leq |x - y|, \forall y \in M$



(20)

(\*)

$$\text{λοχυρισμός: } \left| \frac{x - fl(x)}{x} \right| \leq \frac{1}{2} \theta^{1-t}$$

ανεξάρτητο

του  $x$

$$α) x \in M, fl(x) = x$$

H (\*) λοχύει.

β)  $x \notin M$ . Τότε υπάρχουν δύο διαδοχικά στοιχεία  $x', x'' \in M$  τέτοια ώστε  $x' < x < x''$

Τότε λοχύει:

$$|x - fl(x)| \leq \frac{|x' - x''|}{2}$$

Γνωρίζουμε ότι:

$$\begin{array}{c} x \\ \text{---} \\ a \qquad b \end{array} \quad \min(x-a, b-x) \leq \frac{b-a}{2}$$

Έστω  $x = \alpha \cdot \beta^k$ ,  $\alpha > 0$ :

$$x = \alpha \cdot d_1 d_2 \dots d_t \underbrace{d_{t+1} \dots}_{\theta^t} \cdot \beta^k$$

Τότε:

$$x' = \alpha \cdot d_1 d_2 \dots d_t \cdot \beta^k$$

$$x'' = (\alpha \cdot d_1 d_2 \dots d_t + \theta^t) \cdot \beta^k$$

$$\Rightarrow x'' - x' = \theta^t \cdot \beta^k = \beta^{k-t}$$

$$x \geq .100\dots 0 \cdot \theta^k$$

$$\frac{|x - fl(x)|}{|x|} \leq \frac{\frac{1}{2} \theta^{k-t}}{.100\dots 0 \cdot \theta^k} \iff$$

$$\frac{|x - fl(x)|}{|x|} \leq \frac{1}{2} \frac{\theta^{-t}}{0.1} \iff$$

$$\frac{|x - fl(x)|}{|x|} \leq \frac{1}{2} \frac{\theta^{-t}}{\theta^{-1}} \iff$$

$$\frac{|x - fl(x)|}{|x|} \leq \frac{1}{2} \theta^{1-t}$$

Η διαδικασία για να πετύχουμε την  $\oplus$  λέγεται στρογγύλευση.

Παράδειγμα:  $\theta=10$ ,  $t=5$

$$x = .a_1 a_2 \dots a_5 | a_6 a_7 \dots \cdot 10^k$$

α) Αν  $a_6 \geq 5$  τότε:

$$fl(x) = x'' = (.a_1 a_2 \dots a_5 + 10^{-5}) \cdot 10^k$$

β) Αν  $a_6 < 5$  τότε:

$$fl(x) = x' = .a_1 a_2 \dots a_5 \cdot 10^k$$

(Αν  $a_6 = 5$  και  $a_i = 0$ ,  $i \geq 7$ , τότε μπορούμε ως  $fl(x)$  να επιλέξουμε είτε τον  $x'$  είτε τον  $x''$ )

22

Αποκοπή:  $t=5$

$$x = \alpha_1 \alpha_2 \dots \alpha_5 \cdot 10^k$$

Στην περίπτωση της αποκοπής για το σχετικό σφάλμα λήξει η εκτίμηση:

$$\left| \frac{x - fl(x)}{x} \right| \leq \beta^{1-t}$$

Ενικά:

$$\left| \frac{x - fl(x)}{x} \right| \leq u$$

$$u = \begin{cases} \rightarrow \frac{1}{2} \beta^{1-t} & , \text{ για την στρογγύλιση} \\ \rightarrow \beta^{1-t} & , \text{ για την αποκοπή} \end{cases}$$

$u$ : μοναδιαίο σφάλμα στρογγύλισης

Πράξεις

$$x * y, * \in \{+, -, \cdot, :\}$$

$$z = fl(fl(x) * fl(y))$$

υποθέσαμε ότι η

πράξη αυτή γίνεται

ακριβώς.

Κατα κανόνα η πράξη αυτή γίνεται με δύο ακρίβεια.

Παράδειγμα:  $\beta=10, t=5, u=-L=10$ , στρογγύλιση

$$x = 5891.26, y = 0.0773414$$

$$x + y$$

$$fl(x) = .58913 \cdot 10^4$$

$$fl(y) = .77341 \cdot 10^{-1} = .0000077341 \cdot 10^4$$

$$fl(x) + fl(y) = .\underline{5891377341} \cdot 10^4$$

$$z = fl(fl(x) + fl(y)) = .58914 \cdot 10^4$$

$$z = fl(fl(x) + fl(y)) \neq x + y$$

$$z = fl(fl(x) + fl(y)) \neq fl(x + y)$$

$$z = fl(fl(x) + fl(y)) \neq fl(x) + fl(y)$$

(24)

Παράδοξα:  $\theta = 10$ ,  $t = 5$ ,  $U = -L = 10$ , στρογγυλευση

$$a_1 = 1, a_2 = a_3 = 3 \cdot 10^{-5}, a_1 + a_2 + a_3, a_1, a_2, a_3 \in M$$

$$fl(a_1 + a_2) = fl(1.00003) = 1$$

$$fl(fl(a_1 + a_2) + a_3) = fl(1.00003) = 1$$

$$\Rightarrow (a_1 + a_2) + a_3 = 1$$

$$fl(a_2 + a_3) = 6 \cdot 10^{-5}$$

$$fl(a_1 + fl(a_2 + a_3)) = fl(1.00006) = 1.0001$$

$$\Rightarrow a_1 + (a_2 + a_3) = 1.0001$$

Διαφορετικά αποτελέσματα!

Έχει σημασία η σειρά με την οποία γίνονται οι προσθέσεις σε ένα άθροισμα!

Για κάθε  $0 < |x| < 5 \cdot 10^{-5}$  έχουμε  $fl(1 + fl(x)) = 1$

Γενικά: Για  $0 < |x| < \frac{1}{2} \theta^{1-t}$  έχουμε  $fl(1 + fl(x)) = 1$

Το  $\frac{1}{2} \theta^{1-t}$  λέγεται έφικτον ή μηδέν των μηχανών

$$x = .0 \dots 0 d_t d_{t+1} \dots, d_t < \frac{\theta}{2}$$

$$1+x = \underbrace{1.00 \dots 0}_{t+1 \text{ θέσεις}} d_t d_{t+1} \dots$$

$$fl(1+x) = 1.00 \dots 0 = 1$$



★ Επιρροή των σφαλμάτων στρογγύλευσης στους υπολογισμούς (όταν κάνουμε μία μόνο πράξη)

$$\left| \frac{fl(fl(x) * fl(y)) - x * y}{x * y} \right| = ;$$

Να εκτιμήσουμε το σχετικό σφάλμα στην πράξη  $x * y$

$$\left| \frac{x - fl(x)}{x} \right| \leq u \iff$$

$$\boxed{fl(x) = x(1 + \varepsilon) \text{ και } |\varepsilon| \leq u}, \text{ με } \varepsilon = \varepsilon(x)$$

Ισχυρισμός:  $\varepsilon_i \in [-u, u]$ ,  $i = 1, \dots, m$

$$\implies \text{υπάρχει } \varepsilon \in [-u, u] : \prod_{i=1}^m (1 + \varepsilon_i) = (1 + \varepsilon)^m$$

$$\varphi(x) = (1 + x)^m, \quad x \in [-u, u]$$

$$(1 - u)^m \leq \prod_{i=1}^m (1 + \varepsilon_i) \leq (1 + u)^m \iff$$

$$\varphi(-u) \leq \prod_{i=1}^m (1 + \varepsilon_i) \leq \varphi(u)$$

Σύμφωνα με το θεώρημα της ενδιάμεσης τιμής υπάρχει  $\varepsilon \in [-u, u]$  τέτοιο ώστε:

$$\boxed{\prod_{i=1}^m (1 + \varepsilon_i) = \varphi(\varepsilon) = (1 + \varepsilon)^m}$$

26

$$z = fl(fl(x) * fl(y))$$

$$fl(x) = x(1 + \varepsilon_1), fl(y) = y(1 + \varepsilon_2), |\varepsilon_i| \leq u$$

Πολυπλοκότητα

$$z = fl(fl(x) \cdot fl(y)) \Leftrightarrow$$

$$z = fl(x \cdot y (1 + \varepsilon_1)(1 + \varepsilon_2)) \Leftrightarrow$$

$$z = x \cdot y (1 + \varepsilon_1)(1 + \varepsilon_2)(1 + \varepsilon_3) \Leftrightarrow$$

$$z = x \cdot y (1 + \varepsilon)^3$$

$$\frac{z - xy}{xy} = \frac{xy(1 + \varepsilon)^3 - xy}{xy} = (1 + \varepsilon)^3 - 1 = 3\varepsilon + 3\varepsilon^2 + \varepsilon^3$$

$$\Rightarrow \left| \frac{z - xy}{xy} \right| = |3\varepsilon + 3\varepsilon^2 + \varepsilon^3| \leq 3|\varepsilon| + 3|\varepsilon|^2 + |\varepsilon|^3 \Leftrightarrow$$

$$\left| \frac{z - xy}{xy} \right| \leq 3u + 3u^2 + u^3 \Leftrightarrow$$

$$\left| \frac{z - xy}{xy} \right| \leq 3u + 4u^2 \Leftrightarrow$$

$$\left| \frac{z - xy}{xy} \right| \leq 3u + O(u^2)$$

Λέμε ότι το σχετικό σφάλμα στον πολυπλοκότητα είναι το πολύ τρεις φορές το μοναδιαίο σφάλμα στην πράξη.

Διαίρεση

$$z = fl\left(\frac{fl(x)}{fl(y)}\right) = fl\left(\frac{x(1+\varepsilon_1)}{y(1+\varepsilon_2)}\right) = \frac{x(1+\varepsilon_1)}{y(1+\varepsilon_2)}(1+\varepsilon_3)$$

Γνωρίζουμε ότι:

$$\frac{1}{1+\varepsilon_2} = 1 + \delta \Leftrightarrow$$

$$\delta = \frac{1}{1+\varepsilon_2} - 1 = \frac{-\varepsilon_2}{1+\varepsilon_2}$$

$$\Rightarrow |\delta| \leq \frac{u}{1-u} = u(1+u+\dots) = u + u^2 + u^3 + \dots = u + O(u^2) \Leftrightarrow$$

$$\boxed{|\delta| \leq u + O(u^2)}$$

Άρα, έχουμε:

$$z = \frac{x}{y} (1+\varepsilon_1)(1+\varepsilon_3)(1+\delta) \Leftrightarrow$$

$$z = \frac{x}{y} (1+\varepsilon)^2 (1+\delta)$$

$$\frac{z - \frac{x}{y}}{\frac{x}{y}} = \frac{\frac{x}{y} (1+\varepsilon)^2 (1+\delta) - \frac{x}{y}}{\frac{x}{y}} = (1+\varepsilon)^2 (1+\delta) - 1 \Leftrightarrow$$

$$\frac{z - \frac{x}{y}}{\frac{x}{y}} = 2\varepsilon + \delta + \varepsilon^2 + 2\varepsilon\delta + \delta\varepsilon^2$$

$$\Rightarrow \left| \frac{z - \frac{x}{y}}{\frac{x}{y}} \right| \leq 3u + O(u^2)$$

(28)

## Πρόσθεση - Αφαίρεση

$$z = fl(fl(x) + fl(y)) \Leftrightarrow$$

$$z = fl(x(1+\varepsilon_1) + y(1+\varepsilon_2)) \Leftrightarrow$$

$$z = x(1+\varepsilon_1)(1+\varepsilon_3) + y(1+\varepsilon_2)(1+\varepsilon_3) \Leftrightarrow$$

$$z = x(1+\varepsilon)^2 + y(1+\delta)^2 \Leftrightarrow$$

$$z = x + 2\varepsilon x + y + 2\delta y + x\varepsilon^2 + y\delta^2 \Leftrightarrow$$

$$z = x + 2\varepsilon x + y + 2\delta y + O(u^2)$$

$$\Rightarrow z \approx (x+y) + 2(\varepsilon x + \delta y)$$

$$\Rightarrow \frac{z - (x+y)}{x+y} \approx 2 \frac{\varepsilon x + \delta y}{x+y}$$

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

1<sup>η</sup> Περίπτωση:  $x, y$  ομόσημοι

$$|x+y| = |x| + |y|$$

Σφάλμα:  $\left| \frac{z-(x+y)}{x+y} \right| \leq 2u$

2<sup>η</sup> Περίπτωση:  $x, y$  ετερόσημοι

Στη χειρότερη περίπτωση  $\varepsilon \approx -\delta$ ,  $|\varepsilon| \approx u$ .

Οπότε θα έχουμε:

$$\left| \frac{z-(x+y)}{x+y} \right| \approx 2 \cdot \frac{|x-y|}{|x+y|} u$$

Το κλάσμα  $\frac{|x-y|}{|x+y|}$  μπορεί να είναι μεγάλο.

Αν  $x \approx -y$  τότε ο παράγοντας  $\frac{|x-y|}{|x+y|}$  μπορεί να γίνει πολύ μεγάλος!

Τα σφάλματα στρογγύλευσης μπορούν να έχουν κοζοστροφική επίρροή στην αφαίρεση σχεδόν ίσων αριθμών.

Η αφαίρεση σχεδόν ίσων αριθμών πρέπει να αποφεύγεται (ή να γίνεται με διπλή ακρίβεια)

Σημείωση: Αν  $x, y$  αριθμοί μηχανής, τότε η αφαίρεσή τους γίνεται με ακρίβεια ακόμα και αν  $x \approx y$

$$z = fl(fl(x) + fl(y)) = fl(x+y) = (x+y)(1+\varepsilon)$$

$$\Rightarrow \frac{z-(x+y)}{x+y} = \varepsilon \Rightarrow \left| \frac{z-(x+y)}{x+y} \right| \leq u$$

30

Παράδειγμα:  $\beta=10$ ,  $t=5$ ,  $u=-L=10$ , σε περίπτωση

$$x = .45142|708$$

$$y = -.45115|944$$

$$x+y = .26764 \cdot 10^{-3}$$

$$z = fl(fl(x) + fl(y)) \Leftrightarrow$$

$$z = fl(.45143 - .45116) \Leftrightarrow$$

$$z = .00027 = .27000 \cdot 10^{-3}$$

άσος!

Πώς μπορούμε να αποφύγουμε την αφαίρεση  
σχέδων ίσων αριθμών;

1<sup>ο</sup> Παράδειγμα:  $\theta = 10$ ,  $t = 10$

$$\sqrt{7298} = .88831692926 \cdot 10^2$$

$$\sqrt{7297} = .8883130079 \cdot 10^2$$

$$\sqrt{7298} - \sqrt{7297} = .562847(0000) \cdot 10^{-2}$$

↑ στον υπολογισμό

$$\sqrt{x} - \sqrt{y} = \frac{x-y}{\sqrt{x} + \sqrt{y}}$$

$$\sqrt{7298} - \sqrt{7297} = \frac{1}{\sqrt{7298} + \sqrt{7297}} \iff$$

$$\frac{1}{\sqrt{7298} + \sqrt{7297}} = .5628468294 \cdot 10^{-2} \text{ (πολύ καλή προσέγγιση)}$$

↑ στον υπολογισμό

32

## 2<sup>ο</sup> Παράδειγμα

$$f(x) = x - \sin x$$

Θέλουμε να υπολογίσουμε τιμές της  $f$  για  $|x|$  μικρή.

$$\text{Επειδή: } \lim_{x \rightarrow 0} \frac{\sin x}{x} = 1,$$

θα αναγκαστούμε να απαρίθμησε σχεδόν ίσους αριθμούς.

Ανάπτυξη Taylor (ως προς το μηδέν)

$$\sin x = x - \frac{x^3}{6} + \varepsilon(x),$$

$$\text{με } |\varepsilon(x)| \leq \frac{|x|^5}{180}$$

Επομένως:  $f(x) \approx \frac{x^3}{6}$  (ο υπολογισμός γίνεται χωρίς πρόβλημα!)



Σφάλματα στον υπολογισμό αθροισμάτων

Παράδειγμα:  $S_n = 1 + \sum_{k=1}^n \frac{1}{k^2+k}$ ,  $n \in \mathbb{N}$

$(S_n)_{n \in \mathbb{N}}$  γνήσια αύξουσα ακολουθία

$$S_n = 1 + \sum_{k=1}^n \frac{1}{k^2+k} \Leftrightarrow$$

$$S_n = 1 + \sum_{k=1}^n \frac{1}{k(k+1)} \Leftrightarrow$$

$$S_n = 1 + \sum_{k=1}^n \left( \frac{1}{k} - \frac{1}{k+1} \right) \Leftrightarrow$$

$$S_n = 1 + \left( 1 - \frac{1}{2} \right) + \left( \frac{1}{2} - \frac{1}{3} \right) + \dots + \left( \frac{1}{n} - \frac{1}{n+1} \right)$$

Αυτό είναι ένα τηλεσκοπικό άθροισμα, οπότε:

$$S_n = 2 - \frac{1}{n+1}$$

$$S_n \rightarrow 2, n \rightarrow +\infty$$

$$S_{9999} = 1.9999$$

Αναδρομικός τύπος:

$$S_0 = 1, S_k = S_{k-1} + \frac{1}{k(k+1)}, k = 1, \dots, n$$

(Αθροίζουμε τους όρους από τον μεγαλύτερο προς τον μικρότερο)

Έστω  $\epsilon = 10^{-7}, t = 10 \Rightarrow$  Αποτέλεσμα:  $\tilde{S}_{9999} = \boxed{1.999899972}$

34

Αθροίζοντας από τον πιο μικρό προς τον πιο μεγάλο όρο του ίδιου αθροίσματος έχουμε:

$$(*) \begin{cases} T_0 = \frac{1}{n(n+1)} \\ T_k = T_{k-1} + \frac{1}{(n-k)(n-k+1)}, \quad k=1, \dots, n-1 \\ T_n = T_{n-1} + 1 \end{cases}$$

Προφανώς  $T_n = S_n$

Υλοποιώντας τον αλγόριθμο (\*) με τον υπολογιστή με  $\theta=10$ ,  $t=10$ , παίρνουμε:

$$\tilde{T}_{9999} = \boxed{1.999900\dots 0}$$

Γιατί δίνει ο δεύτερος αλγόριθμος καλύτερα αποτελέσματα από τον πρώτο;

Παρατήρηση:

Έστω  $\lambda, \mu \in \mathbb{R}$  και  $\varepsilon_1, \varepsilon_2 \in [-u, u]$

Τότε, υπάρχει  $\varepsilon_3 \in [-u, u]$  τέτοιο ώστε

$$\lambda \varepsilon_1 + \mu \varepsilon_2 = (|\lambda| + |\mu|) \varepsilon_3 \iff$$

$$\varepsilon_3 = \frac{\lambda \varepsilon_1 + \mu \varepsilon_2}{|\lambda| + |\mu|}$$

$$\implies |\varepsilon_3| \leq \frac{|\lambda| \overset{\leq u}{|\varepsilon_1|} + |\mu| \overset{\leq u}{|\varepsilon_2|}}{|\lambda| + |\mu|} \leq u$$

Πρόβλημα: Έστω  $a_1, a_2, \dots, a_N \in M$

Θέλουμε να υπολογίσουμε το άθροισμα:

$$S_N = \sum_{i=1}^N a_i$$

Αλγόριθμος:

$$S_1 = a_1, S_2 = a_1 + a_2, \dots, S_k = S_{k-1} + a_k, k = 2, \dots, N$$

Τι μπορώ να πω για το σφάλμα, αν υλοποιήσω αυτόν τον αλγόριθμο σε έναν υπολογιστή;

Προσεγγίσεις:

$$\begin{aligned} \tilde{S}_1 &= a_1 \\ \tilde{S}_k &= fl(\tilde{S}_{k-1} + a_k), k = 2, \dots, N \end{aligned}$$

Έχουμε:

$$\tilde{S}_2 = fl(\tilde{S}_1 + a_2) = fl(a_1 + a_2)$$

Γνωρίζουμε ότι:  $|\delta| \leq u$ ,  $|\varepsilon_2| \leq u$ ,  $fl(x) = x(1 + \varepsilon)$ ,  $|\varepsilon| \leq u$

$$\tilde{S}_2 = \underbrace{(a_1 + a_2)}_{S_2} (1 + \delta) = S_2 + S_2 \cdot \delta = S_2 + |S_2| \varepsilon_2$$

(36)

$$\tilde{s}_3 = f(|\tilde{s}_2 + a_3|) = (\tilde{s}_2 + a_3)(1 + \delta'), \quad |\delta'| \leq u$$

$$\tilde{s}_3 = (|s_2| + |s_2| \varepsilon_2 + a_3)(1 + \delta') \Leftrightarrow$$

$$\tilde{s}_3 = (s_3 + |s_2| \varepsilon_2)(1 + \delta')$$

$$\Rightarrow \tilde{s}_3 \approx s_3 + |s_2| \varepsilon_2 + s_3 \delta' = s_3 + (|s_2| + |s_3|) \varepsilon_3$$

Παρόμοια, βρίσκουμε:

$$\tilde{s}_N \approx s_N + (|s_2| + |s_3| + \dots + |s_N|) \varepsilon_N$$

Παραλείπουμε όρους της τάξης του  $u^2$

$$\frac{\tilde{s}_N - s_N}{s_N} \approx \frac{|s_2| + |s_3| + \dots + |s_N|}{s_N} \varepsilon_N$$

Σχετικό σφάλμα:

$$\frac{|\tilde{s}_N - s_N|}{|s_N|} \approx \frac{|s_2| + |s_3| + \dots + |s_N|}{|s_N|} |\varepsilon_N|$$

Εστω;  $|s_2| + |s_3| + \dots + |s_N| = \gamma_N$ , τότε:

$$P_N := \frac{\gamma_N}{|s_N|}$$

Συντελεστής μετάδοσης του σχετικού σφάλματος στην αλγόριθμική μας.

Ειδική περίπτωση:  $a_i > 0, i = 1, \dots, N$

$$|s_2| + \dots + |s_N| = (N-1)a_1 + (N-1)a_2 + (N-2)a_3 + \dots + 2a_{N-1} + a_N$$

Το άθροισμα γίνεται ελάχιστο αν:

$$a_1, a_2 \leq a_3 \leq a_4 \leq \dots \leq a_N$$

Το άθροισμα γίνεται μέγιστο αν:

$$a_1, a_2 \geq a_3 \geq a_4 \geq \dots \geq a_N \text{ (πρέπει να το αποφύγουμε!)}$$

Όταν το  $p_N$  είναι πολύ μεγάλο ο αλγόριθμος είναι ασταθής.

Το  $p_N$  είναι μεγάλο, αν κάποια από τα ενδιάμεσα άθροισματα  $s_k$  έχουν απόλυτη τιμή πολύ μεγαλύτερη από  $|s_N|$ .

Παράδειγμα: Προσέγγιση του  $e^{-x}$  για  $x \gg 1$ .

$$\Gammaνωρίζουμε ότι:  $e^x = 1 + x + \frac{x^2}{2!} + \frac{x^3}{3!} + \dots + \frac{x^N}{N!} + \dots$$$

$$s_N(x) = 1 - x + \frac{x^2}{2!} - \frac{x^3}{3!} + \dots + (-1)^{N-1} \frac{x^{N-1}}{(N-1)!}$$

$$s_N(x) \rightarrow e^{-x}, N \rightarrow +\infty$$

Για αρκετά μεγάλο  $N$  το  $s_N(x)$  είναι καλή προσέγγιση του  $e^{-x}$ .

38

Για  $x=100$  έχουμε  $e^{-100} \approx 0$ , ενώ:

$$s_1=1, s_2=-99, s_3=4901, s_4 \approx -161766 \text{ κλπ.}$$

Ο αριθμός είναι ασταθής. Παράγωδος αποτυχία!

$$e^{-x} = \frac{1}{e^x} \approx \frac{1}{1+x+\frac{x^2}{2!}+\dots+\frac{x^{n-1}}{(n-1)!}}$$

Ο αριθμός αυτός είναι ευσταθής.

## Ευστάθεια αλγορίθμων

Ένας αλγόριθμος λέγεται ασταθής, αν είναι ευαίσθητος σε σφάλματα στρογγύλευσης, δηλαδή αν τα σφάλματα που γίνονται κατά την παράσταση των αριθμών και τις πράξεις σε έναν υπολογιστή, είναι δυνατόν να επιφέρουν μεγάλες μεταβολές στο τελικό αποτέλεσμα.

Ένας αλγόριθμος λέγεται ευσταθής, αν τα τελικά αποτελέσματά του, δεν επηρεάζονται πολύ από τα σφάλματα στρογγύλευσης που γίνονται σε κάθε βήμα του.

### Παραδείγματα:

$$1) e^{-x}, x \gg 1, \text{ ασταθής}$$

$$\frac{1}{e^x}, \text{ ευσταθής}$$

$$2) I_n = \int_0^1 x^n e^{-x} dx, n=1, 2, \dots$$

$$\begin{cases} I_1 = \frac{1}{e} \end{cases}$$

$$\begin{cases} I_n = 1 - n I_{n-1}, n=2, 3, \dots \end{cases}, \text{ ασταθής}$$

$$\begin{cases} I_{20} = 0 \end{cases}$$

$$\begin{cases} I_{n-1} = \frac{1 - I_n}{n}, n=2, 3, \dots \end{cases}, \text{ ευσταθής}$$

40

$$3) y_n = 2^n \sin \frac{\pi}{2^n}, n \in \mathbb{N}$$

$$\begin{cases} y_1 = 2 \end{cases}$$

$$\begin{cases} y_{n+1} = 2^{n+1} \sqrt{\frac{1}{2}(1 - \sqrt{1 - (2^{-n} \cdot y_n)^2})} \end{cases}, n=1, 2, \dots$$

ασταθής

$$\begin{cases} y_1 = 2 \end{cases}$$

$$\begin{cases} y_{n+1} = \sqrt{\frac{2}{1 + \sqrt{1 - (2^{-n} \cdot y_n)^2}}} \cdot y_n \end{cases}, n=1, 2, \dots$$

ευσταθής



## Κατάσταση προβλημάτων

Λέμε ότι ένα πρόβλημα έχει καλή κατάσταση, αν μικρές μεταβολές στα δεδομένα του, έχουν ως αποτέλεσμα μικρές μεταβολές της λύσης του.

Λέμε ότι ένα πρόβλημα έχει κακή κατάσταση, αν είναι δυνατόν μικρές μεταβολές στα δεδομένα του, να οδηγούν σε μεγάλες μεταβολές της λύσης του.

## Παράδειγμα

$$(x-2)^6 = 0 \Rightarrow \text{Λύση: } x^* = 2$$

$$(x-2)^6 = 10^{-6} \Rightarrow \text{Λύσεις: } x_k = 2 + \frac{1}{10} e^{\frac{2\pi i k}{6}},$$

όπου  $k=0, \dots, 5$ , Μεταβολή στα δεδομένα:  $10^{-6}$

Γνωρίζουμε ότι:  $|e^{ix}| = 1$ , για  $x \in \mathbb{R}$

$$\text{Μεταβολή στη λύση: } |x_k - 2| = \frac{1}{10} = 10^{-1}$$

Δηλαδή η μεταβολή στη λύση είναι  $10^5$  επί τη μεταβολή στα δεδομένα. Κακή κατάσταση!

$$x-2=0 \Rightarrow \text{Λύση: } x^* = 2, \quad x-2=10^{-6} \Rightarrow \text{Λύση: } x^* = 2+10^{-6}$$

Η μεταβολή στη λύση είναι  $10^{-6}$ , όσο είναι και η μεταβολή στα δεδομένα. Καλή κατάσταση!

## Ασκήσεις

Άσκηση 1.2

$$\alpha) 1 - \cos x = 2 \sin^2 \frac{x}{2}$$

|x| μικρή (χωρίς Taylor)

$$\beta) e^{x-y} = \frac{e^x}{e^y}$$

$$\gamma) \log x - \log y = \log \frac{x}{y}$$

$$\delta) \sin(a+x) - \sin(a) = 2 \sin \frac{x}{2} \cos(a + \frac{x}{2})$$

|x| μικρή

$$\sin a - \sin b = 2 \sin \left( \frac{a-b}{2} \right) \cos \left( \frac{a+b}{2} \right)$$

2

### Άσκηση 1.3

$$x^2 - 2ax + \beta = 0, \quad a, \beta > 0, \quad a^2 \gg \beta$$

$$x_1 = a + \sqrt{a^2 - \beta}, \quad \text{ευσταθής}$$

$$x_2 = a - \sqrt{a^2 - \beta}, \quad \text{αφαίρεση σχεδόν ίσων αριθμών.}$$

$$x_2 = a - \sqrt{a^2 - \beta} = \frac{(a - \sqrt{a^2 - \beta})(a + \sqrt{a^2 - \beta})}{a + \sqrt{a^2 - \beta}} \Leftrightarrow$$

$$x_2 = \frac{a^2 - (a^2 - \beta)}{a + \sqrt{a^2 - \beta}} = \frac{\beta}{a + \sqrt{a^2 - \beta}} = \frac{\beta}{x_1} \Leftrightarrow$$

$$x_2 = \frac{\beta}{x_1}$$

$$x^2 - 2ax + \beta = (x - x_1)(x - x_2) \Leftrightarrow$$

$$x^2 - 2ax + \beta = x^2 - \underbrace{(x_1 + x_2)}_{=2a}x + \underbrace{x_1 x_2}_{=\beta} = \beta$$

(3)

Άσκηση 1.4

$$x^3 + 3px + 2q = 0, \quad p, q \in \mathbb{R} \quad \text{και} \quad (p^3 + q^2 > 0)$$

α) Η εξίσωση έχει ακριβώς μία πραγματική ρίζα  $p$ , και:

$$p = u - v, \quad \text{με} \quad u = (\sqrt{p^3 + q^2} - q)^{1/3}$$

$$v = (\sqrt{p^3 + q^2} + q)^{1/3}$$

Τύπος του CardanoΛύση

$$f(x) = x^3 + 3px + 2q, \quad x \in \mathbb{R}$$

Υπαρξη ρίζας

$$\lim_{x \rightarrow +\infty} f(x) = +\infty, \quad \lim_{x \rightarrow -\infty} f(x) = -\infty$$

Επομένως η  $f$  παίρνει και θετικές και αρνητικές τιμές. Άρα, σύμφωνα με το θεώρημα της ενδιάμεσης τιμής, (η  $f$  είναι συνεχής) η  $f$  παίρνει και την τιμή μηδέν.

4

### Μοναδικότητα

$$f'(x) = 3x^2 + 3p = 3(x^2 + p)$$

•  $p \geq 0$ :  $f'(x) > 0, \forall x \in \mathbb{R}, x \neq 0$   
 $f'(x) \geq 0, \forall x \in \mathbb{R}$

Η  $f$  είναι γνήσια αύξουσα, οπότε δεν μπορεί να έχει περισσότερες από μία ρίζες.

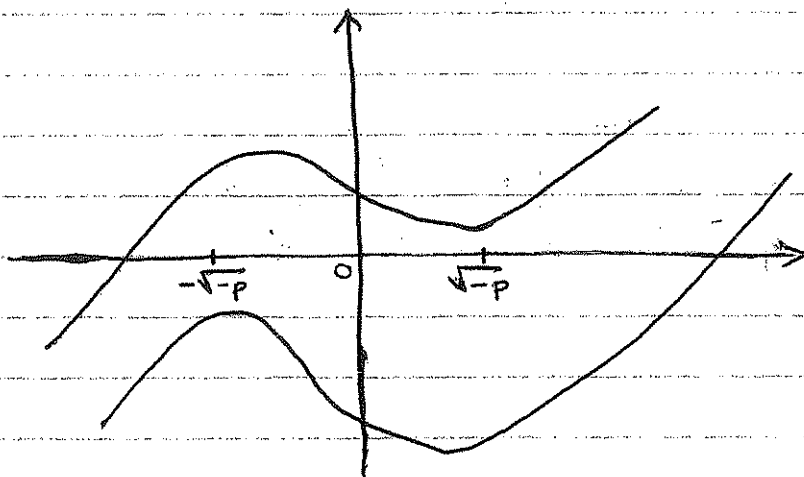
•  $p < 0$ :  $f'(x) = 3(x^2 + p)$   
 $f'(x) = 0 \iff x = \pm \sqrt{-p}$

$f'$	+	-	+
$f$	$\nearrow$	$\searrow$	$\nearrow$

Η  $f$  είναι γνήσια αύξουσα στο  $(-\infty, -\sqrt{-p})$

Η  $f$  είναι γνήσια φθίνουσα στο  $(-\sqrt{-p}, \sqrt{-p})$

Η  $f$  είναι γνήσια αύξουσα στο  $(\sqrt{-p}, +\infty)$



Για να έχουμε μόνο μία ρίζα θα πρέπει οι τιμές  $f(-\sqrt{-p})$  και  $f(\sqrt{-p})$  να είναι ομόσημες.

5

$$\begin{aligned} f(-\sqrt{-p}) &= 2(q - p\sqrt{-p}) \\ f(\sqrt{-p}) &= 2(q + p\sqrt{-p}) \end{aligned}$$

$$\Rightarrow f(-\sqrt{-p}) \cdot f(\sqrt{-p}) = 4(q^2 + p^2) > 0$$

1<sup>η</sup> περίπτωση:  $f(-\sqrt{-p}) > 0$

Τότε η  $f$  έχει μία ρίζα στο  $(-\infty, -\sqrt{-p})$ ,  
καμία ρίζα στο  $(-\sqrt{-p}, \sqrt{-p})$  και  
καμία ρίζα στο  $(\sqrt{-p}, +\infty)$

2<sup>η</sup> περίπτωση:  $f(-\sqrt{-p}) < 0$

Τότε η  $f$  δεν έχει ρίζα στο  $(-\infty, -\sqrt{-p})$ ,  
ούτε στο  $(-\sqrt{-p}, \sqrt{-p})$  έχει ρίζα,  
και έχει ρίζα στο  $(\sqrt{-p}, +\infty)$

6

### Άσκηση 1.4

$$x^3 + 3px + 2q = 0, \quad p, q \in \mathbb{R}, \quad p^3 + q^2 > 0$$

$$f(x) = x^3 + 3px + 2q$$

α) Η εξίσωση έχει ακριβώς μία πραγματική ρίζα  $p$  και λοχύει:

$$p = u - v, \quad \text{με } u = (\sqrt{p^3 + q^2} - q)^{1/3},$$

$$v = (\sqrt{p^3 + q^2} + q)^{1/3}$$

### Τύπος του Cardano

$$\alpha) f(p) = f(u-v) = (u-v)^3 + 3p(u-v) + 2q \Leftrightarrow$$

$$f(p) = (u^3 - v^3) - 3uv(u-v) + 3p(u-v) + 2q \Leftrightarrow$$

$$f(p) = -2q + 3(p-uv)(u-v) + 2q \Leftrightarrow$$

$$f(p) = 3(p-uv)(u-v) = 0$$

$$uv = (\sqrt{p^3 + q^2} - q)^{1/3} \cdot (\sqrt{p^3 + q^2} + q)^{1/3} \Leftrightarrow$$

$$uv = [(\sqrt{p^3 + q^2} - q)(\sqrt{p^3 + q^2} + q)]^{1/3} \Leftrightarrow$$

$$uv = [(\sqrt{p^3 + q^2})^2 - q^2]^{1/3} = [p^3 + q^2 - q^2]^{1/3} \Leftrightarrow$$

$$uv = p$$

β) Αν  $p^3 \gg q^2$ , ο τύπος του Cardano έχει προβλήματα ευστάθειας.

$$\gamma) u \approx \sqrt{p} \approx v$$

Άρα, έχουμε αφαίρεση σχεδόν ίσων αριθμών, οπότε έχουμε προβλήματα ευστάθειας.

$$\delta) p^3 \gg q^2$$

$$u^3 - v^3 = (u-v)(u^2 + uv + v^2)$$

$$\gamma) u-v = \frac{u^3 - v^3}{u^2 + uv + v^2} \iff$$

$$u-v = \frac{-2q}{u^2 + p + v^2}, \quad p > 0$$

Ευσταθής τρόπος.



8

## Άσκηση 1.7

α)  $n \geq 3$

$y_n = n \sin \frac{\pi}{n}$ , η ημiperίμετρος του εγγεγραμμένου στον μοναδιαίο κύκλο

κανονικού  $n$ -γώνου

$Y_n = n \tan \frac{\pi}{n}$ , η ημiperίμετρος του περιγεγραμμένου στον μοναδιαίο κύκλο

κανονικού  $n$ -γώνου.

$$y_n = \pi - \frac{\pi^3}{6n^2} + O\left(\frac{1}{n^4}\right)$$

$$Y_n = \pi + \frac{\pi^3}{3n^2} + O\left(\frac{1}{n^4}\right)$$

Προέκταση κατά Richardson:

$$\left. \begin{aligned} 2y_n &= 2\pi - \frac{\pi^3}{3n^2} + O\left(\frac{1}{n^4}\right) \\ Y_n &= \pi + \frac{\pi^3}{3n^2} + O\left(\frac{1}{n^4}\right) \end{aligned} \right\} \Rightarrow$$

$$2y_n + Y_n = 3\pi + O\left(\frac{1}{n^4}\right) \iff$$

$$z_n = \frac{2y_n + Y_n}{3} = \pi + O\left(\frac{1}{n^4}\right)$$

Άσκηση 1.12

$a \in \mathbb{R}$ ,  $a$  μεγάλος αριθμός

$$y_n = \int_0^1 \frac{x^n}{x+a} dx, \quad n=0, 1, 2, \dots$$

α)  $a > 0 \Rightarrow (y_n)$  γνήσια φθίνουσα και μηδενική (δηλαδή  $\lim_{n \rightarrow +\infty} y_n = 0$ )

$$\alpha) y_{n+1} = \int_0^1 \frac{x^{n+1}}{x+a} dx < \int_0^1 \frac{x^n}{x+a} dx = y_n$$

$$x \in (0, 1), \quad x^{n+1} < x^n$$

Άρα,  $(y_n)_{n \in \mathbb{N}_0}$  γνήσια φθίνουσα.

$$0 \leq y_n = \int_0^1 \frac{x^n}{x+a} dx \leq \int_0^1 \frac{x^n}{a} dx = \frac{1}{a} \frac{1}{n+1}$$

$$\Rightarrow 0 \leq y_n \leq \frac{1}{a} \frac{1}{n+1}$$

$$\lim_{n \rightarrow +\infty} \frac{1}{a} \frac{1}{n+1} = 0$$

$$\text{Άρα, } \lim_{n \rightarrow +\infty} y_n = 0$$

β)...

10

γ)  $y_{n-1} \rightarrow y_n$  ;

$$\gamma) y_n = \int_0^1 \frac{x^n}{x+a} dx \Leftrightarrow$$

$$y_n = \int_0^1 \frac{x^n + ax^{n-1} - ax^{n-1}}{x+a} dx \Leftrightarrow$$

$$y_n = \int_0^1 \frac{x^{n-1}(x+a) - ax^{n-1}}{x+a} dx \Leftrightarrow$$

$$y_n = \int_0^1 x^{n-1} dx - a \int_0^1 \frac{x^{n-1}}{x+a} dx \Leftrightarrow$$

$$y_n = \frac{1}{n} - ay_{n-1}, \quad n=1, 2, \dots$$

$$y_0 = \int_0^1 \frac{1}{x+a} dx = \log(x+a) \Big|_{x=0}^{x=1} \Leftrightarrow$$

$$y_0 = \log(a+1) - \log a = \log \frac{a+1}{a}$$

Άρα έχουμε:

$$\begin{cases} y_0 = \log \frac{a+1}{a} \\ y_n = \frac{1}{n} - ay_{n-1}, \quad n=1, 2, \dots \end{cases}$$

$$a \gg 1, \text{ Ευσιδήσια: } \tilde{y}_n = \frac{1}{n} - a\tilde{y}_{n-1}, \quad n=1, 2, \dots$$

Αφαιρώντας κατά μέλη παίρνουμε:

$$y_n - \tilde{y}_n = -a(y_{n-1} - \tilde{y}_{n-1})$$

$$\Rightarrow \dots y_n - \tilde{y}_n = (-a)^n (y_0 - \tilde{y}_0) \Rightarrow |y_n - \tilde{y}_n| = a^n |y_0 - \tilde{y}_0|$$

μεγαλώνει πολύ γρήγορα αστάθεια!

$$\beta) y_n = \sum_{k=0}^{n-1} (-1)^k \binom{n}{k} a^k \frac{(1+a)^{n-k} - a^{n-k}}{n-k} + (-a)^n \log \frac{a+1}{a}$$

με  $n \geq 1$ , Ευραθία;

$$y_n = \int_0^1 \frac{x^n}{x+a} dx = \int_a^{a+1} \frac{(y-a)^n}{y} dy \Leftrightarrow$$

$$y_n = \int_a^{a+1} \frac{1}{y} \sum_{k=0}^n \binom{n}{k} (-1)^k a^k y^{n-k} dy \Leftrightarrow$$

$$y_n = \sum_{k=0}^{n-1} (-1)^k \binom{n}{k} a^k \int_a^{a+1} y^{n-k-1} dy + (-1)^n a^n \int_a^{a+1} \frac{1}{y} dy \Leftrightarrow$$

$$y_n = \sum_{k=0}^{n-1} (-1)^k \binom{n}{k} a^k \frac{(1+a)^{n-k} - a^{n-k}}{n-k} + (-a)^n \log \frac{a+1}{a}$$

Σημείωση:  $(y-a)^n = (-a+y)^n = \sum_{k=0}^n \binom{n}{k} (-a)^k y^{n-k}$

Για κατάλληλο  $n$  οι όροι  $\binom{n}{k} a^k \frac{(1+a)^{n-k} - a^{n-k}}{n-k}$  είναι μεγάλοι, οπότε αναγκαστικά κάποια ενδιαμέσα αθροίσματα είναι μεγάλα (σε απόλυτη τιμή) ενώ το  $y_n$  είναι μικρό.

Ο αριθμός είναι ασταθής!

(12)

★★ δ) Ζητείται ευσταθής τρόπος για να υπολογίσουμε το  $y_{20}$

$$(*) \quad y_{n-1} = \frac{1}{a} \left( \frac{1}{n} - y_n \right)$$

Έστω ότι έχουμε μια προσέγγιση  $\tilde{y}_n$  του  $y_n$ , τότε:

$$\tilde{y}_{n-1} = \frac{1}{a} \left( \frac{1}{n} - \tilde{y}_n \right)$$

Οπότε:

$$y_{n-1} - \tilde{y}_{n-1} = \left( -\frac{1}{a} \right) (y_n - \tilde{y}_n)$$

↑  
μικρός αριθμός

κωρίζουμε ότι  $0 \leq y_{20} \leq \frac{1}{21a}$

Με την προσέγγιση  $\tilde{y}_{20} = 0$ , το σφάλμα είναι το πολύ  $\frac{1}{21a}$

Εφαρμόζοντας τον αλγόριθμο (\*) βρίσκουμε τις προσεγγίσεις  $\tilde{y}_{19}, \tilde{y}_{18}, \dots, \tilde{y}_{10}$

Άσκηση 1.13

$$\left. \begin{aligned} x + y &= 1 \\ x + (1-a)y &= 0 \end{aligned} \right\}$$

$a \in \mathbb{R}$ , Κατάσταση του συστήματος;

$a = 0$ :

Το σύστημα δεν έχει λύση.

$a \neq 0$ :

$$\left. \begin{aligned} \tilde{x} + \tilde{y} &= 1 + \varepsilon_1 \\ \tilde{x} + (1-a)\tilde{y} &= \varepsilon_2 \end{aligned} \right\}$$

⊖ Έτσι, έχουμε  $u := \tilde{x} - x$  και  $v := \tilde{y} - y$   
και έχουμε:

$$\left. \begin{aligned} u + v &= \varepsilon_1 \\ u + (1-a)v &= \varepsilon_2 \end{aligned} \right\}$$

$$\Rightarrow u = \varepsilon_1 - \frac{\varepsilon_1 - \varepsilon_2}{a}, \quad v = \frac{\varepsilon_1 - \varepsilon_2}{a}$$

Για  $|a|$  μικρή, η κατάσταση του συστήματος είναι κακή!

Για  $|a|$  μεγάλη, η κατάσταση του συστήματος είναι καλή!