

Teoria błędów

1. Błąd względny i bezwzględny:

$$x = 1.23456789$$

$$\bar{x} = 1.2345679$$

$$y = 10^{50} + 1$$

$$\bar{y} = 10^{50}$$

$$|x - \bar{x}| = 10^{-8}$$

$$|y - \bar{y}| = 1$$

← błędy bezwzględne

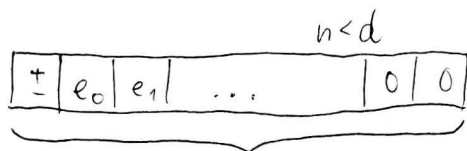
$$\left| \frac{x - \bar{x}}{x} \right| = 0.8 \cdot 10^{-8}$$

$$\left| \frac{y - \bar{y}}{y} \right| = 10^{-50}$$

← błędy względne

2. Reprezentacja liczb:

a) całkowite: $l \in \mathbb{Z} \Rightarrow l = \pm \sum_{i=0}^n e_i 2^i$ ($e_i \in \{0, 1\}$, $e_n = 1$)



$d+1$ bitów ma liczby całkowite ze znakiem

▼ PROBLEM: $n > d \Rightarrow$ nadmiar / błąd reprezentacji ▼

Dodawanie, odejmowanie i mnożenie jest wykonywane w komputerze dokładnie, pod warunkiem, że wynik jest reprezentowalny. Problemy pojawiają się przy dzieleniu, bo wynik jest liczbą rzeczywistą.

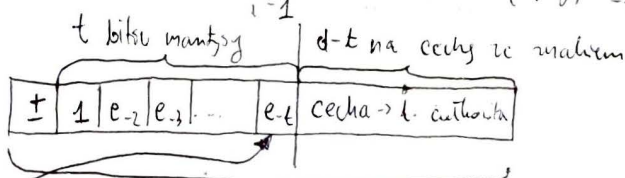
b) liczniki: każdą liczbę można jednoznacznie przedstawić w postaci:

$$x = s \cdot m \cdot 2^c$$

gdzie $s \in \{+1, -1\}$ (znak), $m \in [\frac{1}{2}, 1)$ (mantyza), $c \in \mathbb{Z}$ (cecha)

▼ PROBLEM: wykonywanie dzielenia jest precyzyjnie niedokładne ▼

$$m = \sum_{i=1}^{\infty} e_i 2^{-i} \quad (e_i \in \{0, 1\}, e_{-1} = 1)$$



maksymalna liczba bitów mantyzy

$d+1$ bitów na liczbę rzeczywistą ze znakiem

→ ZACHOWANIE

* Rodzaje zaokrągleń:

(a) obcięcie $\Rightarrow m \approx m_t^c := \sum_{i=1}^t e_{-i} 2^{-i}$ (ignorujemy cyfry od $t+1$)
(chopping)

(b) zaokrąglenie symetryczne $\Rightarrow m \approx m_t^r := \sum_{i=1}^t e_{-i} 2^{-i} + e_{-(t+1)} 2^{-t}$
(symmetric rounding)

Reprezentacja nieskończoności liczby x w narymowaniu liczby:

$\text{chop}(x) := s \cdot m_t^c \cdot 2^c$ (obcięcie)

$\text{rd}(x) := s \cdot m_t^r \cdot 2^c$ (zaokrąglenie symetryczne)

Twierdzenie:

$$\left| \frac{x - \text{chop}(x)}{x} \right| \leq 2 \cdot 2^{-t}$$

$$\left| \frac{x - \text{rd}(x)}{x} \right| \leq 2^{-t}$$

Przykład:

$$\frac{2}{3} = (0.10101010\dots)_2$$

Jeżeli $t=24$, to $\text{chop}(x) = (0.1010\dots 10)_2$
 $\text{rd}(x) = (0.1010\dots 11)_2$

Wtedy mamy: $\left| \frac{x - \text{chop}(x)}{x} \right| = 2 \cdot 2^{-25}$

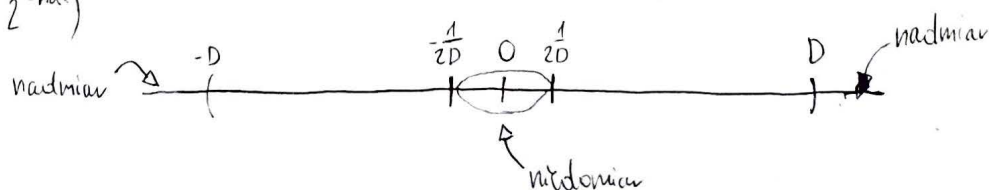
$$\left| \frac{x - \text{rd}(x)}{x} \right| = 2^{-25}$$

Jakie liczby ma komputer?

$$\frac{1}{2D} \leq |x| \leq D$$

$$\begin{cases} c \text{ echa} \in [c_{\min}, c_{\max}] = -c_{\min} = c_{\max} = 2^{d-t-1} - 1 \\ m \in [\frac{1}{2}, 1) \end{cases}$$

($D = 2^{c_{\max}}$)



W obliczeniach komputerowych mogą występować jedynie liczby:

$$X := \left(-D, -\frac{1}{2D}\right] \cup \{0\} \cup \left[\frac{1}{2D}, D\right)$$

W praktyce (nieodmiennie) zbiór dopuszczalny jest postaci:

$$X := (-D, D)$$

Zbiór liczb umiarkowanych

$$X_{fl} := \underbrace{rd(X)}_{\substack{\text{komputer na tyłko} \\ \text{liczy z tego zbiorem}}}$$

Przykład:

$$d=5$$

$$t=3$$

±	1	?	?	±	?
---	---	---	---	---	---

 $\quad z \in \{0, 1\}$

mantyza część ze znakiem

Wsp. przedstawień można: $\pm \left\{ 0, \frac{1}{4}, \frac{5}{16}, \frac{3}{8}, \frac{7}{16}, \frac{1}{2}, \frac{5}{8}, \frac{3}{4}, \frac{7}{8}, 1, \frac{5}{4}, \frac{3}{2}, \frac{7}{4} \right\}$

3. Działania arytmetyczne w ścisłej fl.

$$\left. \begin{array}{l} \cdot x, y \in X_{fl} \\ \cdot o \in \{+, -, *, /\} \\ \cdot x \circ y \in X' \text{ (nie ma nadmiarów)} \end{array} \right\} \begin{array}{l} fl(x \circ y) := (x \circ y)(1 + \varepsilon_{x,y,o}), \\ |\varepsilon_{x,y,o}| \leq 2^{-t} \end{array}$$

(t-liczba bitów na zapamiętanie mantyzy)

Błąd względny:

$$\left| \frac{x \circ y - fl(x \circ y)}{x \circ y} \right| = \left| \frac{(x \circ y) - (x \circ y)(1 + \varepsilon_{x,y,o})}{x \circ y} \right| = \underbrace{|\varepsilon_{x,y,o}|}_{\substack{\text{błąd względny operacji} \\ \circ \text{ wynosi co najwyżej } 2^{-t}}} \leq 2^{-t}$$

błąd względny operacji
o wynosi co najwyżej 2^{-t}

Twierdzenie o kumulacji błędów

Jeśli $|\delta_i| \leq 2^{-t}$ ($i=1, 2, \dots, n$), to zachodzi równość

$$\prod_{i=1}^n (1 + \delta_i) = 1 + \sigma_n$$

gdzie $\sigma_n = \sum_{i=1}^n \delta_i + O(2^{-2t})$.

Jeśli $n2^{-t} < 2$, to $|\sigma_n| \leq \gamma_n := \frac{n2^{-t}}{1 - \frac{1}{2}n2^{-t}} \approx n \cdot 2^{-t}$ (przy pomocy realistywnych rachunków)

W szczególności:

$$|\delta_i| < 2^{-t} \Rightarrow \prod_{i=1}^n (1 + \delta_i) = 1 + \sigma_n, \quad |\sigma_n| \leq n \cdot 2^{-t}$$

Przykład: Symulacja losów

Niech dane będą liczby $x_1, x_2, \dots, x_n \in \mathbb{R}$; oblicz $\sum_{i=1}^n x_i$.

PROGRAM

```

S := x1
FOR i=2 TO n
    S := S + xi
RETURN (S)
    
```

Sprawdźmy co obliczy komputer realizując podany program. Dla uproszczenia (nie zmienia to uogólnienia) założymy, że $x_i = rd(x_i)$ ($1 \leq i \leq n$):
błąd przy każdym dodawaniu

$$f_1(S) = \left(\dots \left((x_1 + x_2)(1 + \varepsilon_1) + x_3 \right)(1 + \varepsilon_2) + x_4 \right)(1 + \varepsilon_3) + \dots + x_{n-1} \right)(1 + \varepsilon_{n-2}) + x_n \underbrace{\left(1 + \varepsilon_{n-1} \right)}_{\text{błąd (n-1)szego dodawania}} =$$

$$= x_1 \overbrace{\left((1 + \varepsilon_1)(1 + \varepsilon_2) \dots (1 + \varepsilon_{n-1}) \right)}^{1 + \sigma_1} + x_2 \overbrace{\left((1 + \varepsilon_1)(1 + \varepsilon_2) \dots (1 + \varepsilon_{n-1}) \right)}^{1 + \sigma_2} + x_3 \overbrace{\left((1 + \varepsilon_2)(1 + \varepsilon_3) \dots (1 + \varepsilon_{n-1}) \right)}^{1 + \sigma_3} + \dots + x_k \overbrace{\left((1 + \varepsilon_{k-1})(1 + \varepsilon_k) \dots (1 + \varepsilon_{n-1}) \right)}^{1 + \sigma_k} + \dots + x_n \overbrace{\left(1 + \varepsilon_{n-1} \right)}^{1 + \sigma_n}$$

$$= \sum_{i=1}^n x_i (1 + \delta_i), \text{ gdzie } \begin{cases} 1 + \delta_i = \prod_{j=i-1}^{n-1} (1 + \varepsilon_j) & (i=2, 3, \dots, n) \\ 1 + \delta_1 = 1 + \varepsilon_1 \end{cases}$$

2 twierdzenie o kumulacji błędów wynika że:

$$|\delta_i| \leq (n-i+1) \cdot 2^{-t} \quad (i=2, 3, \dots, n)$$

$$|\delta_1| \leq (n-1) 2^{-t}$$

WNIOSKI:

1° Zamiast $\sum_{i=1}^n x_i$ otrzymujemy (w skutek fl) $\sum_{i=1}^n x_i (1 + \delta_i)$:

$$\left| \frac{S - fl(S)}{S} \right| = \left| \frac{\sum_{i=1}^n x_i - \sum_{i=1}^n x_i (1 + \delta_i)}{\sum_{i=1}^n x_i} \right| = \left| \frac{\sum_{i=1}^n x_i \delta_i}{\sum_{i=1}^n x_i} \right| \leq$$

$$\leq \left| \frac{\sum_{i=1}^n |x_i| |\delta_i|}{\sum_{i=1}^n |x_i|} \right| \leq \left\{ \Delta := \max |\delta_i| \right\} \leq \Delta \cdot \frac{\sum_{i=1}^n |x_i|}{\sum_{i=1}^n |x_i|}$$

to może być dodane dzie

2° Jeśli $x_i > 0$ $\Rightarrow \frac{\sum_{i=1}^n |x_i|}{\sum_{i=1}^n x_i} = 1 \Rightarrow$ wynik jest na poziomie $(n-1)2^{-t}$ (może być)

3° Jeśli $x_i > 0$, to linijki w skute fl warto sumować od najmniejszej do największej (aby minimalizować błędy).

4°

4. Zjawisko utraty cyfr znaczących

coś odgrywane linijki
przez cyfry znaczące

Operacja + ~~linijki~~ linijki tych samych znaków oraz * i / są uważane za „bezpieczne”. Kłopot występuje np. wtedy, gdy odejmujemy od siebie linijki tych samych znaków:

$$x > y > 0, x \approx y$$

$$rd(x): \begin{array}{|c|c|c|c|c|c|c|c|c|c|} \hline + & 1 & 0 & 1 & 1 & \dots & 1 & 1 & 0 & 1 \\ \hline \end{array} 2^c$$

$$rd(y): \begin{array}{|c|c|c|c|c|c|c|c|c|c|} \hline + & 1 & 0 & 1 & 1 & \dots & 1 & 0 & 0 & 0 \\ \hline \end{array} 2^c$$

> faktycznie same certy, bo $x \approx y$

$$rd(x-y): \begin{array}{|c|c|c|c|c|c|c|c|c|c|} \hline + & 0 & 0 & 0 & 0 & \dots & 0 & 1 & 0 & 1 \\ \hline \end{array} 2^c$$

musimy znormalizować linijki!

$$\begin{array}{|c|c|c|c|c|c|c|c|c|c|} \hline + & 1 & 0 & 1 & 1 & \dots & 1 & 1 & 0 & 1 \\ \hline \end{array} 2^c$$

Nie mamy co tu wpisać, to jest problem utraty cyfr znaczących.

Pytanie:

$$(1) f(x) = \ln(x) - 1, x \approx e \Rightarrow f(x) \approx 0$$

$$f(x) = \ln\left(\frac{x}{e}\right)$$

$$(2) 4038 \cdot \frac{1 - \cos x}{x^2} = 4038 \cdot \frac{1 - \cos x}{x^2} \cdot \frac{1 + \cos x}{1 + \cos x} = 4038 \cdot \frac{\sin^2 x}{(1 + \cos x)x^2}$$

↓
Jeszcze ten przekształcenie
użytkując lewy wzór, bo
unikamy problematycznego dodawania