

Zadanie 1. Jednoznaczne przedstawienie każdej niezerowej liczby rzeczywistej x postaci $x = s \cdot m \cdot B^c$, gdzie $s = \text{sgn } x$, $m \in [\frac{1}{B}, 1)$, $c \in \mathbb{Z}$, $B \in \{2, 3, 4, \dots\}$ w postaci znormalizowanej.

Każdą liczbę można przedstawić jako $s m B^c$:

(1) Skoro s to znak, to możemy go pominąć, wtedy $x = m B^c$ (z dokładnością do $|x|$).

(2) Zróbmy teraz mamy $x = m B^c \Rightarrow m = \frac{x}{B^c}$, więc dla liczby w systemie B -liczbowym (dziesiętny, trójkowy, czwórkowy etc.) możemy przesunąć je o c miejsc w lewo, a wartość c dobierzemy tak, aby $m \in [\frac{1}{B}, 1)$ - w tym celu przesunąć przesuwamy tak, aby przed nim było 0, a za nim liczba w systemie liczbowym określonym przez B .

Stąd każda liczba rzeczywista x ma przedstawienie w tym systemie.

Jednoznaczność reprezentacji x :

(1) Zakładamy, że ~~nie~~ istnieją dwie różne reprezentacje $s m B^c$ tej samej liczby. Znale jest stały, więc go pomijamy:

$$x = s m_1 B^{c_1} = s m_2 B^{c_2} \Rightarrow m_1 B^{c_1} = m_2 B^{c_2} \Rightarrow \frac{m_1}{m_2} = B^{c_2 - c_1}$$

(2) Rozpatujemy przypadki (doprowadzamy do sprzeczności)

$$c_1 = c_2 \rightarrow \text{wtedy } B^0 = 1 = \frac{m_1}{m_2} \Rightarrow m_1 = m_2$$

$$c_1 > c_2 \rightarrow \text{wtedy } B^{c_2 - c_1} \in B^{-1} \Rightarrow \frac{m_1}{m_2} \leq \frac{1}{B} \Rightarrow m_1 \leq \frac{m_2}{B},$$

jednak $m_1, m_2 \in [\frac{1}{B}, 1)$, więc gdyby m_1 było w dobrym zakresie, to musiałoby zachodzić $m_2 \geq 1$.

$$c_1 < c_2 \rightarrow \text{wtedy } B^{c_2 - c_1} \in B^1 \Rightarrow \frac{m_1}{m_2} \geq B \Rightarrow m_1 \geq m_2 B,$$

$m_1, m_2 \in [\frac{1}{B}, 1)$, więc gdyby m_1 było w dobrym przedziale, to m_2 musiałoby być $m_2 \leq \frac{1}{B}$.

Zadanie 2. Wszystkie liczby zmiennopozycyjne postaci:

$$x = \pm (0.1e_2e_3e_4)_2 \cdot 2^{\pm c}; \quad e_i, c \in \{0, 1\}$$

| 2^{-1} | 2^0 | 2^1 |
|----------|--------|--------|
| 0.01000 | 0.1000 | 1.0000 |
| 0.01001 | 0.1001 | 1.0010 |
| 0.01010 | 0.1010 | 1.0100 |
| 0.01011 | 0.1011 | 1.0110 |
| 0.01100 | 0.1100 | 1.1000 |
| 0.01101 | 0.1101 | 1.1010 |
| 0.01110 | 0.1110 | 1.1100 |
| 0.01111 | 0.1111 | 1.1110 |

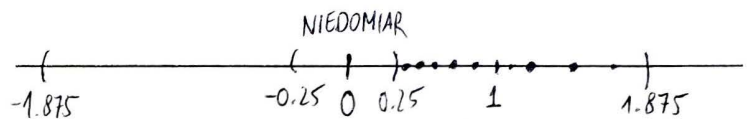
24 liczby dodatnie i tyle
samo liczb ujemnych, więc
łącznie mamy 48 takich liczb.

Najmniejsza / największa wartość co do modułu.

$$\min(|x|) = 0.01000_2 = 0.25$$

$$\max(|x|) = 1.1110_2 = 1.875$$

Stąd mamy, że $x \in [-1.875, 1.875]$,
a niedomiar dla $x \in [-0.25, 0.25]$.



Zadanie 3. Wykazać, że $\frac{|rd(x) - x|}{|x|} \leq 2^{-t}$ dla $x = sm2^c$.

Informacje: $s = \text{sgn } x$, $c \in \mathbb{Z}$, $m \in [\frac{1}{2}, 1)$, $rd(x) = sm_t^r 2^c$, $m_t^r \in [\frac{1}{2}, 1)$,

$$|m - m_t^r| \leq \frac{1}{2} \cdot 2^{-t}$$

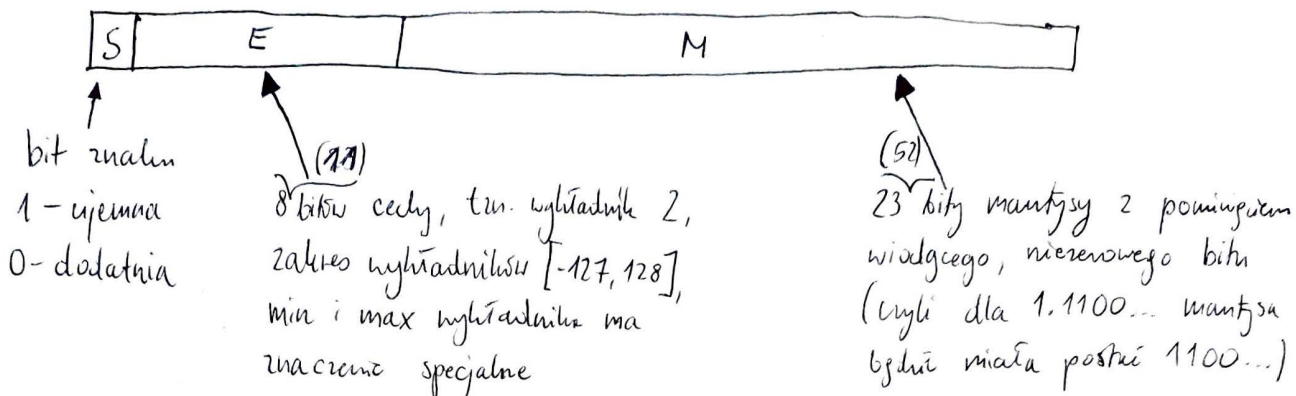
Dowód:

$$\begin{aligned} \frac{|rd(x) - x|}{|x|} &= \frac{|sm_t^r 2^c - sm 2^c|}{|sm 2^c|} = \frac{|(s2^c)(m_t^r - m)|}{|s2^c \cdot m|} = \frac{|m_t^r - m|}{|m|} \leq \\ &\leq \frac{\frac{1}{2} \cdot 2^{-t}}{m} = \frac{2^{-t}}{2m} \leq 2^{-t} \quad \blacksquare \end{aligned}$$

maksymalnej wartości, a skoro $m \in [\frac{1}{2}, 1)$,
to $2m \geq 1$, a dla każdego m większego
od $\frac{1}{2}$ wartość maleje

Zadanie 4. IEEE 754 - standard reprezentacji binarnej i operacji na liczbach rzeczywistych

32 bity (64)



Dokładność $\sim 7-8$ miejsc dziesiętnych, zakres od około $\pm 1.18 \cdot 10^{-38}$ do około $\pm 3.4 \cdot 10^{38}$.

Szczególne przypadki:

- wykładnik jest $\begin{cases} \bullet +0 & - \text{wszystkie bity są zerami} \\ \bullet -0 & - \text{bit znaku ustalony, reszta jest zerami} \end{cases}$
w kodzie z zerami lub zero w kodzie binarnym
- liczby małe (niedominar) - mantysa równa od zera
- $\pm \infty$ - ustawiane wszystkie bity wykładnika (128 w kodzie z nadmiarem lub 255 w kodzie binarnym), mantysa równa zero, może się pojawić jako wynik dzielenia przez 0.
- NaN - Not a Number, wykładnik jak wyżej, może się pojawić jako wynik pierwiastkowania liczby ujemnej

Zadanie 5. x, y - liczby maszynowe, pokażać, że algorytm obliczający

$d := \sqrt{x^2 + y^2}$ postać: $u := x^*x$, $u := u + y^*y$, $d := \text{sqrt}(u)$
może powodować zjawisko nadmiaru.

Niech $X_{fl} \in 2^{32}$, wtedy $x = y = 2^{30}$. Wtedy $\sqrt{x^2 + y^2} = \sqrt{(2^{30})^2 + (2^{30})^2} = \sqrt{2^{60} + 2^{60}} = \sqrt{2^{60}(1+1)} = 2^{30}\sqrt{2} \in X_{fl}$, jednak $2^{60} \notin X_{fl}$. Aby

tam uproszczyć, przekształcamy wzór (dla $x \geq y$, w razie potrzeby swapujemy):

$$\sqrt{x^2 + y^2} = \sqrt{x^2 \left(1 + \frac{y^2}{x^2}\right)} = |x| \cdot \sqrt{1 + \left(\frac{y}{x}\right)^2}$$

Skoro $x \geq y$,

to $\sqrt{1 + \left(\frac{y}{x}\right)^2} \leq 2$,

wg $\sqrt{2} \cdot \max(|x|, |y|) \in X_{fl}$.

Długość euklidesowa: $\|x_n\| = \sqrt{x_1^2 + x_2^2 + \dots + x_n^2}$

Wisec optymalizujemy w następujący sposób (zakładając $x_i \geq x_{i+1}$):

$$\|x_n\| = \sqrt{x_1^2 + x_2^2 + x_3^2 + \dots + x_n^2} =$$

$$= \sqrt{x_1^2 \left(1 + \frac{x_2^2}{x_1^2} + \frac{x_3^2}{x_1^2} + \dots + \frac{x_n^2}{x_1^2} \right)}$$

$$= |x_1| \cdot \sqrt{1 + \underbrace{\frac{x_2^2}{x_1^2}}_{\wedge \atop 1} + \underbrace{\frac{x_3^2}{x_1^2}}_{\wedge \atop 1} + \dots + \underbrace{\frac{x_n^2}{x_1^2}}_{\wedge \atop 1}}$$

$$= |x_1| \cdot \sqrt{n} \equiv \max(x_1, x_2, \dots, x_n) \cdot \sqrt{n}$$

Zadanie 6. Poprawienie $f(x) = 4038 \cdot \frac{\sqrt{x^{11}+1} - 1}{x^{11}}$, aby
mieścił się w przedziale double'a.

$$f(x) = 4038 \cdot \frac{\sqrt{x^{11}+1} - 1}{x^{11}} = 4038 \cdot \frac{(\sqrt{x^{11}+1} - 1)(\sqrt{x^{11}+1} + 1)}{x^{11}(\sqrt{x^{11}+1} + 1)} =$$

$$= 4038 \cdot \frac{x^{11} + 1 - 1}{x^{11}(\sqrt{x^{11}+1} + 1)} = 4038 \cdot \frac{1}{\sqrt{x^{11}+1} + 1} \quad \blacksquare$$

Zadanie 7. $x_{k+1} = 2^k \sqrt{2 \left(1 - \sqrt{1 - \left(\frac{x_k}{2^k}\right)^2}\right)}$ $k = 1, 2, \dots; x_1 = 2$

Czy x_k jest zbliżony do π , dla $k=30$ zaczynamy tracić cyfry znaczące. Podstawiamy więc, aby temu zapobiec:

$$A = 1 - \sqrt{1 - \left(\frac{x_k}{2^k}\right)^2} = \frac{\left(1 - \sqrt{1 - \left(\frac{x_k}{2^k}\right)^2}\right) \left(1 + \sqrt{1 - \left(\frac{x_k}{2^k}\right)^2}\right)}{1 + \sqrt{1 - \left(\frac{x_k}{2^k}\right)^2}} = \frac{\left(\frac{x_k}{2^k}\right)^2}{1 + \sqrt{1 - \left(\frac{x_k}{2^k}\right)^2}}$$

Wtedy $x_{k+1} = 2^k \sqrt{2A} = \sqrt{2 \cdot \frac{x_k^2}{1 + \sqrt{1 - \left(\frac{x_k}{2^k}\right)^2}}}$

Zadanie 8. Dla jakich x słuszny cyfry znaczące?

(a) $f(x) = x^7 + \sqrt{x^{14} + 2019}$ \rightarrow złe, gdy $|x|^7 \approx \sqrt{x^{14} + 2019}$, $x < 0, |x| \gg 2019$

$$f(x) = \frac{(x^7 + \sqrt{x^{14} + 2019})(x^7 - \sqrt{x^{14} + 2019})}{x^7 - \sqrt{x^{14} + 2019}} = \frac{-2019}{x^7 - \sqrt{x^{14} + 2019}}$$

(b) $f(x) = x^{-7} \left(\sin x - x + \frac{x^3}{6} - \frac{x^5}{120} \right)$

to jest problematyczne,
bo dla małych x (do
ok. 3°) $\sin x \approx x$

$$\sin(x) = \sum_{n=0}^{\infty} \frac{(-1)^n}{(2n+1)!} x^{2n+1} = x - \frac{x^3}{6} + \frac{x^5}{120} - \frac{x^7}{7!} + \dots =$$

$$(?) = -\frac{1}{7!} + \frac{x^2}{3!} - \frac{x^4}{11!} + \dots$$