

1. Niech  $B$  będzie liczba naturalna większa od 1. Wykazac, że każda niezerowa liczba rzeczywista  $x$  ma jednoznaczne przedstawienie w postaci znormalizowanej  $x = s m B^c$ , gdzie  $s$  jest znakiem liczby  $x$ ,  $c$  – liczba całkowita (cecha), a  $m$  – liczba z przedziału  $[1, B)$ , zwana mantysa.

1. istnienie:

Niech  $0 \neq x \in \mathbb{R}$ . Wtedy istnieje  $c$  takie, że  $B^c \leq |x| < B^{c+1}$ , czyli  $c = \lfloor \log_B |x| \rfloor$ . Niech  $m = \frac{|x|}{B^c}$ , wtedy  $|x| = B^c \cdot m$ . W końcu, niech  $s = \frac{|x|}{x}$ . Mamy  $x = s B^c m$ .

2. jedynosc:

A. jedynosc  $s$  jest oczywista

B. jedynosc  $c$ : załozmy, nie wprost, że istnieja  $c_1, c_2$  takie, że

$$x = s B^{c_1} m = s B^{c_2} m'$$

Wtedy

$$B^{c_1} m = B^{c_2} m'$$

Jesli  $m = m'$  oczywiste. W przeciwnym wypadku

$$c_1 \log_B m = c_2 \log_B m'$$

mozemy zalozyc, że  $c_1 < c_2$  oraz  $(\exists k \in \mathbb{N}) c_1 + k = c_2$ , czyli

$$c_1 \log_B m = (c_1 + k) \log_B m'$$

$$0 = k \log_B m'$$

W takim razie albo  $m' = 1$ , wtedy  $x = 2_1^c$ , albo  $k = 0$ , czyli  $c_1 = c_2$ .

C. jedynosc  $m$ : załozmy, nie wprost, że istnieja  $m_1, m_2$  takie, że (...), wtedy

$$x = s B^c m_1 = s B^c m_2.$$

$c$  jest jedyne, gdyż  $c = \lfloor \log_B |x| \rfloor$  i to działanie ma jednoznaczny wynik. Czyli

$$s B^c m_1 = s B^c m_2$$

$$m_1 = m_2$$



## 2. Ile jest liczb zmiennopozycyjnych w arytmetyce double w standardzie IEEE754?

Przy 64 bitach mamy  $2^{64}$  mozliwosci ich zapalenia. Liczby NaN to liczby majace wszystkie bity w mantysie zapalone, a takich jest  $2^{53}$ . To daje nam  $2^{64} - 2^{53}$  liczb, ale 0 jest reprezentowane na dwa sposoby, wiec wystarczy  $2^{64} - 2^{53} - 1$ .

Liczby subnormalne maja na pierwszym miejscu mantysy 0, podczas gdy cala reszta zaczyna mantysę od 1 i one juz sie wliczaja.

## 3. Czesc rozw w pliku .jl

```
1 function frst_exp(x, s, t, r)
2   ret = zero(x)
3   ret = (x^3) - (s*(x^2)) + t*x - r
4   print("_", typeof(x), "_wynik:", ret, "\n")
5 end
6
```

```

7 function snd_exp(x, s, t, r)
8     ret = zero(x)
9     ret = ((x - s) * x + t) * x - r
10    print("_", typeof(x), "_wynik:", ret, "\n")
11 end
12
13 function rel_error(val, exp, mess)
14     bez = zero(val)
15     if val > exp
16         bez = val - exp
17     else
18         bez = exp - val
19     end
20
21     println("Błąd względny", mess, "_wynosi:", bez / exp)
22 end

```

-14.636489 - wartosc prawdziwa

Float16	-0.003987568330082385	-0.0002511872895199931
Float32	-7.760455081662309e-7	-5.931504907397523e-8
Float64	-4.85459822885188e-16	0

## Zad 4. Za długa tresc

$$\sum_{k=t+2}^{\infty} \frac{1}{2^k} = 2 \frac{1}{2^{t+2}} = 2^{-t-1}$$

$$\begin{aligned}
 |\text{rd}(x) - x| &= |s(1 + \sum_{k=1}^{\infty} e_{-k} 2^{-k}) 2^c - s(1 + \sum_{k=1}^t e_{-k} 2^{-k} + e_{-t-1} 2^{-t}) 2^c| = \\
 &= 2^c \left| \sum_{k=1}^{\infty} e_{-k} 2^{-k} - \sum_{k=1}^t e_{-k} 2^{-k} - e_{-t-1} 2^{-t} \right| = \\
 &= 2^c \left| \sum_{k=t+1}^{\infty} e_{-k} 2^{-k} - e_{-t-1} 2^{-t} \right| \leq \\
 &\leq 2^c \left| \sum_{k=t+1}^{\infty} 2^{-k} - e_{-t-1} 2^{-t} \right| = \\
 &= 2^c \left| \sum_{k=t+2}^{\infty} 2^{-k} + e_{-t-1} 2^{-t-1} - e_{-t-1} 2^{-t} \right| = \\
 &= 2^c \left| \sum_{k=t+2}^{\infty} 2^{-k} + e_{-t-1} \left( \frac{1}{2^{t+1}} - \frac{1}{2^t} \right) \right| = \\
 &= 2^c \left| \frac{1}{2^{t+1}} - \frac{e_{-t-1}}{2^{t+1}} \right| \leq \\
 &\leq 2^c \cdot 2^{-t-1} = 2^c \cdot u
 \end{aligned}$$

$$\frac{|\text{rd}(x) - x|}{|x|} = \frac{2^c |m - \bar{m}|}{2^c |m|} = \frac{\left| \sum_{k=t+1}^{\infty} e_{-k} 2^{-k} - \dots \right|}{\left| 1 + \sum_{k=1}^{\infty} e_{-k} 2^{-k} \right|} \leq \frac{u}{1} = u$$

## 5. nah

$$\begin{aligned}
 \frac{|\bar{m} - m|}{|m|} &\leq \frac{u}{1 + u} \\
 (1 + u) |\bar{m} - m| &\leq u |m|
 \end{aligned}$$

Z poprzedniego zadania wiemy, że  $|\bar{m}-m| \leq u$ . Wystarczy więc, że rozważymy dwa przypadki:

1.  $(1+u) \leq m$  – oczywiste
2.  $(1+u) > m$

$$\begin{aligned} u &> m - 1 \\ \bar{m}u &> m\bar{m} - \bar{m} > m - \bar{m} \\ \bar{m}u + mu &> m - \bar{m} + mu \\ mu &> m - \bar{m} + mu - \bar{m}u \\ mu &> m(1+u) - \bar{m}(1+u) \\ mu &> (1+u)(m - \bar{m}) \\ \frac{u}{1+u} &> \frac{m - \bar{m}}{m} \end{aligned}$$

## 6. w pliiikuuuuuu

```
1 using Base
2
3 function to_number(str::AbstractString)
4     if length(str) <= 64
5         bias = Float64(2^10 - 1)
6
7         ret = Float64(0)
8         cech = Float64(0) # 11
9         man = Float64(1) # 53
10
11         c = Float64(1)
12
13         for i = 1:11
14             cech += c * (Float64(str[13 - i]) - Float64(48))
15             c *= Float64(2)
16         end
17
18         c = Float64(0.5)
19
20         for i = 1:52
21             man += c * (Float64(str[i+12]) - Float64(48))
22
23             c /= Float64(2)
24         end
25
26         cech -= bias
27
28         ret = man * (2 ^ cech)
29
30         if str[1] == '1'
31             ret *= Float64(-1)
32         end
33
34         return ret
35     else
36         print("this is not a representation of a Float64")
37         return NaN
38     end
39 end
```

## 7.

Liczba 1.000000057228997

```
1 function fl()
```

