

Ćwiczenia z ANALIZY NUMERYCZNEJ (M)

Blok 1: lista M1

10 października 2019 r.

M1.1. 1 punkt Niech B będzie liczbą naturalną większą od 1. Wykazać, że każda niezerowa liczba rzeczywista x ma jednoznaczne przedstawienie w postaci *znormalizowanej* $x = smB^c$, gdzie s jest znakiem liczby x , c – liczbą całkowitą (*cechą*), a m – liczbą z przedziału $[1, B)$, zwaną *mantysą*.

M1.2. 1 punkt Zapoznać się ze standardem IEEE 754 (zob. np. http://en.wikipedia.org/wiki/IEEE_754) Ile jest liczb zmiennopozycyjnych w arytmetyce *single*, a ile w arytmetyce *double* w tym standardzie?

M1.3. 1 punkt Obliczyć wartość $w(x) = x^3 - 6x^2 + 3x - 0.149$ w punkcie $x = 4.71$ używając arytmetyki Float16, Float32 i Float64 w języku Julia. Podać błąd względny wyniku, biorąc pod uwagę wartość dokładną $w(4.71) = -14.636489$. Powtórzyć obliczenia dla równoważnego wyrażenia $w(x) = ((x - 6)x + 3)x - 0.149$. Porównać wyniki.

Podczas prezentacji należy przedstawić plik źródłowy, np. na wydruku.

M1.4. 1 punkt Dla danych: naturalnej liczby t oraz niezerowej liczby rzeczywistej $x = sm2^c$, gdzie s jest znakiem liczby x , c – liczbą całkowitą, a m – liczbą z przedziału $[1, 2)$, o rozwinięciu dwójkowym $m = 1 + \sum_{k=1}^{\infty} e_{-k}2^{-k}$, w którym $e_{-k} \in \{0, 1\}$ dla $k \geq 1$, definiujemy *zaokrąglenie liczby* x do $t + 1$ cyfr za pomocą wzoru

$$\text{rd}(x) := s\bar{m}2^c,$$

gdzie $\bar{m} = 1 + \sum_{k=1}^t e_{-k}2^{-k} + e_{-t-1}2^{-t}$.

Wykazać, że

$$|\text{rd}(x) - x| \leq 2^cu,$$

gdzie $u := 2^{-t-1}$ jest *precyzją arytmetyki*.

Wynioskować stąd, że błąd względny zaokrąglenia liczby x nie przekracza precyzji arytmetyki u .

M1.5. 1 punkt Załóżmy, że $|\alpha_j| \leq u$ i $\rho_j \in \{-1, +1\}$ dla $j = 1, 2, \dots, n$ oraz że $nu < 1$, gdzie $u := 2^{-t-1}$. Wykazać, że zachodzi równość

$$\prod_{j=1}^n (1 + \alpha_j)^{\rho_j} = 1 + \theta_n,$$

gdzie θ_n jest wielkością spełniającą nierówność $|\theta_n| \leq \gamma_n$, gdzie z kolei

$$\gamma_n := \frac{nu}{1 - nu}.$$

M1.6. 1 punkt Napisać w języku Julia funkcję odwrotną do funkcji bibliotecznej `bitstring(...)`, tzn. która dla danego słowa s (łańcuch 64 znaków '0' lub '1') oblicza liczbę rzeczywistą x typu Float64. *Wystarczy, aby program działał dla słów maszynowych reprezentujących liczby normalne.*

M1.7. 1 punkt Znaleźć liczbę maszynową x (`double`, w standardzie IEEE 754) z przedziału $(1, 2)$, dla której $\text{fl}(x \cdot \text{fl}(1/x)) \neq 1$.

3 października 2019 r.

Rafał Nowak