

Obliczenia Naukowe

Lista 1

Paweł Polerowicz

254626

Październik 2021

Zadanie 1.

Zadanie składało się z kilku podpunktów, każdy odpowiadający oddzielnemu programowi. Każdy z nich wyznaczał pewną liczbę dla typów Float16, Float32, Float64.

1. Celem pierwszego programu było znalezienie *macheps*, czyli najmniejszej liczby takiej, że $macheps > 0$, $fl(1.0+macheps) > 1.0$ i $fl(1.0+macheps) = 1 + macheps$. W tym celu zastosowano podejście iteracyjne, dodając do 1.0 stopniowo coraz mniejszą liczbę i sprawdzając, czy nadal spełnione są warunki z założeń. Ostatnia liczba spełniająca warunki to *macheps*.

typ	<i>macheps</i>	eps(typ)	float.h
Float16	0.000977	0.000977	-
Float32	1.1920929×10^{-7}	1.1920929×10^{-7}	1.1920929×10^{-7}
Float64	$2.220446049250313 \times 10^{-16}$	$2.220446049250313 \times 10^{-16}$	$2.220446049250313 \times 10^{-16}$

Dane są spójne dla wszystkich trzech źródeł, z zastrzeżeniem, że float.h nie zawiera informacji o epsilon dla 16 bitowego typu. Ponadto, możemy zauważyć $macheps = 2 * \epsilon$ (precyzja arytmetyki)

typ	ϵ
Float16	4.88×10^{-4}
Float32	5.96×10^{-8}
Float64	1.11×10^{-16}

2. Kolejny program wyznacza liczbę η , czyli najmniejszą liczbę taką, że $\eta > 0$. Ponownie zastosowano podejście iteracyjne, zaczynając od 1 i dzieląc w pętli przez 2 dopóki warunek był spełniony.

Znalezione wartości pozwalają wysnuć wniosek: $\eta \approx Min_{sub}$

typ	η	nextfloat(typ(0.0))	Min_{sub}
Float16	5.96×10^{-8}	5.96×10^{-8}	
Float32	1.40×10^{-45}	1.40×10^{-45}	1.4×10^{-45}
Float64	4.94×10^{-324}	4.94×10^{-324}	4.9×10^{-325}

3. Sprawdziliśmy również wyniki zwracane przez funkcję floatmin

Znalezione wartości Min_{nor} są mniej dokładne, ale pozwalają wysnuć wniosek: $floatmin \approx Min_{nor}$

Zadanie 5.

Naszym zadaniem było obliczenie na różne sposoby iloczynu skalarnego wektorów:

$$x = [2.718281828, -3.141592654, 1.414213562, 0.5772156649, 0.3010299957]$$

$$y = [1486.2497, 878366.9879, -22.37492, 4773714.647, 0.000185049]$$

typ	"w przód"	"w tył"	najw. - najmn.	najmn. - najw.
Float32	-0.4999443	1.1920929×10^{-7}	-0.5	-0.5
Float64	$1.0251881368296672 \times 10^{-10}$	$-1.5643308870494366 \times 10^{-10}$	0.0	0.0

Niestety, wszystkie metody dają wyniki wyraźnie różne od dokładnej wartości $-1.00657107000000 \times 10^{-11}$. W przypadku dwóch ostatnich metod wynika to prawdopodobnie z faktu, że dodajemy do siebie duże sumy częściowe o przeciwnych znakach, bliskie co do wartości bezwzględnej. W pierwszych dwóch metodach dodajemy do siebie liczby o dużej i małej wartości bezwzględnej, co również wprowadza niedokładność.

Zadanie 6.

Naszym zadaniem było policzenie wartości funkcji

$$f(x) = \sqrt{x^2 + 1} - 1$$

$$g(x) = x^2 / (\sqrt{x^2 + 1} + 1)$$

dla $x = 8^{-1}, 8^{-2}, \dots$. Okazało się, że funkcje f i g dają zbliżone wyniki dla wykładnika i : $i \geq -8$. Dla mniejszych i f daje wynik 0, natomiast g daje wiarygodne wyniki aż do $i = 178$. Wynika to z redukcji cyfr znaczących przy odejmowaniu w funkcji f

Zadanie 7.

Zadanie polegało na obliczeniu przybliżonej wartości pochodnej wyrażenia $f(x) = \sin x + \cos 3x$ za pomocą wzoru $f'(x) \approx \frac{f(x+h) - f(x)}{h}$ dla $h = 2^{-n}$ ($n = 0, 1, \dots, 54$). Dla przejrzystości przedstawimy tylko wybrane wartości n .

n	przybliżenie	błąd
5	0.24344307439754687	0.1265007927090087
10	0.12088247681106168	0.0039401951225235265
20	0.11694612901192158	$3.8473233834324105 \times 10^{-6}$
27	0.11694231629371643	$3.460517827846843 \times 10^{-8}$
28	0.11694228649139404	$4.802855890773117 \times 10^{-9}$
29	0.11694222688674927	$5.480178888461751 \times 10^{-8}$
40	0.1168212890625	0.0001209926260381522
54	0.0	0.11694228168853815

Najmniejszy błąd uzyskano dla $n = 28$. Powyżej tej wartości musimy liczyć się z utratą precyzji przy odejmowaniu $f(x+h) - f(x)$, gdyż ich wartości są bliskie sobie, co wpływa na dokładność całego wyniku.