

Sprawozdanie

OBLICZENIA NAUKOWE, LISTA NR 1

ŁUKASZ KLEKOWSKI 229738

1. Zadanie pierwsze

1.1. Opis problemu

Polecenie polegało na napisaniu w języku Julia programów iteracyjnie obliczających epsilon maszynowego (*macheps*), najmniejszej dodatniej liczby rzeczywistej (*eta*) oraz wartości maksymalnej dla trzech liczb typów zmiennopozycyjnych : *Float16*, *Float32*, *Float64*. Należało także porównać otrzymane wartości ze stałymi ze środowiska Julia: *eps()*, *nextfloat(0.0)* oraz *realmax()*. Epsilon maszynowy oraz wartości maksymalne należało porównać także z danymi zawartymi w pliku nagłówkowym *float.h* w dowolnej instalacji języka C.

1.2. Rozwiązanie problemu

Do obliczenia zera maszynowego użyłem pętli *while* która dzieliła liczbę 1.0 dopóki nie był spełniony dany warunek:

$$1.0 + \frac{\text{macheps}}{2} = 1.0$$

Podobnie wyliczyłem najmniejszą dodatnią liczbę rzeczywistą. Dzieliłem w pętli liczbę 1.0 do momentu gdy:

$$\frac{\text{eta}}{2} \leq 0.0$$

Natomiast do obliczenia największej wartości użyłem pętli która mnożyła liczbę 1.0 przez 2 do momentu gdy funkcja *isinf()*, która sprawdza czy dana liczba to nieskończoność zwróci *true*. Następnie przy pomocy funkcji *prevfloat()* odczytujemy poprzednią wartość floata, a w naszym przypadku jest to wartość maksymalna.

1.3. Wyniki

1.3.1. Macheps

Julia

	<i>machepsFloat()</i>	<i>eps()</i>
<i>Float16</i>	9.765625e-4	9.765625e-4
<i>Float32</i>	1.1920928955078125e-7	1.1920928955078125e-7
<i>Float64</i>	2.22044604925031308e-16	2.22044604925031308e-16

Język C

	<i>FLT_EPSILON</i> / <i>DBL_EPSILON</i>
<i>Float (Float32)</i>	1.1920928955078125e-7
<i>Double (Float64)</i>	2.2204460492503131e-16

Epsilon obliczone w pętlach są takie same jak te uzyskane za pomocą funkcji *eps()*. W pliku nagłówkowym *float.h* nie ma informacji na temat epsilon dla arytmetyki *Float16*, natomiast dla *Float32* epsilon jest taki sam jak w Julii, a dla *Float64* jest trochę zaokrąglone.

1.3.2. Eta

	eta()	nextfloat(0.0)
Float16	5.9604644775390625e-8	5.9604644775390625e-8
Float32	1.40129846432e-45	1.40129846432e-45
Float64	4.9e-324	4.9e-324

Obliczone eta zgadza się z liczbą uzyskaną za pomocą funkcji nextfloat(0.0)

1.3.3. MAX

Julia

	inf()	realmax()
Float16	65504.0	65504.0
Float32	3.4028234664e38	3.4028234664e38
Float64	1.7976931348623157e308	1.7976931348623157e308

Język C

	FLT_MAX / DBL_MAX
Float	3.402823466385288600000000000000e38
Double	1.797693134862315700000000000000e308

Wyniki z pętli są takie same jak liczby otrzymane z funkcji realmax(). W pliku nagłówkowym float.h również nie ma informacji co do Float16. Natomiast Float i Double są bardzo zaokrąglone w porównaniu do wyników otrzymywanych w Julii.

1.4. Wnioski

W zadaniu było zawarte pytanie jaki związek ma liczba *macheps* z precyzją arytmetyki oraz jaki związek ma liczba eta z liczbą MIN_{sub} .

Precyzja arytmetyki oraz epsilon maszynowy jest to maksymalny błąd względny, który może zajść podczas zaokrąglania liczby rzeczywistej do liczby zmiennoprzecinkowej.

Liczba MIN_{sub} jest tym samy co epsilon czyli najmniejszą możliwą liczbą większą od 0.

2. Zadanie drugie

2.1. Opis zadania

Polecenie polegało na sprawdzeniu słuszności stwierdzenia Kahana, czy epsilon maszynowy można uzyskać za pomocą wzoru $3 * \left(\frac{4}{3} - 1\right) - 1$ w precyzji Float16, Float32, oraz Float64.

2.2. Rozwiązanie

Napisałem w języku Julia funkcję która liczy podane wyrażenie

2.3. Wyniki

Float16	Float32	Float64
-9.765625e-4	1.1920928955078125e-7	-2.22044604925031308e-16

2.4.Wnioski

Nie otrzymaliśmy zera choć teoretycznie powinniśmy. Komputer nie jest w stanie dokładnie zapisać liczby $4/3$ więc musi zaokrąglić ją do najbliższej liczby w arytmetyce zmiennopozycyjnej. Powstał błąd to właśnie epsilon maszynowy który może powstać podczas zaokrąglania liczby.

3. Zadanie trzecie

3.1.Opis zadania

Polecenie polegało na sprawdzeniu rozmieszczenia liczb zmiennopozycyjnych w precyzji Float64 dla przedziałów [0.5-1], [1-2] [2-4].

3.2.Opis rozwiązania i wyniki

Użyłem tutaj funkcji `bits()` dzięki której dostałem binarną reprezentację danej liczby. W pętli dla danych przedziałów zwiększałem liczbę o dany krok i sprawdzałem jak będzie wyglądała ich reprezentacja binarna, jednocześnie porównując je z `nextfloat()`

[1-2]

[illegible]

[0.5-1]

[illegible]

[2-4]

[illegible]

3.3. Wnioski

Z przeprowadzonego eksperymentu wynika że liczby w przedziale [1-2] są rozłożone równomiernie z krokiem 2-52, w przedziale [0.5-1] z krokiem 2-53, a w przedziale [2-4] z krokiem 2-51.

Wynika to oczywiście z charakterystyki tej arytmetyki ponieważ liczb w takich przedziałach jest tyle samo, a te przedziały się zwiększają więc są one rozłożone bardziej rzadko lub gęściej.

4. Zadanie czwarte

4.1.Opis zadania

Polecenie polegało na znalezienie najmniejszej takiej liczby w przedziale (1-2) w precyzji Float64 która spełnia warunek $x*x!=1$

4.2.Rozwiązanie

W pętli do liczby 1.0 był dodawany epsilon maszynowy i następnie za pomocą powstałej liczby liczony był dany wzór. Gdy wynik był różny od 1 to dana liczba jest drukowana na ekranie.

4.3.Wyniki

Najmniejsza liczba która nie spełnia tego równania to :
1.000000057228997

4.4.Wnioski

Takie liczby są znajdowane dlatego, że komputer nie potrafi zapisać dokładnie każdej liczby w postaci $1/x$ w arytmetyce zmiennopozycyjnej.

5. Zadanie piąte

5.1.Opis zadania

Polecenie polegało na obliczeniu iloczynu skalarnego dwóch wektorów za pomocą czterech algorytmów używając precyzji Float32 oraz Float64.

5.2.Rozwiązanie

Pierwszy algorytm polegał na zwykłym obliczeniu produktu skalarnego sumując po kolei iloczyny. Drugi algorytm robił to w odwrotnej kolejności. Trzeci algorytm polegał na dodawaniu wyników mnożenia składowych wektora od najmniejszej wartości do największej, a czwarty algorytm robił to na odwrót.

5.3.Wyniki

Algorytm	Float32	Float64	Błąd względny F32	Błąd względny F64
I	-4.9994429945e-1	1.02518813e-10	4.996159e10	1.1245e1
II	-4.5434570312500e-1	-1.56433088e-10	4.540473e10	1.463e1
III	-5.000e-1	0	4.996716e10	-
IV	-5.000e-1	0	-4.9967165	-

5.4.Wnioski

Żaden z algorytmów nie dał wyniku choć przybliżonego do wartości poprawnej. W przypadku precyzji Float32 wyniki są bardzo dalekie od rzeczywistych. W przypadku precyzji Float64 pierwszy oraz drugi algorytm dawał wyniki bliskie do wyniku precyzyjnego, a trzeci oraz czwarty algorytm dawał wartości 0. Okazuje się że kolejność wykonywania działań w komputerze ma ogromne znaczenie.

6. Zadanie szóste

6.1. Opis zadania

Polecenie polegało na obliczeniu wartości dwóch funkcji dla argumentów 8^{-1} , 8^{-2} itd.

6.2. Wyniki

x	f(x)	g(x)
8^{-1}	0.0077822185373186414381053	0.0077822185373187064902356
8^{-2}	0.0001220628628286757333399	0.0001220628628287590136193
8^{-3}	0.0000019073468138230964541	0.0000019073468138265659011
8^{-4}	0.0000000298023219436061026	0.0000000298023219436061159
8^{-5}	0.0000000004656612873077393	0.0000000004656612871993190
8^{-6}	0.000000000072759576141834	0.000000000072759576141570
8^{-7}	0.000000000001136868377216	0.000000000001136868377216
8^{-8}	0.0000000000000017763568394	0.0000000000000017763568394
8^{-9}	0.000000000000000000000000	0.000000000000000277555756
8^{-10}	0.000000000000000000000000	0.00000000000000004336809

6.3. Wnioski

Po sprawdzeniu otrzymanych wyników w pakiecie matematycznym wolfram można stwierdzić że wyniki otrzymane za pomocą funkcji g(x) są dokładniejsze. Funkcja f od argumentu 8^{-9} zwraca 0. To zadanie również pokazuje że obliczenia wykonywane w różny sposób mogą dawać różne wyniki.

7. Zadanie siódme

7.1. Opis zadania

Polecenie polegało na obliczeniu wartości pochodnej danej funkcji w punkcie $x_0=1$. Należało użyć wzoru do przybliżonej wartości pochodnej oraz porównać ją z dokładną pochodną. Obliczenia były przeprowadzane w precyzji Float64.

7.2. Wyniki

h	Przybliżona pochodna	Błąd
2^0	2.0179892252685966980152443	1.9010469435800585458196110
2^{-1}	1.8704413979316472094183155	1.7534991162431090572226822
2^{-2}	1.1077870952342974142368348	0.9908448135457592620412015
2^{-3}	0.6232412792975816628882058	0.5062989976090435106925725
2^{-4}	0.3704000662035191737686546	0.2534577845149810215730213
2^{-10}	0.1208824768110616787453182	0.0039401951225235265496849
2^{-11}	0.1189122504688384651672095	0.0019699687803003129715762
2^{-12}	0.1179272337390102620702237	0.0009849520504721098745904
2^{-13}	0.1174347496107657207176089	0.1169422816885381521956333
2^{-14}	0.1171885136209311895072460	0.0002462319323930373116127
2^{-35}	0.1169395446777343750000000	0.0000027370108037771956333
2^{-40}	0.1168212890625000000000000	0.0001209926260381521956333

