

Sprawozdanie z listy 1 na Obliczenia Naukowe

Adrian Herda, 268449

October 22, 2023

1 Rozpoznanie arytmetyki

Zadanie to polegało na obliczeniu *epsilonu maszynowego*, *ety maszynowej* oraz *MAX* dla wszystkich typów zmiennopozycyjnych, w standardzie **IEEE 754**, dostępnych w języku *Julia*. W języku tym mamy dostępne 3 takie typy: Float16, Float32, Float64. Typy te reprezentują precyzje 16-, 32- oraz 64-bitowa.

- Epsilon maszynowy - najmniejsza liczba spełniająca równanie $fl(1.0 + macheps) > 1.0$, która da radę reprezentować w arytmetyce zmiennopozycyjnej
- Eta maszynowa - najmniejsza liczba większa od zera, która da radę reprezentować w arytmetyce zmiennopozycyjnej
- MAX - największa liczba, która można reprezentować w arytmetyce zmiennopozycyjnej

1.1 Rozwiązania

Rozwiązania polegają na iteracyjnym zmniejszaniu (lub powiększaniu w przypadku funkcji MAX) wybranej liczby aż do granicy w której nie da rady mieć bardziej precyzyjnej liczby. Wtedy liczba jest zwracana jako wynik.

1.1.1 Epsilon maszynowy

```
1 # Funkcja znajduje najmniejszą liczbę większą od 1 dla typu T i zwraca jej
   roznice od 1
2 function machine_epsilon(T)
3     result = T(2.0)
4     epsilon = T(1.0)
5
6     while one(T) + epsilon > one(T)
7         result /= T(2.0)
8         epsilon /= T(2.0)
9     end
10
11     return result
12 end
```

Listing 1: Obliczanie epsilonu maszynowego

1.1.2 Eta maszynowa

```
1 # Funkcja znajduje najmniejszą liczbę większą od 0 dla typu T
2 function machine_eta(T)
3     result = T(2.0)
4     eta = T(1.0)
5
6     while eta > zero(T)
7         result /= T(2.0)
8         eta /= T(2.0)
9     end
```

```

10
11     return result
12 end

```

Listing 2: Obliczanie ety maszynowej

1.1.3 MAX

```

1  # Funkcja pomocnicza
2  # Funkcja znajduje największą liczbę mniejszą od liczby x dla typu T
3  function find_previous_float(x, T)
4      x = T(x)
5      diff = T(1.0)
6      while x - diff < x
7          diff /= T(2.0)
8      end
9
10     diff *= T(2.0)
11
12     return x - diff
13 end
14
15 # Funkcja znajduje największą możliwą liczbę
16 function max(T)
17     max = find_previous_float(1.0, T)
18
19     while !isinf(max * T(2.0))
20         max *= T(2.0)
21     end
22
23     return max
24 end

```

Listing 3: Obliczanie MAX

1.2 Wyniki oraz ich interpretacja

1.2.1 Epsilon maszynowy

	Float16	Float32	Float64
machine_eps()	0.000977	1.1920929e-7	2.220446049250313e-16
eps()	0.000977	1.1920929e-7	2.220446049250313e-16
float.h		1.19209e-07	2.22045e-16

Wyniki pokazują że iteracyjny sposób jest dobry, a plik nagłówkowy języka C nie jest aż tak dokładny.

1.2.2 Eta maszynowa

	Float16	Float32	Float64
machine_eta()	6.0e-8	1.0e-45	5.0e-324
nextfloat(zero())	6.0e-8	1.0e-45	5.0e-324

Ponownie widzimy że iteracyjny sposób daje nam poprawne wyniki.

1.2.3 MAX

	Float16	Float32	Float64
max()	6.55e4	3.4028235e38	1.7976931348623157e308
floatmax()	6.55e4	3.4028235e38	1.7976931348623157e308
float.h		3.40282e38	1.79769e308

Po raz trzeci upewniamy się że nasz sposób wyliczania charakterystycznych liczb jest poprawny. Plik nagłówkowy języka C nie daje aż tak dokładnych wyników.

1.2.4 Precyzja arytmetyki

Precyzja arytmetyki jest dokładnością reprezentowania liczb w danym systemie liczbowym. Często jest reprezentowana przez grecką literę ϵ i w naszym przypadku jest obliczana wzorem

$$\epsilon = 2^{-t}$$

gdzie t to liczba cyfr mantysy. Po obliczeniu wartości ϵ dla typów Float16, Float32 oraz Float64 i porównaniu z obliczonymi przez nas epsilonami maszynowymi zauważymy że wyniki są takie same. Wnioskujemy zatem że ϵ jest równy epsilonowi maszynowemu.

1.2.5 MIN_{sub}

Liczba ta jest najmniejsza nieznormalizowana liczba, którą można reprezentować za pomocą danej arytmetyki. Po porównaniu jej z ϵ maszynowym, która jest następną liczbą po zerze, zauważymy że są to te same liczby.

1.2.6 MIN_{nor}

Funkcja `floatmin()` zwraca najmniejszą znormalizowaną liczbę dla danego typu. Jest to definicja liczby MIN_{nor} , więc te dwie liczby są sobie równe.

1.3 Wnioski

Pierwszy i najważniejszy wniosek to że arytmetyka i reprezentacje liczb w komputerach nie zawsze są precyzyjne, liczby nie są ciągłe i nie każdą liczbę rzeczywistą możemy dokładnie zapisać w pamięci komputera. Tu z pomocą przychodzi nam standard IEEE 754 który pozwala na dość dokładne reprezentowanie liczb ale cały czas ma swoje ograniczenia. Należy również pamiętać o tym że liczby mogą być znormalizowane i nieznormalizowane, gdyż czasami może to przypożyczyć problemów.

2 Wyrażenie Kahana

Zadanie 2 polegało na eksperymentalnym przetestowaniu stwierdzenia Kahana. Według niego epsilon maszynowy można obliczyć za pomocą wzoru

$$3 * (\frac{4}{3} - 1) - 1$$

Naszym zadaniem było sprawdzić to dla wszystkich typów zmiennopozycyjnych w języku Julia

2.1 Rozwiązanie

Funkcja napisana przez mnie zmienia wszystkie liczby na podany w parametrze T typ i oblicza wynik podanego przez Kahana wyrażenia.

```
1 # Funkcja znajdujaca epsilon maszynowy metoda Kahana dla typu T
2 function macheps(T)
3     return T(3) * (T(4) / T(3) - T(1)) - T(1)
4 end
```

Listing 4: Stwierdzenie Kahana

2.2 Wyniki oraz ich interpretacja

	Float16	Float32	Float64
Kahan	-0.000977	1.1920929e-7	-2.220446049250313e-16
eps()	0.000977	1.1920929e-7	2.220446049250313e-16

Wartości wyników zgadzają się, problem natomiast jest ze znakami w wynikach. Powodem wydaje się być przybliżenie wyniku dzielenia $\frac{4}{3}$ które dla Float32 jest przybliżeniem w górę a dla pozostałych typów jest przybliżeniem w dół.

2.3 Wnioski

Wniosków z tego zadania może być kilka. Jednym z nich może być że istnieją łatwiejsze sposoby na obliczanie epsilon maszynowego niż w pierwszym zadaniu. Kolejnym może być że nie każde obliczenia, mimo że dla nas proste dadzą dobry wynik przy skończonej dokładności.

3 Równomierne rozmieszczenie liczb

W tym zadaniu mieliśmy sprawdzić eksperymentalnie czy liczby w arytmetyce Float64 w języku Julia liczby są rozmieszczone równomiernie w zakresie $[1, 2)$. Każda liczba w podanym zakresie powinna móc być przedstawiona w postaci $x = 1 + k\delta$ gdzie $k = 1, 2, \dots, 2^{52} - 1$ oraz $\delta = 2^{-52}$.

3.1 Rozwiązanie

Z racji na ogromną ilość liczb w tym przedziale, nie jest możliwe sprawdzenie wszystkich liczb po kolei. Możemy natomiast sprawdzić początek i koniec przedziału aby sprawdzić czy tam zachowana została własność równomiernego rozmieszczenia z różnicami pomiędzy liczbami równymi δ . Innym sposobem sprawdzenia równomierności rozmieszczenia jest porównanie eksponenty, pierwszej i ostatniej liczby z przedziału ponieważ jeśli są takie same to jedyne co mogło się zmieniać to mantysa a wtedy różnice pomiędzy liczbami są na 100% takie same

```
1 # Funkcja sprawdzająca równomiernie rozmieszczone dla 100 pierwszych liczb
   wyznaczanych przez start jako pierwsza liczba oraz func jako funkcja
   wyznaczająca kolejną liczbę
2 # delta to oczekiwana różnica pomiędzy liczbami
3 # funkcja zwraca true jeśli zachowane jest równomiernie rozmieszczenie i false
   jeśli nie jest zachowane
4 function rownomiernie_rozmieszczone_na_ogonach(start, delta, func)
5     rownomiernie_rozmieszczone = true
6     for i in 0:100
7         if func(start) != start + delta
8             rownomiernie_rozmieszczone = false
9             break
10        end
11        start = func(start)
12    end
13
14    return rownomiernie_rozmieszczone
15 end
16
17 # Sprawdzanie początku przedziału [1, 2)
18 println("Równomierne rozmieszczenie na początku przedziału [1, 2): ",
19         rownomiernie_rozmieszczone_na_ogonach(one(Float64), eps(Float64), nextfloat))
20 # Sprawdzanie końca przedziału [1, 2)
21 println("Równomierne rozmieszczenie na końcu przedziału [1, 2): ",
22         rownomiernie_rozmieszczone_na_ogonach(one(Float64) + 1, -eps(Float64),
23         prevfloat))
24
25 # Porównywanie bit w liczb
26 x = one(Float64)
27
28 println(prevfloat(x), " - ", bitstring(prevfloat(x)))
29 println(x, " - ", bitstring(x))
30 println(nextfloat(x), " - ", bitstring(nextfloat(x)))
31 println(nextfloat(nextfloat(x)), " - ", bitstring(nextfloat(nextfloat(x))))
32 println(prevfloat(x+1), " - ", bitstring(prevfloat(x+1)))
33 println(x+1, " - ", bitstring(x+1))
```

Listing 5: Równomierne rozmieszczenie liczb

3.2 Wyniki oraz ich interpretacja

“Równomierne rozmieszczenie na początku przedziału $[1, 2)$: true
Równomierne rozmieszczenie na końcu przedziału $[1, 2)$: true“

[illegible]

Wyniki stworzonej przeze mnie potwierdzają że rozmieszczenie jest równomierne. Podobnie porównanie bitów pokazuje, że eksponenta się nie zmieniła a więc rozmieszczenie jest równomierne.

3.3 Wnioski

Liczby pomiędzy kolejnymi potęgami dwójki są równomiernie rozmieszczone. Im większe potęgi dwójki, tym większe różnice pomiędzy kolejnymi liczbami.

Zakres	δ
$[\frac{1}{2}, 1)$	1.1102230246251565e-16
$[1, 2)$	2.220446049250313e-16
$[2, 4)$	4.440892098500626e-16

4 Liczba nie spełniająca równania

W zadaniu 4 mieliśmy za zadanie znaleźć jakakolwiek liczbe zmiennopozycyjna typu Float64 taka, że:

$$x * \frac{1}{x} \neq 1$$

$$1 < x < 2$$

a także znaleźć najmniejszą taką liczbę.

4.1 Rozwiązanie

Funkcja, która napisałem, za parametry przyjmuje początek i koniec przeszukiwanego przedziału oraz funkcję, która wyznacza następną liczbę do sprawdzenia. Wewnątrz mojej funkcji jest pętla która przechodzi po kolejnych liczbach zmiennopozycyjnych i sprawdza czy spełniają one podane w zadaniu równanie. Jeśli tak to zwraca odpowiedź czy znaleziono oraz znaleziona liczba, jeśli nie znaleziono to zwraca odpowiedź i koniec przedziału.

```

1  # Funkcja znajduje liczbe nie spelniajaca rownania  $x * 1/x = 1$  w przedziale
    (start, end_)
2  # parametr func to funkcja ktora wyznacza kolejna sprawdzana liczbe
3  # zwraca: znaleziono - true jesli znaleziono liczbe nie spelniajaca rownania
    lub false jesli nie znaleziono
4  #         start - jesli znaleziono jest true to bedzie to liczba nie
    spelniajaca rownania
5  function find_ne(start, end_, func)
6      flag = true
7      if start > end_
8          flag = false
9      end
10     znaleziono = false
11     while (flag && start < end_) || (!flag && end_ < start)
12         if Float64(start * Float64(1 / start)) != 1
13             znaleziono = true
14             break

```

```

15         end
16         start = func(start)
17     end
18     return (start, znaleziono)
19 end
20
21 println("Szukanie liczby nie spełniającej działania  $x * 1/x = 1$  na początku
    przedziału (1, 2): ", find_ne(nextfloat(Float64(1)), prevfloat(Float64(2)),
    nextfloat))
22 println("Szukanie liczby nie spełniającej działania  $x * 1/x = 1$  na ko cu
    przedziału (1, 2): ", find_ne(prevfloat(Float64(2)), nextfloat(Float64(1)),
    prevfloat))
23 println("Najmniejsza liczba nie spełniająca równania  $x * 1/x = 1$  w przedziale
    (1, 2): ", find_ne(nextfloat(Float64(1)), prevfloat(Float64(2)), nextfloat))

```

Listing 6: Szukanie liczby x nie spełniającej równaia

4.2 Wyniki i interpretacja

Wystarczy znaleźć najmniejsza liczbę, żeby wykonać oba podpunkty zdania. Liczba taka jest:

$$x = 1.000000057228997$$

4.3 Wnioski

Arytmetyka liczb zmiennopozycyjnych ze względu na swoją niedokładność potrafi źle obliczyć nawet najprostsze działania.

5 Iloczyn skalarny wektorów

W tym zadaniu mieliśmy przetestować różne sposoby obliczania iloczynu skalarnego dla pojedynczej i podwójnej precyzji (typów Float32 oraz Float64). Wybranymi sposobami liczenia iloczynu skalarnego były:

1. "W przód" - $\sum_{i=1}^n x_i y_i$
2. "W tył" - $\sum_{i=n}^1 x_i y_i$
3. Najpierw zsumować liczby dodatnie od największej do najmniejszej a potem liczby ujemne od najmniejszej do największej
4. Najpierw zsumować liczby dodatnie od najmniejszej do największej a potem liczby ujemne od największej do najmniejszej

5.1 Rozwiązanie

Moje funkcje przyjmują 3 argumenty:

- x - pierwszy wektor
- y - drugi wektor
- T - typ danych które ma liczyć funkcja

oblicza iloczyn skalarny a następnie sumuje i zwraca wynik.

5.1.1 Sposób "W przód"

```
1 # Funkcja zwraca iloczyn skalarany x oraz y liczony "w przod"
2 # T to typ zmiennej
3 function a(x, y, T)
4     result = T(0)
5     n = length(x)
6
7     for i in 1:n
8         result += x[i] * y[i]
9     end
10
11     return result
12 end
```

Listing 7: Sposób nr. 1

5.1.2 Sposób "W tył"

```
1 # Funkcja zwraca iloczyn skalarany x oraz y liczony "w ty "
2 # T to typ zmiennej
3 function b(x, y, T)
4     result = T(0)
5     n = length(x)
6
7     for i in n:-1:1
8         result += (x[i] * y[i])
9     end
10
11     return result
12 end
```

Listing 8: Sposób nr. 2

5.1.3 Sposób "Największe -i najmniejszych"

```
1 # Funkcja zwraca iloczyn skalarany x oraz y liczony:
2 #     Dodatnie od największego do najmniejszego
3 #     Ujemne od najmniejszego do największego
4 # T to typ zmiennej
5 function c(x, y, T)
6     to_sum = x .* y
7
8     sum_positive = filter(a -> a >= 0, to_sum)
9     sum_negative = filter(a -> a < 0, to_sum)
10
11     sum_positive = sort(sum_positive, rev=true)
12     sum_negative = sort(sum_negative)
13
14     sum = T(0)
15     for i in eachindex(sum_positive)
16         sum += sum_positive[i]
17     end
18     for i in eachindex(sum_negative)
19         sum += sum_negative[i]
20     end
21
22     return sum
23 end
```

Listing 9: Sposób nr. 3

5.1.4 Sposób "Najmniejsze do największych"

```
1 # Funkcja zwraca iloczyn skalarany x oraz y liczony:
2 #     Dodatnie od najmniejszego do największego
3 #     Ujemne od największego do najmniejszego
4 # T to typ zmiennej
5 function d(x, y, T)
6     to_sum = x .* y
7
8     sum_positive = filter(a -> a >= 0, to_sum)
9     sum_negative = filter(a -> a < 0, to_sum)
10
11     sum_positive = sort(sum_positive)
12     sum_negative = sort(sum_negative, rev=true)
13
14     sum = T(0)
15     for i in eachindex(sum_positive)
16         sum += sum_positive[i]
17     end
18     for i in eachindex(sum_negative)
19         sum += sum_negative[i]
20     end
21
22     return sum
23 end
```

Listing 10: Sposób nr. 4

5.2 Wyniki oraz ich interpretacja

	Float32	Float64	Prawdziwy wynik
Sposób nr. 1	-0.4999443	1.0251881368296672e-10	-1.006571070000000e-11
Sposób nr. 2	-0.4543457	-1.5643308870494366e-10	-1.006571070000000e-11
Sposób nr. 3	-0.39291382	1.4068746168049984e-12	-1.006571070000000e-11
Sposób nr. 4	-0.5	0.0	-1.006571070000000e-11

Precyzja pojedyncza wykazuje większe błędy niż precyzja podwójna. W precyzji podwójnej pojawiły się wyniki które mają w sobie nawet błędny znak. Sposób nr. 4 wydaje się być najgorszym ze wszystkich algorytmów.

5.3 Wnioski

Wyniki działań, mimo że nie powinny, mogą zależeć od kolejności ich wykonywania. Różnice jakie wynikają z tego typu problemu mogą być zadziwiająco duże.

6 Porównanie $f(x)=g(x)$

W zadaniu nr. 6 należało porównywać, w podójnej precyzji, dwie następujące funkcje:

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

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

Funkcje te wydają się dla nas ludzi być takie same, ponieważ takie są. Ale jak w poprzednim zadaniu było widać czasami kolejność działań bardzo wpływa na wyniki działań.

6.1 Rozwiązanie

Funkcje które napisałem są bardzo prosto napisane, na koniec skryptu zostają one porównane dla pierwszych 10 liczb a następnie dla wielokrotności 20 aż do 100.

```
1 # Funkcja podana w zadaniu przyjmuj ca x jako liczb
2 f(x) = sqrt(x^2 + 1.0) - 1.0
3
4 # Funkcja podana w zadaniu przyjmuj ca x jako liczb
5 g(x) = x^2 / (sqrt(x^2 + 1.0) + 1.0)
6
7 # Ta cz pokazuje r nice pomiedzy warto ciami funkcji wy ej
8   zdefiniowanymi
9 for i in 1:10
10     println("f(8^-$i) = ", f(Float64(8)^-i))
11     println("g(8^-$i) = ", g(Float64(8)^-i), "\n")
12 end
13 for i in [20, 40, 60, 80, 100]
14     println("f(8^-$i) = ", f(Float64(8)^-i))
15     println("g(8^-$i) = ", g(Float64(8)^-i), "\n")
16 end
```

Listing 11: Porównanie $f(x)$ oraz $g(x)$

6.2 Wyniki oraz ich interpretacja

x	f(x)	g(x)
8^{-1}	0.0077822185373186414	0.0077822185373187065
8^{-2}	0.00012206286282867573	0.00012206286282875901
8^{-3}	1.9073468138230965e-6	1.907346813826566e-6
8^{-4}	2.9802321943606103e-8	2.9802321943606116e-8
8^{-5}	4.656612873077393e-10	4.6566128719931904e-10
8^{-6}	7.275957614183426e-12	7.275957614156956e-12
8^{-7}	1.1368683772161603e-13	1.1368683772160957e-13
8^{-8}	1.7763568394002505e-15	1.7763568394002489e-15
8^{-9}	0.0	2.7755575615628914e-17
8^{-10}	0.0	4.336808689942018e-19
8^{-20}	0.0	3.76158192263132e-37
8^{-40}	0.0	2.8298997121333476e-73
8^{-60}	0.0	2.1289799200040754e-109
8^{-80}	0.0	1.6016664761464807e-145
8^{-100}	0.0	1.204959932551442e-181

Dla większych wartości x widać podobieństwa w wynikach, nie są one takie same ale są na tyle podobne że ciężko byłoby powiedzieć które są bliżej prawdy. Patrząc na mniejsze x widać już, że funkcja $f(x)$ nie daje sobie rady z tak małymi wartościami i w porównaniu do $g(x)$ szybko traci wiarygodność zwracając wartości 0.0. Dzieje się tak ze względu na odejmowanie 1 od pierwiastka. Funkcja $g(x)$ omija ten problem, dzięki temu zwraca wiarygodne wyniki jeszcze dużo dłużej niż $f(x)$.

6.3 Wnioski

Kolejność działań potrafi sprawić że wyniki stają się szybko dalekie od prawdy. Ale z drugiej strony możemy skonstruować działania z taką kolejnością działań, które będą z dużą dokładnością podawały nam wyniki.

7 Przybliżanie wartości pochodnej

W tym zadaniu naszym celem jest porównanie przybliżonych wartości pochodnej:

$$f'(x_0) \approx f'(x_0) = \frac{f(x_0 + h) - f(x_0)}{h}$$

z prawdziwymi wartościami pochodnej dla $h = 2^{-n}$ ($n = 0, 1, 2, \dots, 54$) Funkcja której pochodną mamy policzyć jest:

$$f(x) = \sin(x) + \cos(3x)$$

7.1 Rozwiązanie

Rozwiązanie zawiera zaimplementowane funkcje podane w treści zadania. Funkcja obliczająca prawdziwą wartość z pochodnej w punkcie x została policzona przeze mnie.

$$\frac{d}{dx}(\sin(x) + \cos(3x)) = \cos(x) - 3\sin(3x)$$

```

1 # Funkcja oblicza wartosc funkcji podanej w zadaniu w punkcie x
2 f(x) = sin(x) + cos(3x)
3 # Funkcja oblicza wartosc pochodnej funkcji podanej w zadaniu w punkcie x
4 dfdx(x) = cos(x) - 3sin(3x)
5 # Aproksymacja pochodnej funkcji podanej w parametrze f w punkcie x i z
  precyzj h
6 dfdx_approx(x, h, f) = (f(x + h) - f(x)) / h
7
8 # Prawdziwa wartosc pochodnej funkcji f w punkcie x = 1
9 actual = dfdx(1)
10
11
12 for n in 0:54
13     println("\nh = 2^-$n:")
14     println("~f'(1)          = ", dfdx_approx(1.0, Float64(2)^-n, f))
15     println(" f'(1)          = ", actual)
16     println("|f'(1) - ~f'(1)| = ", abs(dfdx_approx(1.0, Float64(2)^-n, f) -
      actual))
17 end

```

Listing 12: Obliczanie przybliżonej pochodnej

7.2 Wyniki oraz ich interpretacja

Prawdziwy wynik pochodnej w punkcie $x_0 = 1$ wynosi:

$$f'(x_0) = 0.11694228168853815$$

h	$h + 1$	$f'(x_0)$	$ f'(x_0) - f'(x_0) $
2^0	2.0	2.0179892252685967	1.9010469435800585
2^{-1}	1.5	1.8704413979316472	1.753499116243109
2^{-2}	1.25	1.1077870952342974	0.9908448135457593
2^{-3}	1.125	0.6232412792975817	0.5062989976090435
2^{-4}	1.0625	0.3704000662035192	0.253457784514981
2^{-5}	1.03125	0.24344307439754687	0.1265007927090087
2^{-6}	1.015625	0.18009756330732785	0.0631552816187897
2^{-7}	1.0078125	0.1484913953710958	0.03154911368255764
2^{-8}	1.00390625	0.1327091142805159	0.015766832591977753
2^{-9}	1.001953125	0.1248236929407085	0.007881411252170345
2^{-10}	1.0009765625	0.12088247681106168	0.0039401951225235265
2^{-11}	1.00048828125	0.11891225046883847	0.001969968780300313
2^{-12}	1.000244140625	0.11792723373901026	0.0009849520504721099
2^{-13}	1.0001220703125	0.11743474961076572	0.0004924679222275685
2^{-14}	1.00006103515625	0.11718851362093119	0.0002462319323930373
2^{-15}	1.000030517578125	0.11706539714577957	0.00012311545724141837
2^{-16}	1.0000152587890625	0.11700383928837255	6.155759983439424e-5
2^{-17}	1.0000076293945312	0.11697306045971345	3.077877117529937e-5
2^{-18}	1.0000038146972656	0.11695767106721178	1.5389378673624776e-5
2^{-19}	1.0000019073486328	0.11694997636368498	7.694675146829866e-6
2^{-20}	1.0000009536743164	0.11694612901192158	3.8473233834324105e-6

2^{-20}	1.0000009536743164	0.11694612901192158	3.8473233834324105e-6
2^{-21}	1.0000004768371582	0.1169442052487284	1.9235601902423127e-6
2^{-22}	1.000000238418579	0.11694324295967817	9.612711400208696e-7
2^{-23}	1.0000001192092896	0.11694276239722967	4.807086915192826e-7
2^{-24}	1.0000000596046448	0.11694252118468285	2.394961446938737e-7
2^{-25}	1.0000000298023224	0.116942398250103	1.1656156484463054e-7
2^{-26}	1.0000000149011612	0.11694233864545822	5.6956920069239914e-8
2^{-27}	1.0000000074505806	0.11694231629371643	3.460517827846843e-8
2^{-28}	1.0000000037252903	0.11694228649139404	4.802855890773117e-9
2^{-29}	1.0000000018626451	0.11694222688674927	5.480178888461751e-8
2^{-30}	1.0000000009313226	0.11694216728210449	1.1440643366000813e-7
2^{-31}	1.0000000004656613	0.11694216728210449	1.1440643366000813e-7
2^{-32}	1.0000000002328306	0.11694192886352539	3.5282501276157063e-7
2^{-33}	1.0000000001164153	0.11694145202636719	8.296621709646956e-7
2^{-34}	1.0000000000582077	0.11694145202636719	8.296621709646956e-7
2^{-35}	1.0000000000291038	0.11693954467773438	2.7370108037771956e-6
2^{-36}	1.000000000014552	0.116943359375	1.0776864618478044e-6
2^{-37}	1.000000000007276	0.1169281005859375	1.4181102600652196e-5
2^{-38}	1.000000000003638	0.116943359375	1.0776864618478044e-6
2^{-39}	1.000000000001819	0.11688232421875	5.9957469788152196e-5
2^{-40}	1.0000000000009095	0.1168212890625	0.0001209926260381522

2^{-41}	1.0000000000004547	0.116943359375	1.0776864618478044e-6
2^{-42}	1.0000000000002274	0.1166921875	0.0002430629385381522
2^{-43}	1.0000000000001137	0.1162109375	0.0007313441885381522
2^{-44}	1.0000000000000568	0.1171875	0.0002452183114618478
2^{-45}	1.0000000000000284	0.11328125	0.003661031688538152
2^{-46}	1.0000000000000142	0.109375	0.007567281688538152
2^{-47}	1.000000000000007	0.109375	0.007567281688538152
2^{-48}	1.0000000000000036	0.09375	0.023192281688538152
2^{-49}	1.0000000000000018	0.125	0.008057718311461848
2^{-50}	1.0000000000000009	0.0	0.11694228168853815
2^{-51}	1.0000000000000004	0.0	0.11694228168853815
2^{-52}	1.0000000000000002	-0.5	0.6169422816885382
2^{-53}	1.0	0.0	0.11694228168853815
2^{-54}	1.0	0.0	0.11694228168853815

Po wynikach widać że na początek dokładność przybliżenia rośnie a następnie zaczyna maleć. Najlepsze przybliżenie daje $h = 2^{-28}$. Wynika to z błędów związanych z wartościami h coraz bliższymi do 0.

7.3 Wnioski

Przy prowadzeniu obliczeń przybliżających wynik najlepiej nie korzystać z liczb tak bliskich do zera. Istnieją liczby które wbrew intuicji będą lepiej przybliżać niż te bliskie zeru.