

# Obliczenia naukowe - lista pierwsza

Bartosz Rajczyk

16 października 2019

## Spis treści

|          |                                  |           |
|----------|----------------------------------|-----------|
| <b>1</b> | <b>Zadanie 1.</b>                | <b>3</b>  |
| 1.1      | Opis problemu . . . . .          | 3         |
| 1.2      | Rozwiązanie . . . . .            | 3         |
| 1.3      | Wyniki i interpretacja . . . . . | 3         |
| 1.4      | Wnioski . . . . .                | 4         |
| <b>2</b> | <b>Zadanie 2.</b>                | <b>4</b>  |
| 2.1      | Opis problemu . . . . .          | 4         |
| 2.2      | Rozwiązanie . . . . .            | 4         |
| 2.3      | Wyniki i interpretacja . . . . . | 5         |
| 2.4      | Wnioski . . . . .                | 5         |
| <b>3</b> | <b>Zadanie 3.</b>                | <b>5</b>  |
| 3.1      | Opis problemu . . . . .          | 5         |
| 3.2      | Rozwiązanie . . . . .            | 5         |
| 3.3      | Wyniki i interpretacja . . . . . | 6         |
| 3.4      | Wnioski . . . . .                | 6         |
| <b>4</b> | <b>Zadanie 4.</b>                | <b>6</b>  |
| 4.1      | Opis problemu . . . . .          | 6         |
| 4.2      | Rozwiązanie . . . . .            | 6         |
| 4.3      | Wyniki i interpretacja . . . . . | 7         |
| 4.4      | Wnioski . . . . .                | 7         |
| <b>5</b> | <b>Zadanie 5.</b>                | <b>7</b>  |
| 5.1      | Opis problemu . . . . .          | 7         |
| 5.2      | Rozwiązanie . . . . .            | 7         |
| 5.3      | Wyniki i interpretacja . . . . . | 7         |
| 5.4      | Wnioski . . . . .                | 8         |
| <b>6</b> | <b>Zadanie 6.</b>                | <b>8</b>  |
| 6.1      | Opis problemu . . . . .          | 8         |
| 6.2      | Rozwiązanie . . . . .            | 8         |
| 6.3      | Wyniki i interpretacja . . . . . | 8         |
| 6.4      | Wnioski . . . . .                | 9         |
| <b>7</b> | <b>Zadanie 7.</b>                | <b>9</b>  |
| 7.1      | Opis problemu . . . . .          | 9         |
| 7.2      | Rozwiązanie . . . . .            | 9         |
| 7.3      | Wyniki i interpretacja . . . . . | 10        |
| 7.4      | Wnioski . . . . .                | 11        |
| <b>8</b> | <b>Bibliografia</b>              | <b>12</b> |

# 1 Zadanie 1.

## 1.1 Opis problemu

Problem polega na wyznaczeniu liczb:

1. *macheps*, czyli najmniejszej liczby spełniającej  $fl(1.0 + macheps) > 1.0$
2. *eta*, czyli najmniejszej reprezentowalnej liczby
3. *max*, czyli największej reprezentowalnej liczby

dla typów zmiennoprzecinkowych o precyzji 16-, 32- oraz 64-bitowej.

## 1.2 Rozwiązanie

Zaprezentowane rozwiązania w pliku 1.jl polegają na iteracyjnym dzieleniu lub mnożeniu liczby początkowej (uzyskiwanej najczęściej przez wywołanie `one(type)` zwracające jedynkę w podanym typie danych w Julii) aż do spełnienia określonego warunku, będącego w kolejności:

1.  $1 + aktualna\_liczba/2 \neq 1$
2.  $aktualna\_liczba/2 \neq 0$
3.  $aktualna\_liczba * 2 \neq \infty$

Przy spełnieniu tego warunku pętla jest zakańczana i zwracana jest aktualna liczba.

## 1.3 Wyniki i interpretacja

| typ danych | moja funkcja          | eps                   | float.h      |
|------------|-----------------------|-----------------------|--------------|
| Float16    | 0.000977              | 0.000977              | niedostępne  |
| Float32    | 1.1920929e-7          | 1.1920929e-7          | 1.192093e-07 |
| Float64    | 2.220446049250313e-16 | 2.220446049250313e-16 | 2.220446e-16 |

Tabela 1: wartości epsilon maszynowego

| typ danych | moja funkcja | nextfloat(0) |
|------------|--------------|--------------|
| Float16    | 6.0e-8       | 6.0e-8       |
| Float32    | 1.0e-45      | 1.0e-45      |
| Float64    | 5.0e-324     | 5.0e-324     |

Tabela 2: wartości eta

| typ danych | moja funkcja         | floatmax               |
|------------|----------------------|------------------------|
| Float16    | 3.277e4              | 6.55e4                 |
| Float32    | 1.7014118e38         | 3.4028235e38           |
| Float64    | 8.98846567431158e307 | 1.7976931348623157e308 |

Tabela 3: wartości max

Z powyższych danych wynika, że zaimplementowane przeze mnie funkcje poprawnie zwracały wyniki poza przypadkiem wyliczania liczby maksymalnej, gdzie uzyskane rezultaty były dwukrotnie za małe. Wynika to najprawdopodobniej z tego, że mnożąc liczbę za każdym krokiem dwukrotnie wkraczamy w pewnym momencie na dokładną wartość nieskończoności, kiedy poprawnym rozwiązaniem jest jedna liczba przed nią, stąd ta rozbieżność.

## 1.4 Wnioski

Standardowa reprezentacja liczb zgodna z IEEE754 w komputerze ma skończoną dokładność i posiada skończone zakresy liczb, które potrafi reprezentować. Jednocześnie istnieją dwie osobne i rozróżnialne reprezentacje bardzo małych liczb - znormalizowane i zdenormalizowane. Wyznaczone przez nas liczby *eta* oraz liczby zwracane przez *nextfloat(0)* dla różnych typów danych mają postać zdenormalizowaną (aczkolwiek nie pokrywają się idealnie z wartościami podawanymi przez źródła [1], najpewniej przez zaokrąglenia). Inne funkcje, jak chociażby *floatmin* zwracają wartości znormalizowane. Należy więc pamiętać o istnieniu obu reprezentacji i odpowiednim używaniu ich - także ze względu na używany procesor, bo ze źródeł wynika, że nie każde CPU obsługuje liczby zdenormalizowane całkowicie poprawnie.

## 2 Zadanie 2.

### 2.1 Opis problemu

Problem polega na odnalezieniu epsilonu maszynowego używając wyrażenia:

$$3 * (4/3 - 1) - 1$$

w arytmetyce zmiennoprzecinkowej.

### 2.2 Rozwiązanie

Zaprezentowane rozwiązanie w pliku 2. j1 polega na użyciu funkcji one do uzyskania wartości 1 w danym typie liczbowym, a następnie wykonanie podanego działania korzystając z liczb opartych na tej jedynce.

## 2.3 Wyniki i interpretacja

| typ danych | moja funkcja           | eps                   |
|------------|------------------------|-----------------------|
| Float16    | -0.000977              | 0.000977              |
| Float32    | 1.1920929e-7           | 1.1920929e-7          |
| Float64    | -2.220446049250313e-16 | 2.220446049250313e-16 |

Tabela 4: obliczone wartości  $3*(4/3 - 1) - 1$

Widzimy, że wartość bezwzględna z uzyskanych przez nas wyników pokrywa się z epsilonem zwracanym przez funkcję. Zmiana znaku spowodowana jest najpewniej tym, że liczba bitów znaczących dla kolejnych danych wynosi:

- Float16 - 10
- Float32 - 23
- Float64 - 52

A rozwinięciem binarnym ułamka  $4/3$  jest  $1.(10)$ , więc dla typów Float16 i Float32 ostatnią cyfrą mantysy będzie 0, a dla typu Float64 - 1, co powinno decydować o znaku odejmowania.

## 2.4 Wnioski

Przez skończoną dokładność reprezentacji, niektóre równania dające w normalnej arytmetyce zero, w arytmetyce zmiennoprzecinkowej mogą dawać inne wyniki.

# 3 Zadanie 3.

## 3.1 Opis problemu

Problem polega na sprawdzeniu rozmieszczenia liczb zmiennoprzecinkowych w arytmetyce IEEE 754 podwójnej precyzji w danych przedziałach liczbowych.

## 3.2 Rozwiązanie

Zaprezentowane rozwiązania w pliku 3. j1 to dwa różne podejścia do tego problemu. Pierwsze polega na ręcznym przeiterowaniu przez wszystkie istniejące liczby podanego przedziału i sprawdzenie, czy odległości między kolejnymi rzeczywście zawsze wynoszą sprawdzaną wartość. Drugie to bardziej analityczne rozwiązanie bazujące na naszej znajomości standardu IEEE754 i analizie rozmieszczenia bitów zwracanego przez funkcję `bitstring`. W tym podejściu porównujemy pierwszą i ostatnią liczbę z przedziału; jeżeli ich eksponenty są inne,

wyklucza to równomierny rozkład pomiędzy nimi. Jeżeli są równe, możemy obliczyć jak bardzo zmienia się liczba przy powiększeniu o jeden mantysy używając wzoru:

$$2^{\text{eksponenta}-1023} * 2^{-52}$$

ponieważ biasem dla eksponenty Float64 jest 1023, a mantysa ma 52 bity znaczące.

### 3.3 Wyniki i interpretacja

Utworzona przez mnie funkcja potwierdziła, że w przedziale  $[1, 2]$  liczby rozmieszczone są co  $2^{-52}$ . Dodatkowo udało się wyznaczyć rozmieszczenie liczb dla innych przedziałów:

| przedział  | odległości             |
|------------|------------------------|
| $[0.5, 1]$ | 1.1102230246251565e-16 |
| $[1, 2]$   | 2.220446049250313e-16  |
| $[2, 4]$   | 4.440892098500626e-16  |

Tabela 5: przedziały i odległości między liczbami

Te wyniki pokrywają się z moim rozumieniem standardu IEEE754 - ze względu na powiększenie się eksponenty o jeden, odległości pomiędzy kolejnymi liczbami rosną dwukrotnie (ponieważ eksponenta jest wykorzystywana jako wykładnik w  $2^{\text{eksponenta}}$ ).

### 3.4 Wnioski

Liczby w IEEE754 są reprezentowane z określoną dokładnością różniącą się zależnie od przedziału, w którym się znajdują.

## 4 Zadanie 4.

### 4.1 Opis problemu

Problem polega na znalezieniu dwóch liczb:

1. najmniejszego  $1 < x < 2$ , takiego że  $x * \frac{1}{x} \neq 1$
2. najmniejszego  $0 < x$ , takiego że  $x * \frac{1}{x} \neq 1$

### 4.2 Rozwiązanie

Zaprezentowane rozwiązania w pliku 4.jl polegają na rozpoczęciu od jednej liczby powyżej dolnego ograniczenia dla  $x$  i powiększaniu jej do następnej aż do uzyskania szukanej nierówności.

### 4.3 Wyniki i interpretacja

Algorytm zwrócił następujące wyniki:

1.  $x = 1.000000057228997 \rightarrow x * \frac{1}{x} = 0.9999999999999999$

2.  $x = 1.0e - 323 \rightarrow x * \frac{1}{x} = \infty$

### 4.4 Wnioski

Arytmetyka zmiennoprzecinkowa ze względu na swoją skończoną dokładność nie zawsze zwraca poprawne wyniki nawet najprostszych działań.

## 5 Zadanie 5.

### 5.1 Opis problemu

Problem polega na obliczaniu iloczynu skalarnego dwóch podanych wektorów na kilka różnych sposobów:

1. w przód, zaczynając dodawanie od pierwszych indeksów
2. w tył, zaczynając dodawanie od ostatnich indeksów
3. od największego iloczynu
4. od najmniejszego iloczynu

### 5.2 Rozwiązanie

Zaprezentowane rozwiązania w pliku 5. j1 polegają na dosłownej implementacji podanych algorytmów.

### 5.3 Wyniki i interpretacja

| algorytm | wynik                   |
|----------|-------------------------|
| 1        | 1.0251881368296672e-10  |
| 2        | -1.5643308870494366e-10 |
| 3        | 0.0                     |
| 4        | 0.0                     |

Tabela 6: Wyniki dla Float64

| algorytm | wynik               |
|----------|---------------------|
| 1        | -0.3472038161853561 |
| 2        | -0.3472038162872195 |
| 3        | -0.5                |
| 4        | -0.5                |

Tabela 7: Wyniki dla Float32

## 5.4 Wnioski

Kolejność wykonywanych obliczeń może drastycznie wpłynąć na otrzymywany wynik.

## 6 Zadanie 6.

### 6.1 Opis problemu

Problem polega na obliczeniu wartości dwóch równoważnych matematycznie funkcji w dla kolejnych wartości  $8^{-1}, 8^{-2}, 8^{-3}, \dots$ . Funkcje to:

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

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

### 6.2 Rozwiązanie

Zaprezentowane rozwiązania w pliku 6. j1 polegają na dosłownej implementacji podanych równań.

### 6.3 Wyniki i interpretacja

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



|     |     |                         |
|-----|-----|-------------------------|
| 80  | 0.0 | 1.6016664761464807e-145 |
| 100 | 0.0 | 1.204959932551442e-181  |
| 120 | 0.0 | 9.065110999561118e-218  |
| 140 | 0.0 | 6.819831532519088e-254  |
| 160 | 0.0 | 5.1306710016229703e-290 |
| 180 | 0.0 | 0.0                     |

Tabela 8: porównanie wyników  $f$  i  $g$

Widzimy, że funkcja  $f$  bardzo szybko daje wyniki równe zero, natomiast funkcja  $g$  pozwala obliczyć swoją wartość bardzo blisko minimalnych liczb w zakresie typu danych Float64. Widzimy dodatkowo, że dla wartości  $x$ , dla których obie funkcje nadal dają rezultaty, są one dosyć podobne, chociaż nie identyczne. Funkcja  $f$  radzi sobie gorzej ze względu na to, że operuje ona na wartościach bardzo bliskich zero ze względu na odejmowanie jedynki od pierwiastka - w ten sposób traci dużo cyfr znaczących z wyniku pierwiastka. Funkcja  $g$  obchodzi ten problem i dlatego pozwala na otrzymanie wyników przy znacznie mniejszym wejściu.

## 6.4 Wnioski

Należy wykonywać obliczenia w ten sposób, aby liczba cyfr znaczących przy kolejnych działaniach zbyt nie różniła, ponieważ pozwala to na znaczne poprawienie dokładności obliczeń.

# 7 Zadanie 7.

## 7.1 Opis problemu

Problem polega na sprawdzeniu dokładności pochodnej funkcji liczonej przy pomocy wzoru

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

dla  $h \rightarrow 0$  i porównanie jej z wartością matematycznej pochodnej funkcji w punkcie  $x_0 = 1$ . Badaną funkcją jest

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

$$f'(x) = \cos(x) - 3\sin(3x)$$

## 7.2 Rozwiązanie

Zaprezentowane rozwiązanie w pliku 7. j1 polega na dosłownej implementacji podanych równań.

### 7.3 Wyniki i interpretacja

Obliczoną wartością pochodnej w punkcie  $x + 0 = 1$  było 0.11694228168853815. Tabela w kolumnie "pochodna" zawiera obliczoną pochodną ze wzoru dla danego  $h$ , kolumna "różnica" jest odległością obliczonej w ten sposób pochodnej od jej prawdziwej wartości.

| <b>h</b>  | <b>1+h</b>         | <b>pochodna</b>     | <b>różnica</b>         |
|-----------|--------------------|---------------------|------------------------|
| $2^{-0}$  | 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  |
| $2^{-21}$ | 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  |

|           |                    |                    |                       |
|-----------|--------------------|--------------------|-----------------------|
| $2^{-37}$ | 1.0000000000007276 | 0.1169281005859375 | 1.4181102600652196e-5 |
| $2^{-38}$ | 1.0000000000003638 | 0.116943359375     | 1.0776864618478044e-6 |
| $2^{-39}$ | 1.0000000000001819 | 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 | 0.11669921875      | 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   |

Tabela 9: porównanie wyników f i g

Możemy zauważyć, że wbrew intuicji matematycznej najlepsze przybliżenie wartości pochodnej w tym punkcie jest uzyskiwane dla  $h = 2^{-28}$ , gdzie błąd spada do rzędu wielkości  $10^{-9}$ , natomiast po tej wartości już tylko rośnie, aby na sam koniec wynieść 100%. Jest to spowodowane faktem, że bardzo małe liczby zmiennoprzecinkowe posiadają niewielką liczbę cyfr znaczących w swoim zapisie, więc wraz z ich maleniem tracą dokładność obliczeń, aż do momentu, w którym stają się one bezużyteczne.

## 7.4 Wnioski

W prowadzonych obliczeniach najlepiej unikać wartości bardzo bliskich zeru.

## 8 Bibliografia

### Literatura

- [1] W. Kahan *Lecture Notes on the Status of IEEE Standard 754 for Binary Floating-Point Arithmetic*. 1997