

# Laboratorium 2

## Arytmetyka komputerowa (cd.)

Dominik Marek

12 marca 2024



**AKADEMIA GÓRNICZO-HUTNICZA  
IM. STANISŁAWA STASZICA  
W KRAKOWIE**

### ***1. Treść zadań.***

1. Napisać algorytm do obliczenia funkcji wykładniczej  $e^x$  przy pomocy nieskończonych szeregów

$$e^x = 1 + x + \frac{x^2}{2!} + \frac{x^3}{3!} + \dots$$

- a) Wykonując sumowanie w naturalnej kolejności, jakie kryterium zakończenia obliczeń przyjmiesz ?
- b) Proszę przetestować algorytm dla:

$x = +1, +5, +10$

i porównać wyniki z wynikami wykonania standardowej funkcji  $\exp(x)$

- c) Czy można posłużyć się szeregami w tej postaci do uzyskania dokładnych wyników dla  $x < 0$  ?
- d) Czy możesz zmienić wygląd szeregu lub w jakiś sposób przegrupować składowe żeby uzyskać dokładniejsze wyniki dla  $x < 0$  ?

2. Które z dwóch matematycznie ekwiwalentnych wyrażeń  $x^2 - y^2$  oraz  $(x - y) \cdot (x + y)$  może być obliczone dokładniej w arytmetyce zmienna-przecinkowej ? Dlaczego ? Dla jakich wartości  $x$  i  $y$ , względem siebie, istnieje wyraźna różnica w dokładności dwóch wyrażeń ?

3.

Zakładamy że rozwiązujemy równanie kwadratowe  $ax^2 + bx + c = 0$  z  $a = 1.22$ ,  $b = 3.34$  i  $c = 2.28$ , wykorzystując znormalizowany system zmienna-przecinkowy z podstawą  $\beta = 10$  i dokładnością  $p=3$

(a) ile wyniesie obliczona wartość  $b^2 - 4ac$  ?

(b) jaka jest dokładna wartość wyróżnika w rzeczywistej (dokładnej) arytmetyce ?

(c) jaki jest względny błąd w obliczonej wartości wyróżnika ?

## 2. Rozwiązania

1.

W celu obliczenia wartości  $e^x$ , wykorzystującą rozwinięcie w nieskończony szereg, do póki wartość bezwzględna kolejnego wyrazu szeregu jest większą niż przyjęte epsilon, zwiększam wartość zmiennej `exponents`, reprezentującej wartość obliczanej funkcji, o dany wyraz. Następnie obliczam kolejny wyraz szeregu poprzez przemnożenie poprzedniego przez iloraz podanej wartości  $x$  i kolejnej liczby naturalnej począwszy od 1.

*Funkcja obliczająca wartość funkcji  $e^x$ , według powyższego algorytmu:*

```
def find_exp(x: int, epsilon: float) -> float:
    exponents = 0
    element = 1
    i = 1
    while abs(element) > epsilon:
        exponents += element
        element *= x / i
        i += 1

    return exponents
```

a) Dla sumowania w naturalnej kolejności, jak kryterium końca obliczeń przyjmuje maszynowe epsilon wyliczone za pomocą poniższej funkcji:

```
def find_machine_epsilon(basis: int, precision: int) -> float:
    return basis**(1-precision)
```

W tym zadaniu posługując się powyższą funkcją dla systemu dwójkowego z precyzją równą 53, otrzymuję  $\varepsilon = 2.220446049250313e-16$ , co jest zgodne ze standardem IEEE 754 dla języka Python.

b)

*Zestawienie wyników otrzymanych za pomocą powyższego algorytmu:*

Argument	Wynik algorytmu	Dokładna wartość	Błąd bezwzględny
-1	0.3678794411714423	0.36787944117144233	1.508949536687701e-16
1	2.7182818284590455	2.718281828459045	1.6337129034990842e-16
-5	0.006737946999084603	0.006737946999085469	1.2847043982832115e-13
5	148.4131591025766	148.41315910257657	1.9150397176546985e-16
-10	4.5399929670295957e-05	4.5399929762484875e-05	2.030596046863973e-09
10	22026.465794806714	22026.465794806703	4.954919469581169e-16

### Wnioski:

Analizując powyższe wyniki, można zauważyć, iż dla dodatnich wartości wykładnika otrzymane wyniki są bliskie rzeczywistym wartościom. Natomiast w przypadku ujemnego wykładnika należy zwrócić uwagę na dostrzegalną różnicę względem oczekiwanej wartości. Ponadto różnica ta rośnie wraz ze zmniejszaniem wykładnika.

c) W przypadku dla  $x < 0$ , w szeregu będą pojawiały się na przemian wyrazy dodatnie i ujemne, zatem będziemy na przemian wykonywać operację dodawania i odejmowania.

Postać rozwinięcia  $e^x$  w nieskończony szereg dla  $x < 0$ :

$$e^{-x} = 1 - x + \frac{x^2}{2!} - \frac{x^3}{3!} + \frac{x^4}{4!} + \dots$$

Wówczas będziemy mieli do czynienia ze zjawiskiem zwanym „catastrophic cancellation”, polegającym na tym, iż przy odejmowaniu dwóch bliskich liczb wynik tej operacji po normalizacji może znacząco różnić się względem rzeczywistej różnicy tych liczb. Dlatego chcąc otrzymać odpowiednią dokładność wyników nie należy korzystać z podstawowej postaci tego szeregu.

d) W celu poprawy dokładności wyników, można skorzystać z następującej własności:

$e^{-x} = \frac{1}{e^x}$ , wówczas nasz szereg przyjmie następującą postać:

$$e^{-x} = \frac{1}{1 + x + \frac{x^2}{2!} + \frac{x^3}{3!} + \dots}$$

Korzystając z wyżej otrzymanej postaci, wystarczy obliczyć wartość szeregu jak w przypadku dodatnich wartości  $x$ , a następnie odwrócić wynik. Zastosowanie powyższych kroków pozwala pozbyć problematycznego odejmowania, co pozytywnie wpłynie na dokładność wyników.

*Kod funkcji realizującej obliczanie wartości  $e^x$ , po uwzględnieniu powyższej modyfikacji:*

```
def find_exp(x: int, epsilon: float) -> float:
    exponents = 0
    element = 1
    i = 1
    is_negative_index = False
    if x < 0:
        is_negative_index = True
        x = abs(x)
    while element > epsilon:
        exponents += element
        element *= x / i
        i += 1

    return 1/exponents if is_negative_index else exponents
```

Zestawienie wyników po zastosowaniu powyższej modyfikacji:

Argument	Wynik algorytmu	Dokładna wartość	Błąd bezwzględny
-1	0.3678794411714423	0.36787944117144233	1.508949536687701e-16
1	2.7182818284590455	2.718281828459045	1.6337129034990842e-16
-5	0.006737946999085467	0.006737946999085469	2.5745579123912054e-16
5	148.4131591025766	148.41315910257657	1.9150397176546985e-16
-10	4.5399929762484854e-05	4.5399929762484875e-05	4.477714137545078e-16
10	22026.465794806714	22026.465794806703	4.954919469581169e-16

### Wnioski:

Po przekształceniu wyrażenia na szereg w którym nie występują wyraz ujemne, zostało zredukowane ryzyko wystąpienia zjawiska „catastrophic cancellation”. Zmiana ta znacząco przyczyniła się do poprawy dokładności wyników w przypadku ujemnych wykładników. Stosując tą modyfikację jesteśmy w stanie uzyskać zbliżoną dokładność wyników nie zależnie do znaku wykładnika.

## 2.

Rozważając dwa ekwiwalentne matematycznie wyrażenia  $x^2 - y^2$  oraz  $(x - y) \cdot (x + y)$  w arytmetyce zmiennoprzecinkowej dokładniejszy wynik otrzymamy stosując wyrażenie drugie, gdyż podczas obliczania wyrażenia postaci  $x^2 - y^2$ , wielkości  $x^2$  i  $y^2$  są obarczone błędami zaokrągleń, ponieważ są wynikami mnożenia zmiennoprzecinkowego. Różnica tych liczb może spowodować wyzerowanie wielu miejsc znaczących wyniku, tak więc spotkamy się ze wspomniany już wcześniej zjawiskiem „catastrophic cancellation”. Natomiast stosując wyrażenie drugie obliczymy różnicę i sumę liczb, nie obarczonych błędami wcześniejszych zaokrągleń. Obliczając iloczyn wcześniej wyliczonej sumy i różnicy, nie narażamy się na wystąpienie tego samego zjawiska co dla pierwszego wyrażenia, a zatem otrzymamy bardziej dokładny wynik.

Szczególnie zauważalna różnica w dokładności obu wyrażeń, będzie widoczna jeśli  $x$  i  $y$  będą bliskimi liczbami, gdyż wówczas podczas odejmowania dwóch liczb zaokrąglonych podczas potęgowania, dojdzie do wyzerowania się wielu znaczących miejsc w otrzymanej różnicy („catastrophic cancellation”). Tak otrzymana wartość będzie znacząco odbiegać od rzeczywistej wartości wyrażenia.

## 3.

Jeśli wykorzystujemy znormalizowany system zmiennie-przecinkowy z podstawą  $\beta = 10$  i dokładnością  $p=3$ , to liczby reprezentowane w tym systemie będą normalizowane do 3 cyfry znaczących.

$ax^2 + bx + c = 0$ , gdzie  $a = 1.22$ ,  $b = 3.34$  i  $c = 2.28$ ,

- a) Przyjmując zadane wartości współczynników równania kwadratowego wartość wyrażenia  $b^2 - 4ac$ , wynosi:

$$\Delta_p = b^2 - 4ac = (3.34)^2 - 4 \cdot 1.22 \cdot 2.28 = 11.1556 - 11.1264 \approx 11.2 - 11.1 \approx 0.1$$

Wartości wyrażeń  $b^2 = 11.1556$  oraz  $4ac = 11.1264$  zostały zgodnie z zadaną precyzją ( $p=3$ ) znormalizowane do odpowiednio 11.2 oraz 11.1.

- b) Dokładna wartość wyróżnika jest równa :

$$\Delta = b^2 - 4ac = (3.34)^2 - 4 \cdot 1.22 \cdot 2.28 = 11.1556 - 11.1264 = 0.0292$$

c) Obliczając błąd względny dla delty powyższego równania otrzymujemy:

$$\frac{|\Delta - \Delta_p|}{\Delta} = \frac{|0.0292 - 0.1|}{0.0292} = 2.42465753 \approx 2.41 = 241\%$$

#### ***4.Bibliografia:***

- Katarzyna Rycerz: Wykład z przedmiotu Metody Obliczeniowe w Nauce i Technice
- Michael T. Heath: Scientific Computing: An Introductory Survey
- [https://en.wikipedia.org/wiki/IEEE\\_7541985](https://en.wikipedia.org/wiki/IEEE_7541985)
- [https://en.wikipedia.org/wiki/Machine\\_epsilon](https://en.wikipedia.org/wiki/Machine_epsilon)
- [https://en.wikipedia.org/wiki/Catastrophic\\_cancellation](https://en.wikipedia.org/wiki/Catastrophic_cancellation)