

Sprawozdanie z zadań z analizy numerycznej

Patryk Blacha, Radosław Szepielak

10 marca 2025

Spis treści

1	Zadanie 1: Obliczanie pochodnej funkcji	2
1.1	Wprowadzenie	2
1.2	Metoda różnicy do przodu	2
1.2.1	Błędy	2
1.2.2	Optymalne h_{\min}	2
1.3	Metoda różnicy centralnej	2
1.3.1	Błędy	2
1.3.2	Optymalne h_{\min}	3
1.4	Implementacja	3
1.4.1	Definicje funkcji	3
1.4.2	Obliczenia	3
1.4.3	Różnice do przodu	3
1.4.4	Różnice centralne	3
1.5	Wyniki	4
1.6	Wykresy	4
1.6.1	Metoda różnicy do przodu	4
1.6.2	Metoda różnicy centralnej	4
1.7	Wnioski	4
2	Zadanie 2: Sumowanie liczb zmiennoprzecinkowych	4
2.1	Wstęp	4
2.2	Implementacja	5
2.2.1	Definicje funkcji	5
2.2.2	Obliczenia	5
2.3	Wykresy	6
2.4	Dyskusja	6
2.5	Wnioski	6
3	Zadanie 3: Stabilność numeryczna wyrażeń	7
3.1	Rozwiązania	7
3.1.1	(a) $\sqrt{x+1} - 1, x \approx 0$	7
3.1.2	(b) $x^2 - y^2, x \approx y$	7
3.1.3	(c) $1 - \cos x, x \approx 0$	7
3.1.4	(d) $\cos^2 x - \sin^2 x, x \approx \frac{\pi}{4}$	7
3.1.5	(e) $\ln x - 1, x \approx e$	8
3.1.6	(f) $e^x - e^{-x}, x \approx 0$	8
4	Zadanie 4: Porównanie sprawności kolektorów słonecznych S1 i S2	8
4.1	Wstęp	8
4.2	Błędy pomiarowe	8
4.3	Propagacja błędów	8
4.4	Analiza porównawcza	9
4.5	Wniosek	9

1 Zadanie 1: Obliczanie pochodnej funkcji

1.1 Wprowadzenie

W tym zadaniu obliczamy przybliżoną wartość pochodnej funkcji $f(x) = \tan(x)$ w punkcie $x_0 = 1.0$ przy użyciu dwóch metod: różnicy do przodu (forward difference) oraz różnicy centralnej (central difference). Następnie porównujemy wyniki z dokładną wartością pochodnej, korzystając z tożsamości:

$$\tan'(x) = 1 + \tan^2(x)$$

1.2 Metoda różnicy do przodu

Wzór na różnicę do przodu jest następujący:

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

1.2.1 Błędy

W metodzie różnicy do przodu występują dwa główne źródła błędów:

- Błąd metody (truncation error):

$$E_{\text{trunc}} \approx \frac{|f''(x)|}{2} \cdot h$$

- Błąd obliczeniowy (roundoff error):

$$E_{\text{round}} \approx \frac{2\epsilon_{\text{mach}}}{h}$$

1.2.2 Optymalne h_{\min}

Teoretycznie optymalna wartość h_{\min} dla metody różnicy do przodu jest dana wzorem:

$$h_{\min} \approx 2\sqrt{\frac{\epsilon_{\text{mach}}}{M}}$$

gdzie $M \approx |f''(x)|$.

1.3 Metoda różnicy centralnej

Wzór na różnicę centralną jest następujący:

$$f'(x) \approx \frac{f(x+h) - f(x-h)}{2h}$$

1.3.1 Błędy

W metodzie różnicy centralnej błędy są następujące:

- Błąd metody (truncation error):

$$E_{\text{trunc}} \approx \frac{|f'''(x)|}{6} \cdot h^2$$

- Błąd obliczeniowy (roundoff error):

$$E_{\text{round}} \approx \frac{\epsilon_{\text{mach}}}{h}$$

1.3.2 Optymalne h_{\min}

Teoretycznie optymalna wartość h_{\min} dla metody różnicy centralnej jest dana wzorem:

$$h_{\min} \approx \left(\frac{3\epsilon_{\text{mach}}}{M} \right)^{1/3}$$

gdzie $M \approx |f'''(x)|$.

1.4 Implementacja

Poniżej znajduje się implementacja obu metod w Pythonie:

1.4.1 Definicje funkcji

```
1 import numpy as np
2 import matplotlib.pyplot as plt
3
4 def f(x):
5     return np.tan(x)
6
7 def f_analytical_prime(x):
8     return 1 + np.tan(x)**2
9
10 def f_second(x):
11     return 2*(1/np.cos(x)**2)*np.tan(x)
12
13 def f_third(x):
14     return 2*(1/np.cos(x)**2) + 6*(np.sin(x)**2)/(np.cos(x)**4)
```

1.4.2 Obliczenia

```
1 x0 = 1.0
2 d_true = f_analytical_prime(x0)
3 eps = np.finfo(float).eps
4 k_vals = np.arange(0, 17)
5 h_vals = 10.0**(-k_vals)
```

1.4.3 Różnice do przodu

```
1 d_forward = (f(x0 + h_vals) - f(x0)) / h_vals
2 error_forward = np.abs(d_forward - d_true)
3 E_trunc_forward = 0.5 * np.abs(f_second(x0)) * h_vals
4 E_round_forward = 2 * eps / h_vals
5
6 idx_min_forward = np.argmin(error_forward)
7 h_min_forward_emp = h_vals[idx_min_forward]
8 E_min_forward_emp = error_forward[idx_min_forward]
9
10 M_forward = np.abs(f_second(x0))
11 h_min_forward_theor = 2 * np.sqrt(eps / M_forward)
```

1.4.4 Różnice centralne

```
1 d_central = (f(x0 + h_vals) - f(x0 - h_vals)) / (2 * h_vals)
2 error_central = np.abs(d_central - d_true)
3 E_trunc_central = (np.abs(f_third(x0)) / 6) * h_vals**2
4 E_round_central = eps / h_vals
5
6 idx_min_central = np.argmin(error_central)
7 h_min_central_emp = h_vals[idx_min_central]
8 E_min_central_emp = error_central[idx_min_central]
9
10 M_central = np.abs(f_third(x0))
11 h_min_central_theor = (3 * eps / M_central)**(1/3)
```

1.5 Wyniki

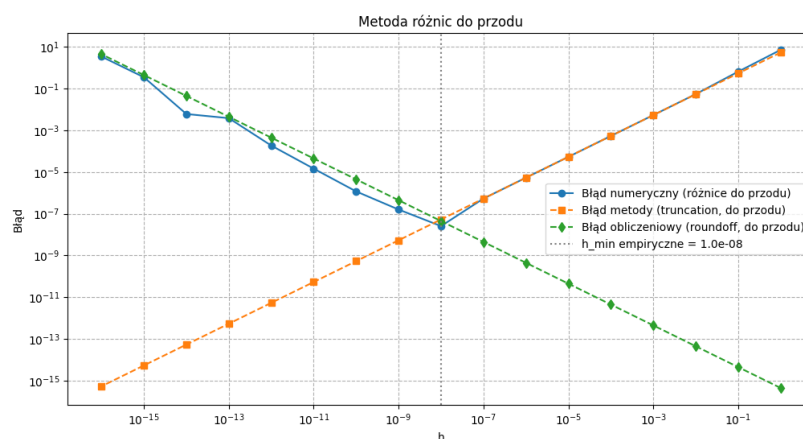
Poniżej przedstawiamy wyniki dla obu metod:

```
1 print("Dla metody r   nic do przodu:")
2 print(f"   Empiryczne h_min: {h_min_forward_emp:.3e}, E(h_min): {E_min_forward_emp:.3e}")
3 print(f"   Teoretyczne h_min (wz r (2)): {h_min_forward_theor:.3e}")
4 print()
5
6 print("Dla metody r   nic centralnych:")
7 print(f"   Empiryczne h_min: {h_min_central_emp:.3e}, E(h_min): {E_min_central_emp:.3e}")
8 print(f"   Teoretyczne h_min (wz r (4)): {h_min_central_theor:.3e}")
9 print()
10
11 if E_min_central_emp < E_min_forward_emp:
12     print("Metoda r   nic centralnych jest dok adniejsza.")
13 else:
14     print("Metoda r   nic do przodu jest dok adniejsza.")
```

1.6 Wykresy

Poniżej znajdują się wykresy błędów dla obu metod:

1.6.1 Metoda różnicy do przodu



Rysunek 1: Błędy dla metody różnicy do przodu

1.6.2 Metoda różnicy centralnej

1.7 Wnioski

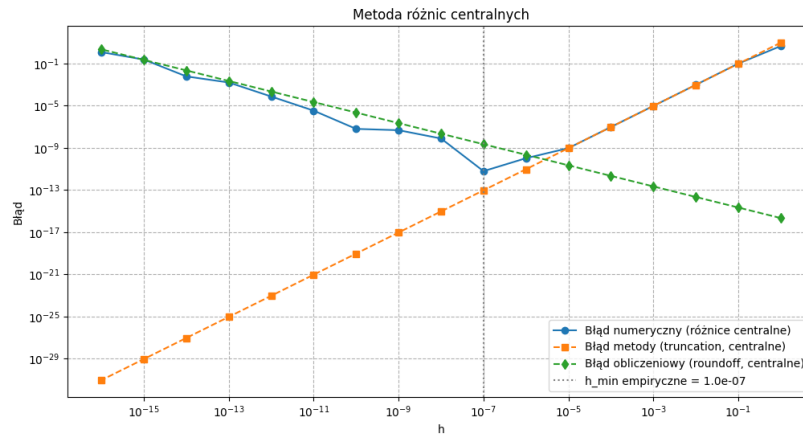
Na podstawie przeprowadzonych obliczeń i analizy błędów można stwierdzić, że metoda różnic centralnych jest dokładniejsza niż metoda różnic do przodu, co potwierdzają zarówno wyniki empiryczne, jak i teoretyczne.

2 Zadanie 2: Sumowanie liczb zmiennoprzecinkowych

2.1 Wstęp

Celem zadania było zbadanie wpływu różnych metod sumowania na dokładność obliczeń przy użyciu liczb zmiennoprzecinkowych pojedynczej precyzji. W szczególności, porównano pięć różnych metod sumowania:

- (a) Sumowanie w kolejności generowania z akumulatorem podwójnej precyzji,



Rysunek 2: Błędy dla metody różnicy centralnej

- (b) Sumowanie w kolejności generowania z akumulatorem pojedynczej precyzji,
- (c) Sumowanie z użyciem algorytmu Kahana,
- (d) Sumowanie w porządku rosnącym,
- (e) Sumowanie w porządku malejącym.

Dla każdej metody obliczono względny błąd sumowania w zależności od liczby elementów $n = 10^k$, gdzie $k = 4, 5, 6, 7, 8$. Jako wartość referencyjną przyjęto sumę obliczoną za pomocą funkcji `math.fsum`, która zapewnia wysoką precyzję.

2.2 Implementacja

Poniżej znajduje się implementacja w Pythonie:

2.2.1 Definicje funkcji

```
1 import numpy as np
2 import matplotlib.pyplot as plt
3 import math
4
5 # Funkcja realizująca sumowanie Kahana w pojedynczej precyzji
6 def kahan_sum(x):
7     suma = np.float32(0.0)
8     komp = np.float32(0.0)
9     for xi in x:
10         y = np.float32(xi - komp)
11         temp = np.float32(suma + y)
12         komp = np.float32((temp - suma) - y)
13         suma = temp
14     return suma
```

2.2.2 Obliczenia

```
1 # Lista wartości n: n = 10^k, k = 4,5,6,7,8
2 n_values = [10 ** k for k in range(4, 9)]
3
4 # Listy do przechowywania wyników dla poszczególnych metod
5 errors_a = [] # (a) podwójna precyzja
6 errors_b = [] # (b) pojedyncza precyzja
7 errors_c = [] # (c) Kahan
8 errors_d = [] # (d) sortowanie rosnące
9 errors_e = [] # (e) sortowanie malejące
10
```

```

11 for n in n_values:
12     # Generujemy n liczb pojedynczej precyzji
13     x = np.random.uniform(0, 1, n).astype(np.float32)
14     # Prawdziwa suma (wysoka precyzja)      wykorzystujemy math.fsum
15     true_sum = math.fsum(x.tolist())
16
17     # (a) Sumowanie w kolejno ci generowania z akumulatorem podwójnej precyzji
18     sum_a = np.cumsum(x, dtype=np.float64)[-1]
19
20     # (b) Sumowanie w kolejno ci generowania z akumulatorem pojedynczej precyzji
21     sum_b = np.cumsum(x, dtype=np.float32)[-1]
22
23     # (c) Sumowanie Kahana      algorytm sumowania z kompensacją, akumulacja w
24     # pojedynczej precyzji
25     sum_c = kahan_sum(x)
26
27     # (d) Sumowanie w kolejno ci rosnących (od najmniejszych do największych)
28     x_sorted_asc = np.sort(x)
29     sum_d = np.cumsum(x_sorted_asc, dtype=np.float32)[-1]
30
31     # (e) Sumowanie w kolejno ci malejących (od największych do najmniejszych)
32     x_sorted_desc = np.sort(x)[::-1]
33     sum_e = np.cumsum(x_sorted_desc, dtype=np.float32)[-1]
34
35     # Obliczamy względną błąd dla każdej metody: |suma_metody - true_sum| / |
36     # true_sum|
37     err_a = abs(sum_a - true_sum) / abs(true_sum)
38     err_b = abs(sum_b - true_sum) / abs(true_sum)
39     err_c = abs(sum_c - true_sum) / abs(true_sum)
40     err_d = abs(sum_d - true_sum) / abs(true_sum)
41     err_e = abs(sum_e - true_sum) / abs(true_sum)
42
43     errors_a.append(err_a)
44     errors_b.append(err_b)
45     errors_c.append(err_c)
46     errors_d.append(err_d)
47     errors_e.append(err_e)
48
49     print(f"n = {n:>10}:")
50     print(f"  (a) podwójna precyzja: suma = {sum_a:.8e}, błąd = {err_a:.8e}")
51     print(f"  (b) pojedyncza precyzja: suma = {sum_b:.8e}, błąd = {err_b:.8e}")
52     print(f"  (c) Kahan:      suma = {sum_c:.8e}, błąd = {err_c:.8e}")
53     print(f"  (d) sort. rosnący:   suma = {sum_d:.8e}, błąd = {err_d:.8e}")
54     print(f"  (e) sort. malejący:  suma = {sum_e:.8e}, błąd = {err_e:.8e}")
55     print()

```

2.3 Wykresy

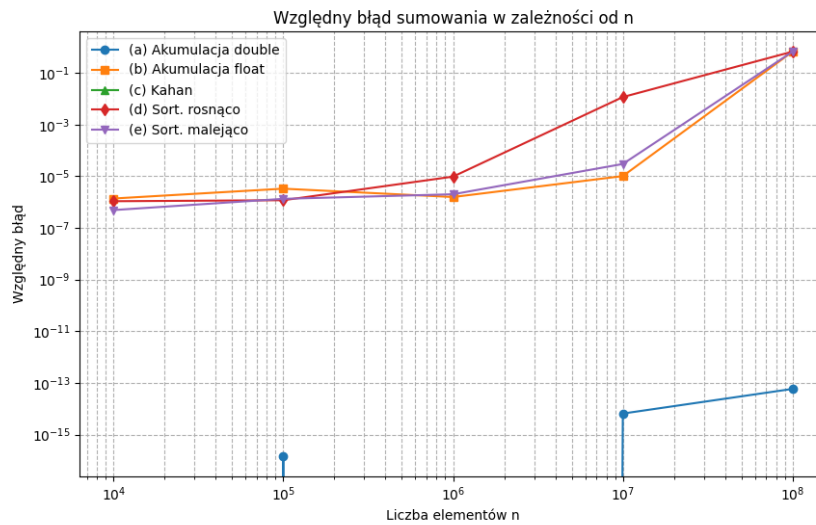
Poniżej znajduje się wykres względnego błędu w zależności od liczby elementów n :

2.4 Dyskusja

- Metoda (a) z użyciem akumulatora podwójnej precyzji charakteryzuje się najmniejszym błędem, co wynika z większej precyzji obliczeń.
- Metoda (b) z użyciem akumulatora pojedynczej precyzji wykazuje większy błąd niż metoda (a), co jest spowodowane ograniczoną precyzją liczb pojedynczej precyzji.
- Algorytm Kahana (metoda c) pozwala na znaczące zmniejszenie błędu w porównaniu do standardowego sumowania w pojedynczej precyzji, dzięki kompensacji błędów zaokrągleń.
- Sumowanie w porządku rosnącym (metoda d) daje lepsze wyniki niż sumowanie w porządku malejącym (metoda e), co jest zgodne z oczekiwaniami, ponieważ sumowanie od najmniejszych do największych wartości minimalizuje błędy zaokrągleń.

2.5 Wnioski

- Użycie akumulatora podwójnej precyzji (metoda a) jest najbardziej efektywnym sposobem na zmniejszenie błędu sumowania.



Rysunek 3: Względny błąd sumowania w zależności od liczby elementów n dla różnych metod sumowania.

- Algorytm Kahana (metoda c) jest skutecznym narzędziem do poprawy dokładności sumowania w pojedynczej precyzji.
- Sumowanie w porządku rosnącym (metoda d) jest lepsze niż sumowanie w porządku malejącym (metoda e) pod względem dokładności.

3 Zadanie 3: Stabilność numeryczna wyrażeń

3.1 Rozwiązania

3.1.1 (a) $\sqrt{x+1} - 1$, $x \approx 0$

Aby uniknąć zjawiska kancelacji, mnożymy przez sprzężenie:

$$\sqrt{x+1} - 1 = \frac{(\sqrt{x+1} - 1)(\sqrt{x+1} + 1)}{\sqrt{x+1} + 1} = \frac{x}{\sqrt{x+1} + 1}.$$

3.1.2 (b) $x^2 - y^2$, $x \approx y$

Różnicę kwadratów zapisujemy jako iloczyn:

$$x^2 - y^2 = (x - y)(x + y).$$

Dzięki temu eliminujemy problem z odejmowaniem bliskich wartości.

3.1.3 (c) $1 - \cos x$, $x \approx 0$

Stosujemy tożsamość trygonometryczną:

$$1 - \cos x = 2 \sin^2 \frac{x}{2}.$$

3.1.4 (d) $\cos^2 x - \sin^2 x$, $x \approx \frac{\pi}{4}$

Korzystamy z tożsamości:

$$\cos^2 x - \sin^2 x = \cos 2x.$$

Dla $x \approx \frac{\pi}{4}$, mamy $\cos 2x = -\sin(2(x - \frac{\pi}{4}))$, co jest bardziej stabilne numerycznie.

3.1.5 (e) $\ln x - 1$, $x \approx e$

Przekształcamy wyrażenie:

$$\ln x - 1 = \ln \frac{x}{e}.$$

Dzięki temu unikamy odejmowania wartości bliskich sobie.

3.1.6 (f) $e^x - e^{-x}$, $x \approx 0$

Używamy funkcji hiperbolicznej:

$$e^x - e^{-x} = 2 \sinh x.$$

Dla małych x rozwijamy $\sinh x$ w szereg Taylora:

$$\sinh x = x + \frac{x^3}{3!} + \mathcal{O}(x^5).$$

To pozwala na stabilniejsze obliczenia numeryczne.

4 Zadanie 4: Porównanie sprawności kolektorów słonecznych S1 i S2

4.1 Wstęp

Efektywność kolektora słonecznego dana jest wzorem:

$$\eta = \frac{K Q T_d}{I},$$

gdzie:

- K – stała znana z dużą dokładnością,
- Q – objętość przepływu,
- T_d – różnica temperatur,
- I – natężenie promieniowania.

Dla dwóch kolektorów obliczono:

$$\eta_{S1} = 0.76 \quad \text{oraz} \quad \eta_{S2} = 0.70.$$

4.2 Błędy pomiarowe

Wielkości Q , T_d oraz I zmierzono z następującymi względnymi błędami:

Wielkość	S1	S2
Q	1.5%	0.5%
T_d	1.0%	1.0%
I	3.6%	2.0%

4.3 Propagacja błędów

Zakładając, że K jest znane bardzo dokładnie, względny błąd efektywności η wyraża się przy pomocy zasady dla iloczynu i ilorazu:

$$\frac{\Delta \eta}{\eta} \approx \frac{\Delta Q}{Q} + \frac{\Delta T_d}{T_d} + \frac{\Delta I}{I}.$$

Kolektor S1

Dla S1:

$$\frac{\Delta\eta_{S1}}{\eta_{S1}} = 1.5\% + 1.0\% + 3.6\% = 6.1\%.$$

Błąd bezwzględny wynosi:

$$\Delta\eta_{S1} = 0.76 \times 0.061 \approx 0.046.$$

Przyjmując, że rzeczywista efektywność mieści się w przedziale:

$$\eta_{S1} \in [0.76 - 0.046, 0.76 + 0.046] \approx [0.714, 0.806].$$

Kolektor S2

Dla S2:

$$\frac{\Delta\eta_{S2}}{\eta_{S2}} = 0.5\% + 1.0\% + 2.0\% = 3.5\%.$$

Błąd bezwzględny wynosi:

$$\Delta\eta_{S2} = 0.70 \times 0.035 \approx 0.0245.$$

Przyjmując przedział niepewności:

$$\eta_{S2} \in [0.70 - 0.0245, 0.70 + 0.0245] \approx [0.6755, 0.7245].$$

4.4 Analiza porównawcza

Dla kolektora S1 efektywność mieści się w przedziale $[0.714, 0.806]$, natomiast dla S2 w przedziale $[0.6755, 0.7245]$. Zauważamy, że:

- Górna granica dla S2 wynosi około 0.7245,
- Dolna granica dla S1 wynosi około 0.714.

Przedziały te częściowo się nakładają, co oznacza, że w wyniku błędów pomiarowych nie możemy z całą pewnością stwierdzić, że rzeczywista efektywność S1 jest większa niż S2.

4.5 Wniosek

Pomimo nominalnie wyższej sprawności S1 (0.76) w porównaniu do S2 (0.70), niepewności pomiarowe (szczególnie większy błąd w I dla S1) powodują, że przedziały niepewności dla obu kolektorów się nakładają. Dlatego na podstawie dostępnych danych nie możemy być pewni, że kolektor S1 ma większą sprawność niż kolektor S2.