

Algorytmy numeryczne

Zadanie 1

Aleksander Kosma

grupa 1 tester-programista

16 Październik 2017

1 Sumowanie szeregów potęgowych

Opracowanie prezentuje próby obliczenia funkcji \arctan , przy pomocy języka java. Wykorzystano w tym celu dwa sposoby na obliczenie tej wartości. Są to:

-zsumowanie elementów wyliczonych z szeregu potęgowego:

$$\arctg(x) = \sum_{n=0}^{\infty} \frac{(-1)^n}{2n+1} x^{2n+1}$$

-zsumowanie elementów wyliczonych na podstawie poprzednika:

$$a_{n+1} = a * -\frac{x^2 * (2n+1)}{2n+3}$$

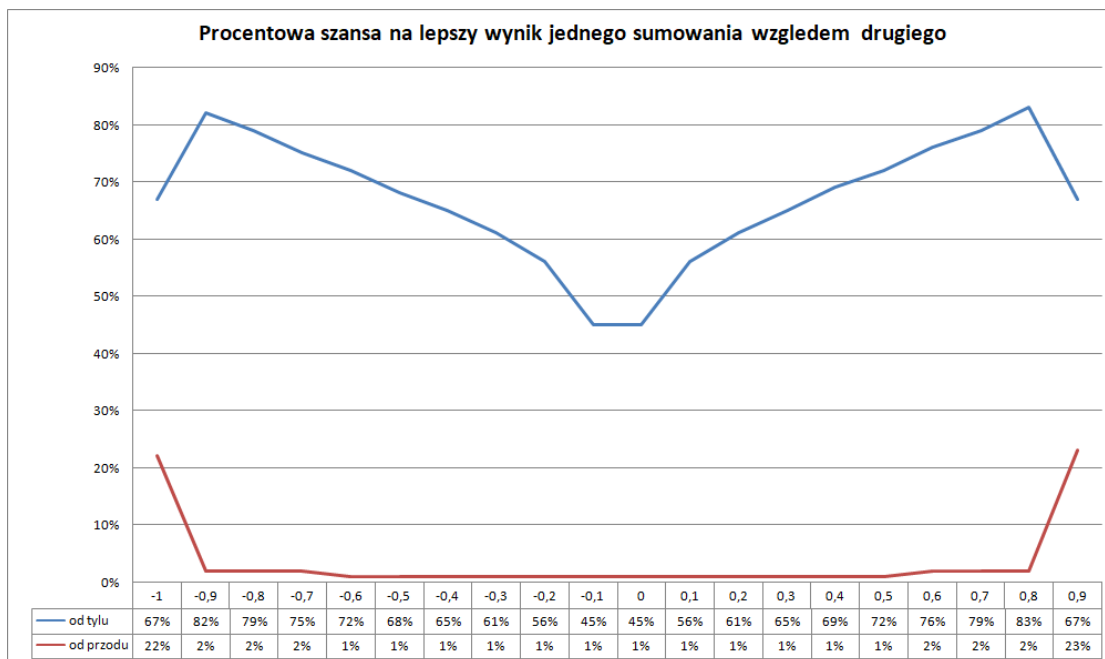
Ze względu na charakter obliczeń podzieliłem owe zsumowania na kolejne dwa:

-zsumowanie elementów licząc od tyłu

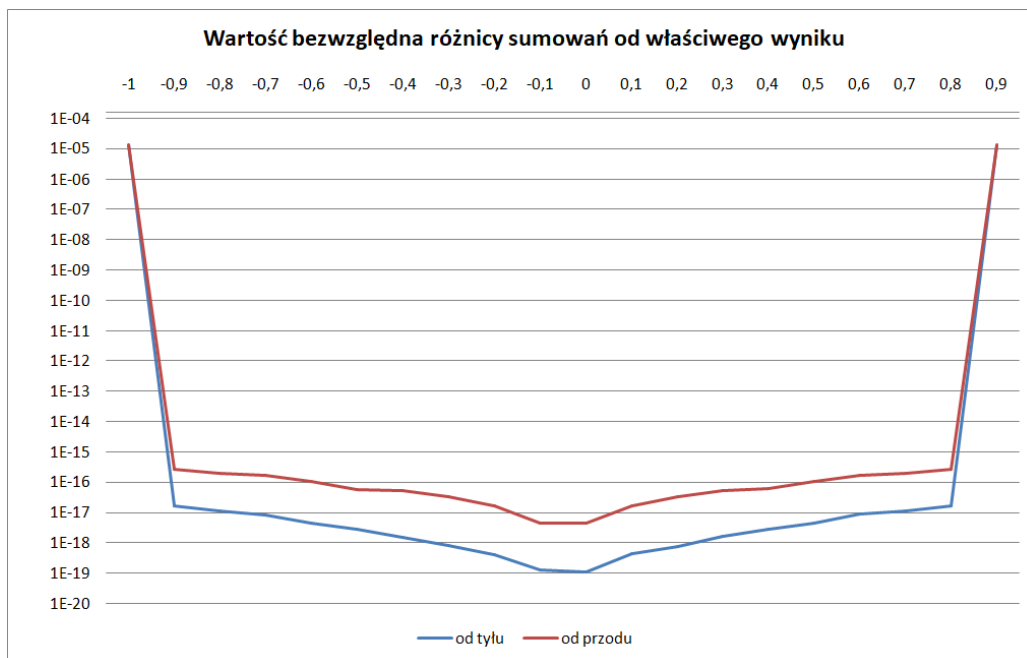
-zsumowanie elementów licząc od przodu

Wnioski zostały obliczone na zmiennej typu double. Również ze względu na zbieżność szeregu w $[-1, 1]$ liczyłem tylko dla tego zakresu.

Wykres przedstawia który rodzaj sumowania ma większą szansę na precyzyjniejszy wynik względem drugiego. Każda kolumna wartości została wyliczona dla 100 tysięcy liczb między podanymi przedziałami. Razem 2 miliony próbek.



Drugi wykres prezentuje bezwzględną różnicę między wynikiem biblioteki a sumowaniami. Środkowa część wykresu ukazuje wyższość sumowania od tyłu. Jest ono precyzyjniejsze nawet o 2 rzędy wielkości.



Takie rezultaty wynikają z natury działania zmiennej double. W przypadku sumowania liczby większej do liczby mniejszej (gdzie różnica wynosi parę rzędów wielkości), zmienna odrzuca końcowe cyfry mniejszej liczby by pomieścić najistotniejsze wartości. Stąd utrata precyzji. W momencie kiedy dwie liczby mają podobną wartość, ich suma potrzebować będzie podobnej precyzji do zapisania wyniku. W konsekwencji nie odrzuci danych. Podsumowując kwestie precyzji wyniku. Opierając się o zmienną double, można stwierdzić, że precyzja wyniku waha się między 14, a 16 miejscem po przecinku.

Na obrazku poniżej znajduje się praktyczny przykład. Dla uproszczenia zagadnienia jest on prezentowany na zmiennej float(32 bity) która zazwyczaj mieści w sobie 8 cyfr dziesiętnych.

<p>tak jest</p> <p>8 cyfr</p> $ \begin{array}{r} 10,0010000 \\ + 0,0000001 \\ \hline 10,0010000 \end{array} $		<p>tak powinno być</p> <p>8 cyfr</p> $ \begin{array}{r} 10,0010000 \\ + 0,0000001 \\ \hline 10,0010001 \end{array} $
--	--	---

By rozwiązać problem niskiej precyzji, który zaistniał w przypadku double, trzeba mieć kontrolę nad precyzją obliczeń. Rozwiązaniem może okazać się klasa BigDecimal. Klasa ta może przechować nieograniczoną wartość liczbową, ogranicza ją jedynie skala która osiąga wartość 32-bitowego integera, czyli trochę ponad dwa miliardy miejsc.

Udało mi się osiągnąć praktycznie nieograniczoną precyzję. Wyniki oparte o BigDecimal porównywałem z kalkulatorem online Wolfram(computational knowledge engine). Przykładowy wypis programu wraz ze wszystkimi rodzajami obliczeń i ich wynikami. W tym przypadku podana precyzja BigDecimal i Wolfram to 64 miejsca po przecinku.

```

Output of Math library:      0.2303949468807241
Output of Wolfram:          0.2303949468807240960351954018340492626054069551810387153331776640

Sum forward McLourin BigDecimal: 0.2303949468807240960351954018340492626054069551810387153331776640
Sum backward McLourin BigDecimal: 0.2303949468807240960351954018340492626054069551810387153331776640

Sum forward McLourin Double:    0.23039494688072407
Sum backward McLourin Double:   0.2303949468807241

Sum forward next BigDecimal:    0.2303949468807240960351954018340492626054069551810387153331776640
Sum backward next BigDecimal:   0.2303949468807240960351954018340492626054069551810387153331776640

Sum forward next Double:       0.23039494688072407
Sum backward next Double:      0.2303949468807241

```