

Algorytm Molera-Morissona

Aleksander Czeszejko-Sochacki

12 listopada 2017

1 Wprowadzenie

Algorytm Molera-Morissona to iteracyjna metoda obliczania wyrażeń postaci $\sqrt{a^2 + b^{2^1}}$. Tego typu funkcje pojawiają się nadzwyczaj często, między innymi przy obliczaniu:

- przeciwprostokątnej
- modułu liczby zespolonej
- normy wektora
- zmiany układu współrzędnych z euklidesowych na biegunowe

1.1 Pseudokod

Algorithm 1 Algorytm iteracyjny Molera-Morissona

```
p := max(|a|, |b|)
q := min(|a|, |b|)
while wartość q jest znacząca do
  r := ( $\frac{q}{p}$ )2
  s :=  $\frac{r}{4+r}$ 
  p := p + 2 * s * p
  q := s * q
end while
return p
```

Od razu możemy zauważyć zalety algorytmu - nie wymaga on używania skomplikowanych operacji arytmetycznych, a każda iteracja składa się z kilku elementarnych działań i przypisów.

¹a, b ∈ ℝ, najlepiej gdy a, b ∈ ℚ, ponieważ przybliżona reprezentacja liczb niewymiernych może powiększać ewentualny błąd

1.2 Liczba iteracji

Dowodzi się, że po 3 iteracjach algorytmu błąd względny jest na poziomie $0.5 \cdot 10^{-20}$, co uznamy za wystarczającą dokładność w naszym artykule. W przypadku chęci szerszej analizy zbieżności i numerycznej poprawności odsyłamy Czytelnika do rozdziału czwartego https://blogs.mathworks.com/images/cleve/moler_morrison.pdf

1.3 Opis użytych metod

W swoim sprawozdaniu zaprezentuję obliczenia wykonane w języku Julia. Przeprowadzone zostały w dwóch arytmetykach:

1. double, tzn. Float64
2. BigFloat(128)

Dlaczego takie arytmetyki? Precyzja arytmetyki Float64 może okazać się niewystarczająca, ponieważ wynosi ona

$$\frac{1}{2}2^{-52} = \frac{1}{8}(2^{-10})^5 \approx 10^{-16}$$

Numer iteracji	Błąd względny dla pary (3, 4)	Błąd względny dla pary (-5, 12)
2.	5.162349481224737e-9,	5.241619104568739e-13
3.	1.7763568394002506e-16,	1.3664283380001927e-16

Powyższe wyniki zgadzają się z naszym wnioskiem. Powyższe błędy po trzeciej iteracji powinny być rzędu najwyżej 10^{-21} . Powinniśmy dobrać taką arytmetykę, której precyzja będzie mniejsza niż $\frac{1}{2}2^{-20}$. Zakładając, że zwiększając w Julii liczbę bitów przydzielanych dla arytmetyki dwukrotnie, w przybliżeniu dwukrotnie zwiększa się liczba bitów przydzielanych na mantysę, precyzja arytmetyki BigFloat(128) wynosi około:

$$\frac{1}{2}2^{-100} = \frac{1}{2}(2^{-10})^{10} \approx 10^{-31}$$

co zdecydowanie wystarczy na wychwycenie ewentualnego błędu algorytmu.

Za wartość dokładną \sqrt{a} będę uznawał wartość funkcji `sqrt(a)` w arytmetyce BigFloat(128)

2 Analiza Algorytmu

2.1 Testownie algorytmu dla wybranych danych

Poniżej przeprowadzimy testy dla następujących danych:

(3, 4), (-5, 12), (7, -24), (1, 1), (1000000000, 2), (71075075103, 1000000000)

Testy (trzy iteracje:)

[illegible][illegible][illegible][illegible][illegible][illegible]

(1, 1), dokładna wartość: 1.414213562373095048801688724209698078569
Równe obie dane, nie dają całkowitego wyniku

[illegible]

Duża pierwsza dana, duża względna różnica, nie dają całkowitego wyniku

Numer iteracji	Wynik algorytmu (Float64)	Wynik algorytmu (BigFloat(128))
1	1.0e9	1.00000000000000000019999999999999998001e+09
2	1.0e9	1.00000000000000000001999999999999998001e+09
3	1.0e9	1.00000000000000000001999999999999998001e+09

7.108210956982840179871838428090505048092e+10
Duże obie dane, nie dają całkowitego wyniku

Numer iteracji	Wynik algorytmu (Float64)	Wynik algorytmu (BigFloat(128))
1	7.108210956981117e10	7.108210956981117607726957458289901508523e+10
2	7.10821095698284e10	7.108210956982840179871838428090505048092e+10
3	7.10821095698284e10	7.108210956982840179871838428090505048132e+10

Możemy w powyższych przykładach zaobserwować następujące zjawiska:

- Zmniejszenie precyzji arytmetyki powoduje zwiększenie liczby cyfr znaczących
- Algorytm nie wyliczył oczekiwanych wartości całkowitych dla pierwszych trzech par

Więcej obserwacji będziemy w stanie wyciągnąć po analizie błędów względnych.

Błędy względne ²[illegible]

²dla przejrzystości zmniejszyłem liczbę nieznaczących zer w błędach w arytmetyce BigFloat(128); \tilde{y} - wynik algorytmu, y - wartość obliczona przy pomocy funkcji sqrt()

Obserwacje (w powyższych przykładach):

- Błędy w arytmetyce Float64 są najwyżej rzędu 10^{-16} , co zgadza się z przybliżoną wartością precyzji tej arytmetyki obliczoną we wprowadzeniu
- W niektórych przypadkach błędy są zerowe (lub mniejsze niż precyzja arytmetyki)
- W arytmetyce BigFloat(128) wszystkie błędy są zerowe (lub mniejsze niż precyzja arytmetyki)
- Błędy w arytmetyce BigFloat(128) są niewiększe niż we Float64

Hipoteza 1 *Dla dostatecznie dużej precyzji algorytm zwraca matematycznie dokładny wynik dla dowolnych danych (nie wiadomo, czy taka precyzja jest osiągalna na każdej maszynie).*

2.2 Ciekawostka

Algorytm zachowuje się ciekawie, gdy opuścimy wartości bezwzględne w inicjalizacji zmiennych. Po trzeciej iteracji

$$\sqrt{(-7)^2 + 24^2} \neq \sqrt{7^2 + (-24)^2}$$

natomiast po czwartej

$$\sqrt{(-7)^2 + 24^2} = \sqrt{7^2 + (-24)^2}$$

Zachęcam Czytelnika do przeprowadzenia własnych eksperymentów.

2.3 Zalety stosowania algorytmu

Algorytm iteracyjny Molera-Morissona ma bardzo przyzwoity wykładnik zbieżności - z każdą iteracją liczba cyfr znaczących zwiększa się około trzykrotnie. Niewątpliwą jego przewagą nad liczeniem wyrażenia przy pomocy funkcji `sqr()` jest fakt, że w algorytmie nie podnosimy danych do kwadratu, przez co ma dużo mniejszą złożoność pamięciową. Złożoność czasowa zależy od czasu wykonywania elementarnych operacji (chyba że chcemy wykonać więcej niż trzy iteracje).

3 Zasosowanie algorytmu Molera-Morissona do obliczania normy wektora

Normą wektora x nazywamy wyrażenie:

$$\|x\| = \sqrt{\sum_{i=1}^n x_i^2}$$

Możemy ją jednak zapisać następująco:

$$\begin{aligned}
\sqrt{\sum_{i=1}^n x_i^2} &= \sqrt{\underbrace{x_1^2 + x_2^2 + \dots + x_n^2}_n} \\
&= \sqrt{\sqrt{(\sqrt{x_1^2 + x_2^2})^2 + x_3^2} + \dots + x_n^2} \\
&= \sqrt{\sqrt{\sqrt{\sqrt{(\sqrt{x_1^2 + x_2^2})^2 + x_3^2} + \dots + x_{n-1}^2} + x_n^2}}
\end{aligned} \tag{1}$$

3.1 Pseudokod obliczania normy wektora

Powyższe rozpisanie umożliwia nam zapisanie algorytmu:

Algorithm 2 Algorytm iteracyjny Molera-Morrissona

```

result =  $x_1$ 
for i in 2, 3, ..., n do
    result = Moler-Morrison(result,  $x_i$ )
end for
return result

```

Dokładność otrzymanego wyniku będzie zależna od dokładności samego algorytmu, jak i kolejności sumowania, jednak problem kolejności sumowania dla uzyskania najmniejszego błędu to temat obszerny i nie będzie on rozważany w tym artykule.

Literatura

- [1] C. Moler, D. Morrison: Replacing Square Roots by Pythagorean Sums