

Mnożenie macierzy

Maciej Borowiec

26.10.2025

1 Algorytmy

1.1 Rekurencyjny algorytm Binét'a

1.1.1 Przebieg algorytmu

- Jeśli dowolny wymiar, dowolnej macierzy wejściowej, jest równy 1 - mnożymy macierze iteratywnie
- W przeciwnym wypadku:
 - Dzielimy macierze na podmacierze:

$$A = \begin{bmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \end{bmatrix}, B = \begin{bmatrix} B_{11} & B_{12} \\ B_{21} & B_{22} \end{bmatrix}$$

- Obliczamy podmacierze macierzy wynikowej (mnożenie macierzy wykonujemy rekurencyjnie):

$$C_{11} = A_{11}B_{11} + A_{12}B_{21}$$

$$C_{12} = A_{11}B_{12} + A_{12}B_{22}$$

$$C_{21} = A_{21}B_{11} + A_{22}B_{21}$$

$$C_{22} = A_{21}B_{12} + A_{22}B_{22}$$

- Zwracamy macierz wynikową:

$$C = \begin{bmatrix} C_{11} & C_{12} \\ C_{21} & C_{22} \end{bmatrix}$$

1.1.2 Kod

```
1 def binet(A, B):
2     flop_count = 0
3
4     def f(A, B):
5         nonlocal flop_count
6         # initialize result matrix as zeros
7         C = np.zeros((A.shape[0], B.shape[1]))
8         # if matrices can be split further, split them
9         if A.shape[0] > 1 and A.shape[1] > 1 and B.shape[1] > 1:
10             # split matrices
11             A11, A12, A21, A22 = split_matrix(A)
12             B11, B12, B21, B22 = split_matrix(B)
13             # calculate each part of the result recursively
14             C11 = f(A11, B11) + f(A12, B21)
15             C12 = f(A11, B12) + f(A12, B22)
16             C21 = f(A21, B11) + f(A22, B21)
17             C22 = f(A21, B12) + f(A22, B22)
18             flop_count += el_count(C11) + el_count(C12) + el_count(
19                 C21) + el_count(C22)
20             # write the result into new matrix
21             half_r, half_c = C11.shape
22             C[:half_r, :half_c] = C11
23             C[:half_r, half_c:] = C12
24             C[half_r:, :half_c] = C21
25             C[half_r:, half_c:] = C22
26             # if you can no longer split matrices, do normal matrix
27             # multiplication
28             else:
29                 for r in range(C.shape[0]):
30                     for c in range(C.shape[1]):
31                         for i in range(A.shape[1]):
32                             C[r,c] += A[r,i]*B[i,c]
33                             flop_count += 2
34             # return resulting matrix
35             return C
36
37     return f(A, B), flop_count
```

Listing 1: Algorytm Binét'a

1.1.3 Oszacowanie złożoności obliczeniowej

Dla każdego wywołania funkcji wykonujemy 4 dodawania macierzy i 8 rekurencyjnych mnożeń macierzy o 2 razy mniejszych wymiarach. Zatem możemy zapisać złożoność obliczeniową jako:

$$T(n) = 8T(n/2) + 4n^2$$

$$T(1) = 1$$

Możemy zatem zapisać tę złożoność jako sumę:

$$T(n) = 1 \cdot 8^{\log_2 n} + \sum_{i=0}^{\log_2 n - 1} 8^i \cdot 4\left(\frac{n}{2^i}\right)^2$$

Używając $a^{\log_c b} = b^{\log_c a}$ możemy uprościć równanie:

$$T(n) = n^{\log_2 8} + 4 \sum_{i=0}^{\log_2 n - 1} 8^i \cdot \frac{n^2}{4^i} = n^3 + 4n^2 \sum_{i=0}^{\log_2 n - 1} 2^i$$

Możemy teraz użyć wzoru na sumę ciągu geometrycznego:

$$T(n) = n^3 + 4n^2 \frac{1 - 2^{\log_2 n}}{1 - 2} = n^3 + 4n^2 \frac{1 - n}{-1} = 5n^3 - 4n^2$$

Zatem złożoność obliczeniowa wynosi $O(n^3)$.

1.2 Rekurencyjny algorytm Strassena

1.2.1 Przebieg algorytmu

- Dokonujemy paddingu - dodajemy zera do macierzy wejściowych, by miały one wymiar $2^n \times 2^n$
- Jeśli dowolny wymiar, dowolnej macierzy wejściowej, jest równy 1 - mnożymy macierze iteratywnie
- W przeciwnym wypadku:
 - Dzielimy macierze na podmacierze:

$$A = \begin{bmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \end{bmatrix}, \quad B = \begin{bmatrix} B_{11} & B_{12} \\ B_{21} & B_{22} \end{bmatrix}$$

- Obliczamy pomocnicze zmienne (mnożenie macierzy wykonujemy rekurencyjnie):

$$M_1 = (A_{11} + A_{22})(B_{11} + B_{22})$$

$$M_2 = (A_{21} + A_{22})B_{11}$$

$$M_3 = A_{11}(B_{12} - B_{22})$$

$$M_4 = A_{22}(B_{21} - B_{11})$$

$$M_5 = (A_{11} + A_{12})B_{22}$$

$$M_6 = (A_{21} - A_{11})(B_{11} + B_{12})$$

$$M_7 = (A_{12} - A_{22})(B_{21} + B_{22})$$

- Obliczamy podmacierze macierzy wynikowej:

$$C_{11} = M_1 + M_4 - M_5 + M_7$$

$$C_{12} = M_3 + M_5$$

$$C_{21} = M_2 + M_4$$

$$C_{22} = M_1 - M_2 + M_3 + M_6$$

- Zwracamy macierz wynikową:

$$C = \begin{bmatrix} C_{11} & C_{12} \\ C_{21} & C_{22} \end{bmatrix}$$

1.2.2 Kod

```
1 def strassen(A, B):
2     flop_count = 0
3     # pad matrices
4     n = A.shape[0]
5     size_2n = 1
6     while size_2n < n:
7         size_2n *= 2
8     A_padded = np.zeros((size_2n, size_2n))
9     B_padded = np.zeros((size_2n, size_2n))
10    A_padded[:n, :n] = A
11    B_padded[:n, :n] = B
12
13    def f(A, B):
14        nonlocal flop_count
15        n = A.shape[0]
16        # initialize result matrix as zeros
17        C = np.zeros((n,n))
18        # if matrices can be split further, split them
19        if n > 1:
20            # split matrices
21            A11, A12, A21, A22 = split_matrix(A)
22            B11, B12, B21, B22 = split_matrix(B)
23            # calculate needed values recursively
24            M1 = f(A11+A22, B11+B22)
25            M2 = f(A21+A22, B11)
26            M3 = f(A11, B12-B22)
27            M4 = f(A22, B21-B11)
28            M5 = f(A11+A12, B22)
29            M6 = f(A21-A11, B11+B12)
30            M7 = f(A12-A22, B21+B22)
31            flop_count += (n//2)**2 * 10
32            # add values to create each part of result matrix
33            C11 = M1 + M4 - M5 + M7
34            C12 = M3 + M5
35            C21 = M2 + M4
36            C22 = M1 - M2 + M3 + M6
37            flop_count += (n//2)**2 * 8
38            # write the result into new matrix
39            half_r, half_c = C11.shape
40            C[:half_r, :half_c] = C11
41            C[:half_r, half_c:] = C12
42            C[half_r:, :half_c] = C21
43            C[half_r:, half_c:] = C22
44            # if you can no longer split matrices, it's 1x1 by 1x1
45        else:
46            C[0,0] = A[0,0]*B[0,0]
47            flop_count += 1
48        # return resulting matrix
49        return C
50
51    return f(A_padded, B_padded)[:n, :n], flop_count
```

Listing 2: Algorytm Strassena

1.2.3 Oszacowanie złożoności obliczeniowej

Dla każdego wywołania funkcji wykonujemy 18 dodawań/odejmowań macierzy i 7 rekurencyjnych mnożeń macierzy o 2 razy mniejszych wymiarach. Zatem możemy zapisać złożoność obliczeniową jako:

$$T(n) = 7T(n/2) + 18n^2$$

$$T(1) = 1$$

Możemy zatem zapisać tę złożoność jako sumę:

$$T(n) = 1 \cdot 7^{\log_2 n} + \sum_{i=0}^{\log_2 n - 1} 7^i \cdot 18\left(\frac{n}{2^i}\right)^2$$

Używając $a^{\log_c b} = b^{\log_c a}$ możemy uprościć równanie:

$$T(n) = n^{\log_2 7} + 18 \sum_{i=0}^{\log_2 n - 1} 7^i \cdot \frac{n^2}{4^i} = n^{\log_2 7} + 18n^2 \sum_{i=0}^{\log_2 n - 1} \left(\frac{7}{4}\right)^i$$

Możemy teraz użyć wzoru na sumę ciągu geometrycznego:

$$T(n) = n^{\log_2 7} + 18n^2 \frac{1 - \left(\frac{7}{4}\right)^{\log_2 n}}{1 - \frac{7}{4}} = n^{\log_2 7} - 18n^2(1 - n^{\log_2 7/4})$$

$$T(n) = n^{\log_2 7} - 18n^2(1 - n^{\log_2 7/4}) = n^{\log_2 7} + 18n^{\log_2 7} - 18n^2 = 19n^{\log_2 7} - 18n^2$$

Zatem złożoność obliczeniowa wynosi $O(n^{\log_2 7}) \approx O(n^{2.8074})$.

1.3 Algorytm wykorzystujący faktoryzacje znalezione przez AI

1.3.1 Przebieg algorytmu

Algorytm dokonuje mnożenia rekurencyjnie za pomocą metody Binét'a, natomiast warunek stopu jest osiągany dla dowolnej faktoryzacji, którą znalazł Alpha Tensor. Wtedy wykorzystujemy dane na temat faktoryzacji, by obliczyć dane podmacierze, i co za tym idzie macierz wynikową.

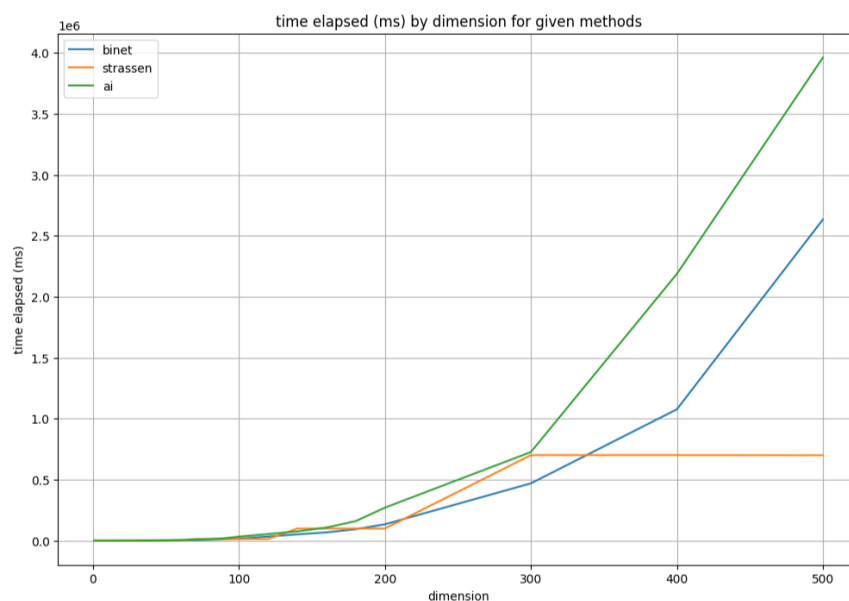
1.3.2 Kod

```
1 def ai(A, B):
2     flop_count = 0
3
4     def f(A, B):
5         nonlocal flop_count
6         # initialize result matrix as zeros
7         C = np.zeros((A.shape[0], B.shape[1]))
8         shape_string = f"{A.shape[0]},{A.shape[1]},{B.shape[1]}"
9         # if matrix sizes are in factorizations, multiply them
10        if shape_string in FACTS:
11            u, v, w = FACTS[shape_string]
12            # create vectors representing matrices
13            A_vec = A.reshape(-1)
14            B_vec = B.reshape(-1)
15            # calculate Ms
16            A_part, f1 = mat_by_vec(u.T, A_vec)
17            B_part, f2 = mat_by_vec(v.T, B_vec)
18            M = A_part * B_part
19            # create a vector of result and reshape it
20            C_vec, f3 = mat_by_vec(w, M)
21            flop_count += f1 + f2 + f3 + len(M)
22            C = C_vec.reshape(A.shape[0], B.shape[1], order='F')
23            # if they cannot be split further, multiply them normally
24            elif A.shape[0] == 1 or A.shape[1] == 1 or B.shape[1] == 1:
25                for r in range(C.shape[0]):
26                    for c in range(C.shape[1]):
27                        for i in range(A.shape[1]):
28                            C[r,c] += A[r,i]*B[i,c]
29                            flop_count += 2
30            # if not, split them using classic binet
31            else:
32                # split matrices
33                A11, A12, A21, A22 = split_matrix(A)
34                B11, B12, B21, B22 = split_matrix(B)
35                # calculate each part of the result recursively
36                C11 = f(A11, B11) + f(A12, B21)
37                C12 = f(A11, B12) + f(A12, B22)
38                C21 = f(A21, B11) + f(A22, B21)
39                C22 = f(A21, B12) + f(A22, B22)
40                flop_count += el_count(C11) + el_count(C12) + el_count(
41                C21) + el_count(C22)
42                # write the result into new matrix
43                half_r, half_c = C11.shape
44                C[:half_r, :half_c] = C11
45                C[:half_r, half_c:] = C12
46                C[half_r:, :half_c] = C21
47                C[half_r:, half_c:] = C22
48            # return resulting matrix
49            return C
50        return f(A, B), flop_count
```

Listing 3: Algorytm AI

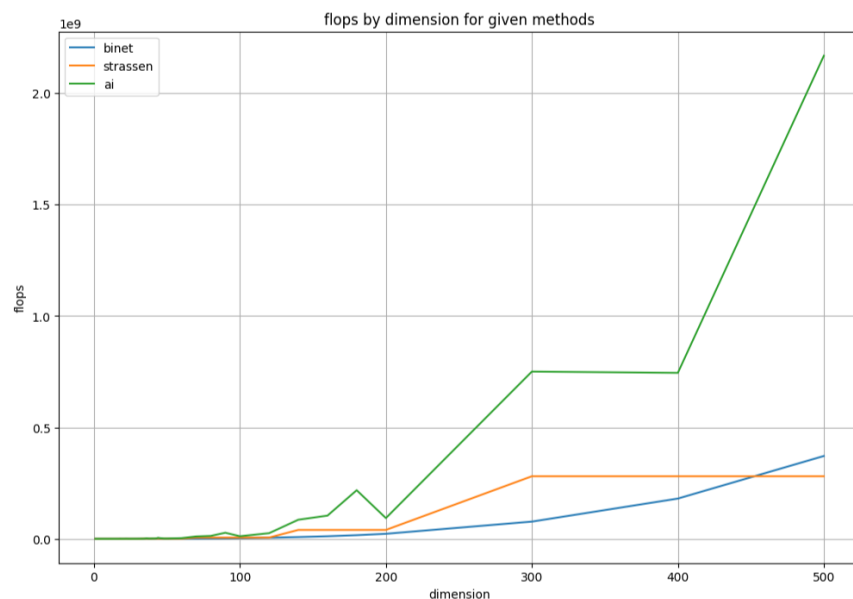
2 Porównanie algorytmów

Poniżej znajdują się 3 wykresy pokazujące działanie algorytmów w zależności od wymiarów macierzy. Są to wykresy pokazujące czas wykonania, liczbę operacji zmiennoprzecinkowych oraz zużyta pamięć dla każdego z 3 algorytmów. Algorytmy zostały uruchomione dla wszystkich wymiarów od 1 do 50, co dziesiątego wymiaru od 60 do 100, co dwudziestego wymiaru od 120 do 200 oraz wymiarów 300, 400 i 500.



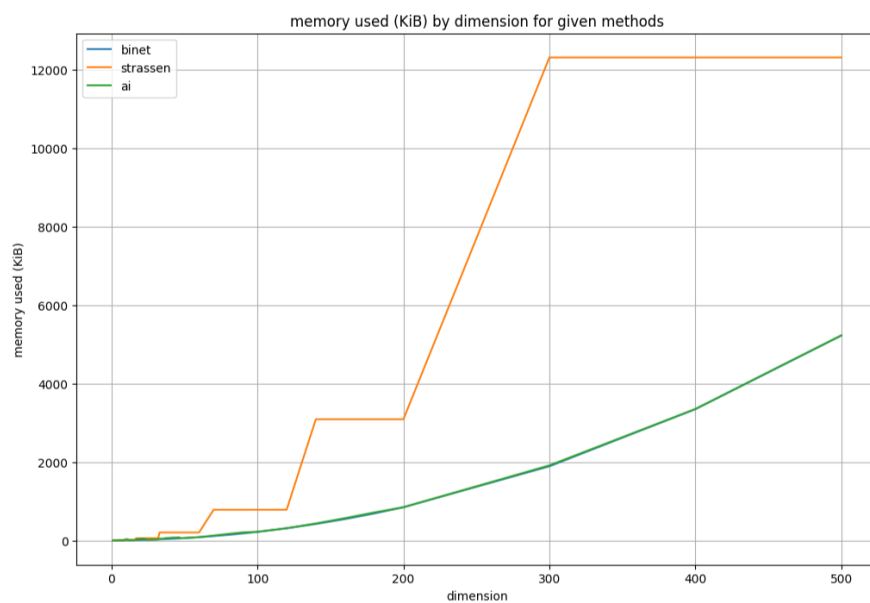
Rysunek 1: Czas wykonania algorytmów w zależności od rozmiaru macierzy

Czas wykonania dla największych macierzy okazał się być najlepszy dla algorytmu Strassena - jednakże widać, że on "skacze", ponieważ paduje macierze. Implementacja algorytmu Binét'a zaskakująco okazała się szybsza od algorytmu, który tak samo dzieli macierze, ale wykorzystuje lepsze faktoryzacje. Najprawdopodobniej jest to spowodowane nieoptymalnym mnożeniem rzadkiej macierzy przez wektor.



Rysunek 2: Liczba operacji zmiennoprzecinkowych algorytmów w zależności od rozmiaru macierzy

Liczba operacji zmiennoprzecinkowych dla algorytmu AI jest największa co nie jest zaskakujące patrząc na to, że był najwolniejszy. Ciekawe jest to, że algorytm Strassena dla wszystkich macierzy oprócz macierzy 500×500 wykonywał więcej operacji zmiennoprzecinkowych niż algorytm Binét'a, a dla tej największej różnica jest znacznie mniejsza od różnicy w czasie. Jest tak zapewne przez to, że algorytm Strassena dokonuje znacznie więcej dodawań niż mnożeń w porównaniu do algorytmu Binét'a, a dodawania są lżejsze dla sprzętu od mnożenia.



Rysunek 3: Zużycie pamięci algorytmów w zależności od rozmiaru macierzy

Dla algorytmów AI i Binét'a nie ma zauważalnej różnicy w użyciu pamięci, co nie jest zaskakujące. Oczekiwanym wynikiem również jest dużo większy przydział pamięci przez algorytm Strassena, gdyż dokonuje on paddingu macierzy.