

Wstęp do bioinformatyki

3 Dopasowanie lokalne

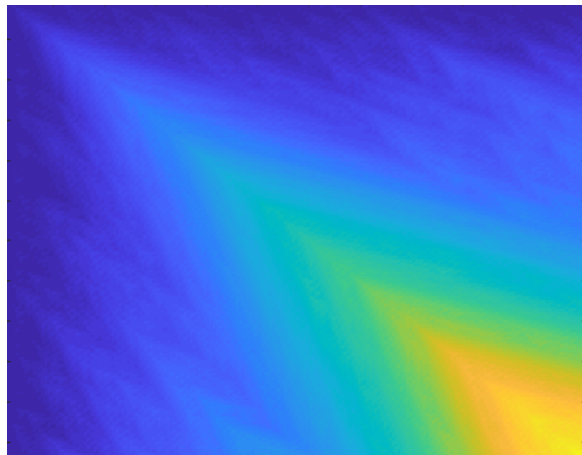
Link do repozytorium : <https://github.com/monikaRegula/Bioinformatics>

1. Porównanie par sekwencji ewolucyjnie powiązanych i niepowiązanych

Na rys. 2 porównanie sekwencji cytochromu b nosorożca i mastodona dla gap = -1.

```
#Seq1: GAAATTTTCGGCTCACTACTAGGAGCATG
#Seq2: GAAATTTTGGCTCTCTACTAGGAATCTG
#Length: 268
#Gap: -1
#Identity: 188/268 70%
#Gaps: 80/268 30%
GAAATTT-CGGCTC-CTACTAGGAGCAT--GCCT
||||| ||||| ||||||| || |||
GAAATTTT-GGCTCT-CTACTAGGA--ATCTGCCT
```

Rysunek 1 1 Plik końcowy z danymi
statystycznymi

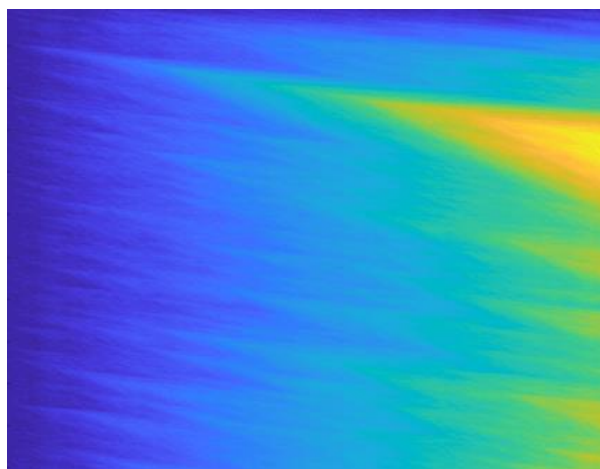


Rysunek 2 Porównanie sekwencji: Mastodon – Rhino gap = -1

Na rys. 4 porównanie sekwencji cytochromu b nosorożca i człowieka dla gap = -1.

```
1 #Seq1: ATGACCCCAATAAGCAAAATTAACCCCTAATAAAATTAA
2 #Seq2: GAAATTTTCGGCTCACTACTAGGAGCATGCCTAATTACCCA
3 #Length: 277
4 #Gap: -1
5 #Identity: 175/277 63%
6 #Gaps: 102/277 37%
7 GAAACTT-CGGCTCACT-CCT-TGG-CGC-CTGCCTGA-T--CCTCC
8 |||| || |||||||| || || || |||| || || ||
9 GAAA-TTTCGGCTCACTAC-TA-GGA-GCA-TGCCT-AATTACC-C-
10
11 #Seq1: ATGACCCCAATAAGCAAAATTAACCCCTAATAAAATTAA
12 #Seq2: GAAATTTTCGGCTCACTACTAGGAGCATGCCTAATTACCCA
13 #Length: 280
14 #Gap: -1
15 #Identity: 176/280 63%
16 #Gaps: 104/280 37%
17 GAAACTT-CGGCTCACT-CCT-TGG-CGC-CTGCCTGA-T--CCTCC
18 |||| || |||||||| || || || |||| || || ||
19 GAAA-TTTCGGCTCACTAC-TA-GGA-GCA-TGCCT-AATTACC-C-
20
21 #Seq1: ATGACCCCAATAAGCAAAATTAACCCCTAATAAAATTAA
22 #Seq2: GAAATTTTCGGCTCACTACTAGGAGCATGCCTAATTACCCA
23 #Length: 283
24 #Gap: -1
25 #Identity: 177/283 63%
26 #Gaps: 106/283 37%
27 GAAACTT-CGGCTCACT-CCT-TGG-CGC-CTGCCTGA-T--CCTCC
```

Rysunek 3 Plik końcowy z danymi statystycznymi



Rysunek 4 Porównanie sekwencji Human- Rhino gap = -1

2. Przykładowe działanie programu

Program korzysta z punktacji zgodności i niezgodności zdefiniowanej jako plik tekstowy punctuation.txt.

Dla wprowadzonych sekwencji wywoływane są funkcje odpowiedzialne między innymi za wyszukanie krótszych odcinków w obu sekwencjach, które są do siebie dobrze dopasowane i zapisuje je w pliku tekstowym matching.txt.

W celu podkreślenia najlepszych dopasowań lokalnych, komórki macierzy z punktacją ujemną są ustawiane na zero. Procedura śledzenia rozpoczyna się od komórki o najwyższym wyniku i trwa dopóki napotka komórkę z wynikiem równym zero.

```
#Seq1: GACTTAC
#Seq2: CGTGAATTCAT
#Length: 8
#Gap: -1
#Identity: 5/8 63%
#Gaps: 3/8 38%
GA-CTTAC
|| || |
GAA-TT-C
```

	-	C	G	T	G	A	A	T	T	C	A	T
-	0	0	0	0	0	0	0	0	0	0	0	0
G	0	0	2	1	2	1	0	0	0	0	0	0
A	0	0	1	0	1	4	3	2	1	0	2	1
C	0	2	1	0	0	3	2	1	0	3	2	1
T	0	1	0	3	2	2	1	4	3	2	1	4
T	0	0	0	2	1	1	0	3	6	5	4	3
A	0	0	0	1	0	3	3	2	5	4	7	6
C	0	2	1	0	0	2	2	1	4	7	6	5

Rysunek 6 Wynik działania programu – jedna z możliwych ścieżek dopasowania lokalnego

```
#Seq1: GACTTAC
#Seq2: CGTGAATTCAT
#Length: 8
#Gap: -1
#Identity: 5/8 63%
#Gaps: 3/8 38%
GA-CTT-A
|| || |
GAA-TTCA
```

	-	C	G	T	G	A	A	T	T	C	A	T
-	0	0	0	0	0	0	0	0	0	0	0	0
G	0	0	2	1	2	1	0	0	0	0	0	0
A	0	0	1	0	1	4	3	2	1	0	2	1
C	0	2	1	0	0	3	2	1	0	3	2	1
T	0	1	0	3	2	2	1	4	3	2	1	4
T	0	0	0	2	1	1	0	3	6	5	4	3
A	0	0	0	1	0	3	3	2	5	4	7	6
C	0	2	1	0	0	2	2	1	4	7	6	5

Rysunek 7 Wynik działania programu – kolejna z możliwych ścieżek dopasowania lokalnego

Rysunek 5 Wynik działania programu dla porównanie sekwencji wpisanej z klawiatury dla gap = -1

Przyjęta punktacja do obliczania dopasowania lokalnego zdefiniowana jest w pliku tekstowym punctuation.txt. Wartość gap jest definiowana przez użytkownika.

```
# A C G T
A 2 -7 -5 -7
C -7 2 -7 -5
G -5 -7 2 -7
T -7 -5 -7 2
```

Rysunek 3 Zawartość pliku z punktacją zgodności/niezgodności

3. Analiza złożoności obliczeniowej czasowej i pamięciowej

```
localMatching.m
1 function scoredMatrix = localMatching(seq1,seq2,gap,punctuationMatrix)
2 %Funkcja dla podanych sekwencji generuje macierz kosztów dopasowania
3 %lokalnego (algorytm Simtha- Watermana)
4 s1 = length(seq1);
5 s2 = length(seq2);
6
7 scoredMatrix = zeros(s1+1,s2+1);
8 scoredMatrix(1,2:end) = 0;
9 scoredMatrix(2:end,1) = 0;
10
11 for i = 2:s1+1 %iteracja po wierszach
12     for j = 2:s2+1 %iteracja po kolumnach
13         help = seq1(i-1);
14         help2 = seq2(j-1);
15         |
16         score = findPunctuation(punctuationMatrix,help,help2);
17
18         if(help == help2)
19             diagonal = scoredMatrix(i-1,j-1) + score;
20         else
21             diagonal = scoredMatrix(i-1,j-1) + score;
22         end
23         %POZIOM
24         left = scoredMatrix(i-1,j) + gap;
25         %PION
26         up = scoredMatrix(i,j-1) + gap;
27
28         %wybranie maksimum z 4 opcji score: diagonal,left,up,zero
29         %maksimum to odległość edycyjna pomiędzy seq1 a seq2
30         maxScore = max([diagonal left up 0]);
31         scoredMatrix(i,j)= maxScore;
32
33     end
34 end
35
36 end
37
38
```

Pamięciowa

$1+1+(s1+1) * (s2+1) + 1+1+1+1+1+1+1+1+1+1 = 11 + nm$
gdzie $s1+1 = m$, $s2+1 = n$

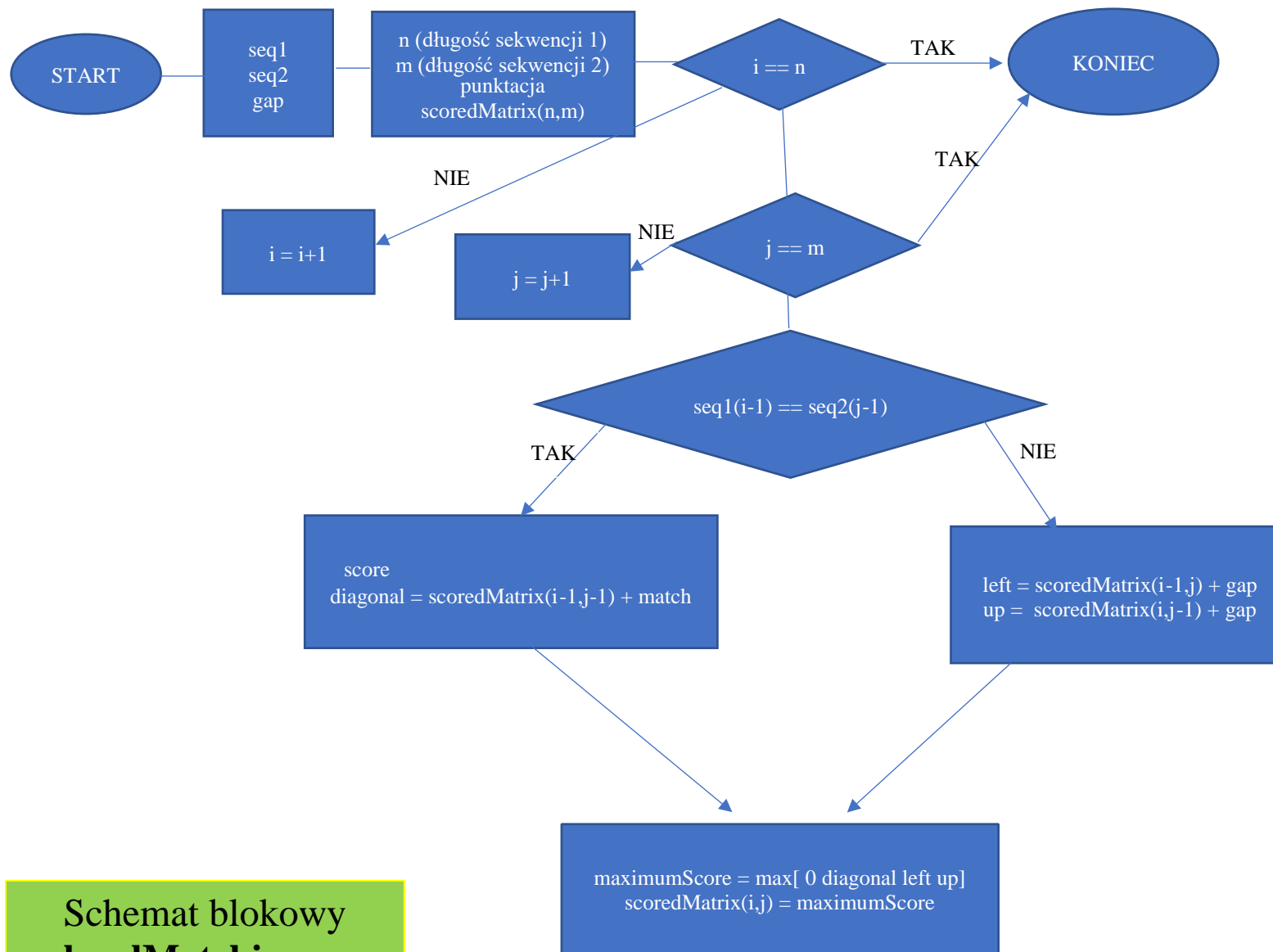
Czasowa

Pętla wykonuje się $(m-1)(n-1)$ razy z powtarzaniem czynności tj:

- Przypisanie licznika pętli i,j
- Przypisanie zmiennych pomocniczych help, help2
- Przypisanie zmiennej score
- Sprawdzenie warunku
- Przypisanie zmiennych diagonal, left, up, maxScore
- Przypisanie wartości do macierzy scoredMatrix
- Inkrementacja licznika pętli

$$(m-1)*(n-1)(1+1+1+1+1+1+1+1+1+1+1+1) = 12(m-1)(n-1)$$

$$O(m) = (m-1)(n-1) \sim O(m^2)$$



**Schemat blokowy
localMatching.m**

```

1 function [localPaths,optimalPaths] = traceback(scoredMatrix,seq1,seq2,gap,punctuationMatrix)
2 %Funkcja generuje dla znalezionych wartości maksymalnych w macierzy kosztów
3 %scoredMatrix ścieżki optymalne dopasowania lokalnego
4
5 %wymiary macierzy:
6 m = size(scoredMatrix,1);
7 n = size(scoredMatrix,2);
8 %szukam maximum w macierzy punktacji
9 maximumScore = max(scoredMatrix(:));
10 %prealokacja tabeli dla optymalnej ścieżki
11 optimalPath = zeros(m,n);
12 %komórek zawierających wartość równą maximumScore może być więcej, dlatego
13 %zapamiętuję lokalizację wszystkich możliwych
14 allMaxes = (scoredMatrix(:) == maximumScore);
15 %w miejscu maksima wstawiam 1
16 optimalPath(allMaxes) = 1;
17 %x zwraca numery komórek z maksimami
18 x = find(optimalPath == 1);
19 %memorise to pewnego rodzaju nawigacja; tu gromadzone są współrzędne, gdzie:
20 %- Columns - zapamiętuje nr kolumn, w których znajdują się maximumScore
21 %- Rows - zapamiętuje nr wierszy, w których znajdują się maximumScore
22 [Rows,Columns] = find(optimalPath == 1);
23 howManyMaxes = length(x);
24 memorise = [Rows,Columns];
25
26 %pętla jest wykonywana tyle razy ile jest możliwych maksimów w macierzy
27 %punktacji (scoredMatrix)
28 for i = 1:length(x)
29     %currentCell pobiera aktualny nr komórki maksima
30     currentCell = x(i);
31     path = zeros(m,n);
32     %steps to mapa kroków
33     steps = [];
34     alignment1 = '';
35     alignment2 = '';
36     aligner = '';
37     identity = 0;
38     gaps = 0;
39
40     row = Rows(i);
41     column = Columns(i);
42
43     while scoredMatrix(currentCell)>0
44         navigator = scoredMatrix(currentCell);
45         score = findPunctuation(punctuationMatrix,seq1(row -1),seq2(column -1));

```

3. Analiza złożoności obliczeniowej czasowej i pamięciowej

Pamięciowa

$$1+1+1+m*n + 1+1+1+1+1+1+1+m*n+1+1+1+1+1+1+1+1+1 = 19 + 2mn$$

Czasowa

Ilość wykonywanej pętli uzależniona jest od ilości występujących maximów w macierzy scoredMatrix, ta ilość to (length(x)) oraz od długości obliczonych optymalnych ścieżek.

Pętla while wykonywana jest dopóki element macierzy nie jest równa 0

Powtarzające się czynności:

- Przypisanie zmiennych: licznik i, currentCell, path, steps, row, column
- Sprawdzenie warunku pętli while
- Przypisanie zmiennych: nawigator,score
- Sprawdzenie warunków linie: 46,56,66,69
- Przypisanie zmiennej path, alignment1, alignment2, aligner
- Zmiana wartości zmiennych currentCell, row, column, steps

Założenie: length(x) = k, ilość wykonania pętli while dla jednego maximum = s
Ilość sprawdzanych warunków maksymalnie 4

$$k*(1+1+1+1+1) * s(1+1+1+1+3+4) = 5k*11s$$

$$O(n) = 55ks \sim O(n^2)$$

```

46 -         if scoredMatrix(currentCell - 1) == navigator - gap
47 -             alignment1 = strcat(alignment1,seq1(row-1));
48 -             alignment2 = strcat(alignment2,'-');
49 -             aligner = strcat(aligner," ");
50 -             path(currentCell) = 1;
51 -             currentCell = currentCell - 1;
52 -             row = row -1;
53 -             steps = [steps,3];
54 -             gaps = gaps + 1;
55 -
56 -         elseif scoredMatrix(currentCell - m) == navigator - gap
57 -             alignment1 = strcat(alignment1,'-');
58 -             alignment2 = strcat(alignment2,seq2(column-1));
59 -             aligner = strcat(aligner," ");
60 -             path(currentCell) = 1;
61 -             currentCell = currentCell - m;
62 -             column = column -1;
63 -             steps = [steps,1];
64 -             gaps = gaps+1;
65 -
66 -         elseif scoredMatrix(currentCell - m -1) == navigator - score
67 -             alignment1 = strcat(alignment1,seq1(row-1));
68 -             alignment2 = strcat(alignment2,seq2(column-1));
69 -             if(seq2(column-1)== seq1(row-1))
70 -                 aligner = strcat(aligner,'|');
71 -                 identity = identity + 1;
72 -             else
73 -                 aligner = strcat(aligner,' ');
74 -             end
75 -             path(currentCell ) = 1;
76 -             currentCell = currentCell - m - 1;
77 -             row = row -1;
78 -             column = column -1;
79 -             steps = [steps,2];
80 -             identity = identity + 1;
81 -         end
82 -     end

```

Rysunek 4 Kod źródłowy pliku traceback.m

