

PROJEKT I

09.01.2019

TERMIN ROZLICZENIA: 21.01.2019

WYSZUKIWANIE OGRANICZONE MOTYWU W ŁAŃCUCHACH DNA

```
'tagtggctctttgagtgtagatctgaagggaaagtattccaccagttcggggtcaccagcagggcagggcgacttaat',
'cgcgactcggcgctcacagttatcgacgcttagacaaaaacggagttggatccgaaactggagtttaacggagtcctt',
'gttacttgtgagcctgggttagaccgaaatataattgttggctgcatagcggagctgacatacgagtaggggaaatcgct',
'aacatcaggctttgattaacaatttaagcacgtaaatccgaattgacctgatgacaatacggaacatgccggctccggg',
'accaccggataggctgcttattaggccaaaaggtagatcgaataatggctcagccatgtcaatgtcggcattccac',
'tagattcgaatcgatcgtgtttctcctctgtgggttaacgaggggtccgaccttgctcgcgatgtgccgaacttgatcc',
'gaaatggttcggtcgatatcaggccgttcttaacttggcgggtgcagatccgaactctctggaggggtcgtcgcgta',
'atgtatactagacatttaacgctcgtctattggcggagaccatttgctccactacaagaggctactgtgtagatccgta',
'ttcttacacccttcttagatccaaactgttggcggcattctcttttcgagtccttgtaacctccatttgctctgatgac',
'ctacctatgtaaaacaacatctactaactgtagtcggctcttctctgatctgccctaacctacaggtcgatccgaaattcg'
```

Dane jest 10 sekwencji DNA. Przygotować implementację wyznaczającą łańcuch konsensusu długości 7 dla podanego zestawu DNA, wykorzystując algorytmy z opisanej biblioteki operujące na strukturze drzewa danych.

- 1) Opracować własne testy do wszystkich stosowanych funkcji bibliotecznych
- 2) Zaimplementować :
 - 1) przeszukiwanie wyczerpujące
 - 2) przeszukiwanie z ograniczeniami

• Given $s = (s_1, \dots, s_t)$ and DNA:

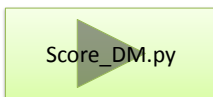
	a	g	g	t	a	c	t
C	c	a	t	a	c	g	t
a	c	g	t	a	a	g	t
a	c	g	t	c	a	a	t
C	c	g	t	a	c	g	g

$$Score(s, DNA) = \sum_{i=1}^l \left(\max_{k \in \{A, C, G, T\}} count(k, i) \right)$$

A	3	0	1	0	3	1	1	0
C	2	4	0	0	1	4	0	0
G	0	1	4	0	0	0	3	1
T	0	0	0	5	1	0	1	4

Consensus a c g t a c g t

Score 3+4+4+5+3+4+3+4=30



- **Goal:** Given a set of DNA sequences, find a set of l -mers, one from each sequence, that maximizes the consensus score
- **Input:** A $t \times n$ matrix of DNA, and l the length of the pattern to find
- **Output:** An array of t starting positions $s = (s_1, s_2, \dots, s_t)$ maximizing $Score(s, DNA)$

łańcuch konsensusu

Lecture 5:
Finding Regulatory Motifs
Within DNA Sequences

COMP 555 Bioalgorithms (Fall 2014)

- DNA - array of sequence fragments
- t - number of sample DNA sequences
- n - length of each DNA sequence

- l - length of the motif (l -mer)
- s_i - starting position of an l -mer in sequence i
- $s = (s_1, s_2, \dots, s_t)$ - array of motif's starting positions

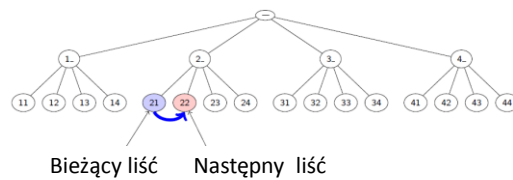
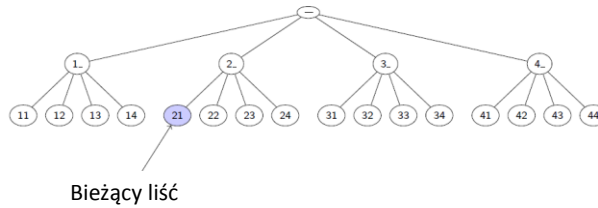
Przejdź do następnego liścia

NextLeaf(a, L, k)

```

1: for  $i \leftarrow L$  to 1 do
2:   if  $a_i < k$  then
3:      $a_i \leftarrow a_i + 1$ 
4:   return  $a$ 
5:    $a_i \leftarrow 1$ 
6: return  $a$ 

```



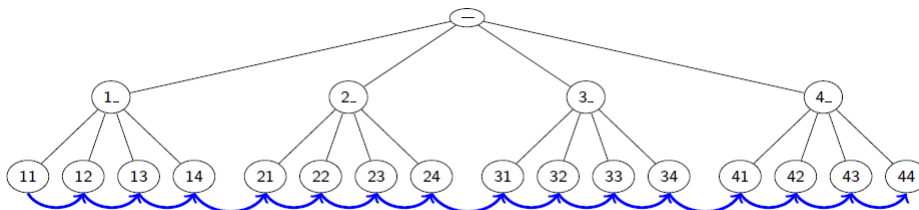
Odwiedź wszystkie liście w porządku rosnącym

AllLeaves(L, k)

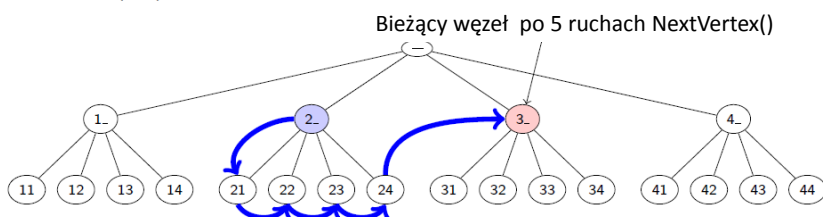
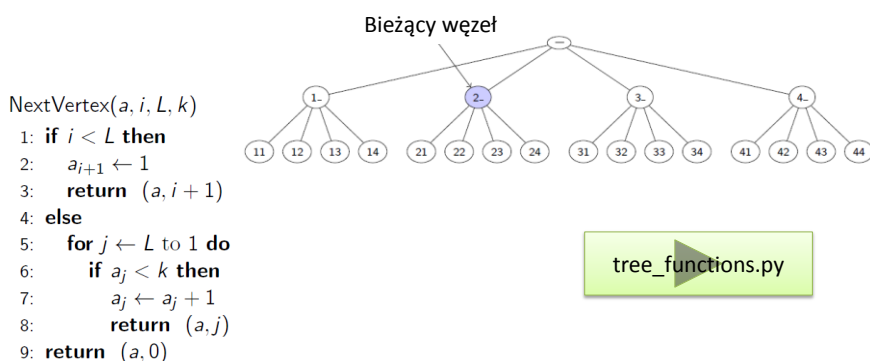
```

1:  $a \leftarrow (1, \dots, 1)$ 
2: while forever do
3:   output  $a$ 
4:    $a \leftarrow \text{NextLeaf}(a, L, k)$ 
5:   if  $a = (1, 1, \dots, 1)$  then
6:     return

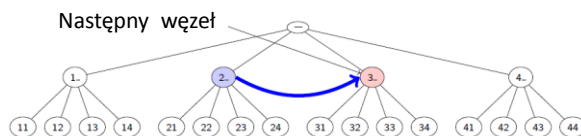
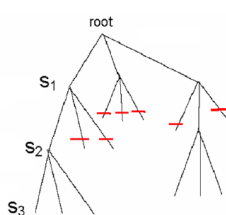
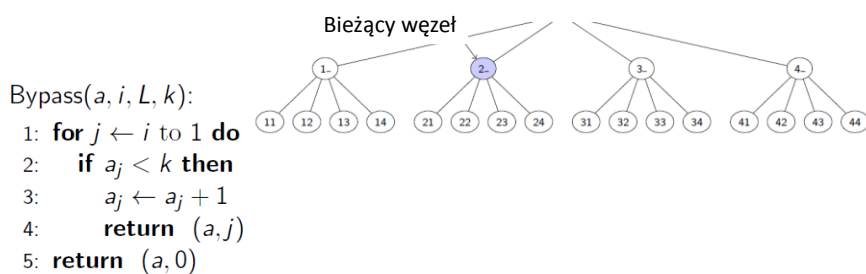
```



Przejdź do następnego węzła w drzewie



Przejdźcie na skróty: dla danego węzła wewnętrznego wskaź węzeł następny taki, który pomija wszystkie dzieci danego węzła.





BruteForceMotifSearchAgain(DNA, t, n, ℓ)

```

1:  $\mathbf{s} \leftarrow (1, 1, \dots, 1)$ 
2:  $bestScore \leftarrow score(\mathbf{s}, DNA)$ 
3: while forever do
4:    $\mathbf{s} \leftarrow \text{NextLeaf}(\mathbf{s}, t, n - \ell + 1)$ 
5:   if  $Score(\mathbf{s}, DNA) > bestScore$  then
6:      $bestScore \leftarrow score(\mathbf{s}, DNA)$ 
7:      $bestMotif \leftarrow (s_1, s_2, \dots, s_t)$ 
8:   if  $\mathbf{s} = (1, 1, \dots, 1)$  then
9:     return  $bestMotif$ 

```

Poprawiamy score() na „optimistyczne” score()

1. Score() zlicza jedynie do pozycji i w sekwencji l -meru

$$Score(s, i, DNA) = \sum_{j=1}^{\ell} \max_{c \in \{A, C, G, T\}} count_i(c, j)$$

Score() dla węzłów wewnętrznych:

pozycje nieobsadzone zastępujemy maksymalnie możliwym zliczeniem co daje:

$$optimisticScore \leftarrow Score(s, i, DNA) + (t - i) \cdot \ell$$

```

BranchAndBoundMotifSearch(DNA, t, n,  $\ell$ )
1:  $s \leftarrow (1, 1, \dots, 1)$ 
2:  $bestScore \leftarrow 0$ 
3:  $i \leftarrow 1$ 
4: while  $i > 0$  do
5:   if  $i < t$  then
6:      $optimisticScore \leftarrow Score(s, i, DNA) + (t - i) \cdot \ell$ 
7:     if  $optimisticScore < bestScore$  then
8:        $(s, i) \leftarrow Bypass(s, i, n - \ell + 1)$ 
9:     else
10:       $(s, i) \leftarrow NextVertex(s, i, t, n - \ell + 1)$ 
11:   else
12:     if  $Score(s, DNA) > bestScore$  then
13:        $bestScore \leftarrow score(s, DNA)$ 
14:        $bestMotif \leftarrow (s_1, s_2, \dots, s_t)$ 
15:      $(s, i) \leftarrow NextVertex(s, i, t, n - \ell + 1)$ 
16: return  $bestMotif$ 

```