

Homework 4

Elliott Pryor

14 November 2021

Problem 1 Consider the following set of sequences:

$$\begin{aligned} S_1 &= ACTCTCGATC \\ S_2 &= ACTTCGATC \\ S_3 &= ACTCTCTATC \\ S_4 &= ACTCTCTAATC \end{aligned}$$

Compute the MSA using the center star method.

We give the pairwise distances

	S_1	S_2	S_3	S_4
S_1	0	1	1	2
S_2	1	0	2	3
S_3	1	2	0	1
S_4	2	3	1	0

We get that S_1 or S_3 would be equally good center stars. We pick S_1 .

We compute alignments with S_1 :

From this we get the MSA:

$$\begin{array}{ll} S_1 = ACTCTCGATC & ACTCTCG_ATC \\ & = ACTCTCGATC \\ S_2 = ACTCTCGATC & ACT_TCG_ATC \\ & = ACT_TCGATC \\ S_3 = ACTCTCGATC & ACTCTCT_ATC \\ & = ACTCTCTATC \\ S_4 = ACTCTCG_ATC & ACTCTCTAATC \\ & = ACTCTCTAATC \end{array}$$

Problem 2 Compute the ClustalW alignment of the previous problem:

	S_1	S_2	S_3	S_4
S_1	0	0	0.1	0.1
S_2	0	0	0.1	0.1
S_3	0	0	0	0
S_4	0	0	0	0

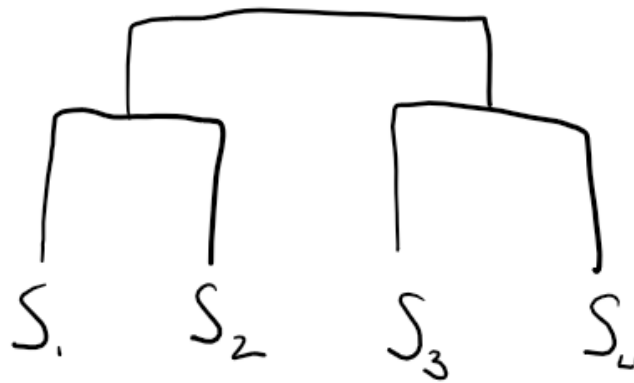


Figure 1: Guide tree

Aligning S_1, S_2 produces $(ACTCTCGATC, ACT_TCGATC)$,

aligning S_3, S_4 produces $(ACTCTCT_ATC, ACTCTCTAATC)$.

Then using the profile-profile alignment we get:

$ACTCTCGATC_$
 $ACT_TCGATC_$
 $ACTCTCTATC_$
 $ACTCTCTAATC$

Problem 3 Given a set of k sequences: S_1, \dots, S_k find k substrings such that the optimal SP score of multiple sequence alignment is maximized. Give a Dynamic Programming algorithm to solve this. (Note, when $k = 2$ this is same as solving local alignment problem)

We build a generalized version of Smith-Waterman algorithm. We base this on the generalized version of Needleman-Wunch algorithm from class:

$$V(i_1, \dots, i_k) = \max_{b_1, \dots, b_k \in \{(0,1)\}^k \setminus (0,0)} V(i_1 - b_1, \dots, i_k - b_k) + SP(S_1[i_1 b_1], \dots S_k[i_k b_k])$$

We note that the only difference in the Smith-Waterman algorithm from the Needleman-Wunch algorithm is the option to take a score of 0 always. So we have a generalized Smith-Waterman:

$$V(i_1, \dots, i_k) = \max_{b_1, \dots, b_k \in \{(0,1)\}^k \setminus (0,0)} (V(i_1 - b_1, \dots, i_k - b_k) + SP(S_1[i_1 b_1], \dots S_k[i_k b_k]), 0)$$

With a base case of $V(0, \dots, 0) = 0$ and $i_j = 0 \implies b_j = 0$.

Then we have to find the local alignment with the best score. This is the entry with the greatest value in the matrix. So the optimal score is $\max_{i_1 \dots i_k} V(i_1 \dots i_k)$.

Building the matrix V takes $O(n^k)$ space. For running time we have $O(n^k)$ entries to fill out, $2^k - 1$ options for b , and $\binom{k}{2}$ evaluations of SP . Thus we have $O(n^k 2^k k^2)$ time (same as generalized Needleman-Wunch) to build V . Then we need to find the max in V : which is just searching through all n^k entries. Thus our running time is: $O(n^k 2^k k^2)$

Problem 4 Code Center Star Algorithm.

```
PS D:\pryor\Documents\GitHubProjects\CompBio\HW4> python .\Center_Star.py 1 1 example1.fasta
Opening: example1.fasta
Center (string, index): ('bananan', 0)
MSA:
banana_n_
ba__n__n_
aannna_na
_ananan_a
banana_n_
```

Figure 2: Example showing execution on example1

```
PS D:\pryor\Documents\GitHubProjects\CompBio\HW4> python .\Center_Star.py 1 1 example2.fasta
Opening: example2.fasta
Center (string, index): ('ACTCTCGATC', 0)
MSA:
ACTCTCG_ATC
ACT_TCG_ATC
ACTCTCT_ATC
ACTCTCTAATC
```

Figure 3: Example showing execution on strings from Problem 1

```
> EXAMPLE 1 seq 1
bananan
> seq 2
bann
> seq 3
aannnana
> seq 4
ananana
> seq 5
bananan
```

```
> EXAMPLE 2 seq1
ACTCTCGATC
>seq2
ACTTCGATC
>seq3
ACTCTCTATC
>seq4
ACTCTCTAATC
```

```
1 from EditDistance import *
2 import numpy as np
3
4
5 def compute_indx(Sc,s2):
6     indx = []
7     for i in range(len(s2)):
8         if i >= len(Sc) or (s2[i] == '_' and Sc[i] != '_'):
9             indx.append(i)
10            Sc = Sc[:i] + '_' + Sc[i:]
```

```

11     indx += [len(Sc) + i for i in range(len(s2) - len(Sc))]
12     return indx
13
14
15 def center_star(strings, f=delta):
16     """
17     :param strings: List of strings to compute
18     :return: center string, index of center string in strings array (param)
19     """
20
21     mat = np.zeros((len(strings), len(strings))) # make k x k matrix
22     for i, S in enumerate(strings):
23         for j, T in enumerate(strings):
24             if j <= i: # skip the lower diagonal elements
25                 continue
26                 D = needleman(S, T)[0]
27                 mat[i][j] = D[len(S)][len(T)] # get the edit distance
28                 mat[j][i] = mat[i][j]
29
30     # the string that minimizes sum of distances to all other strings
31     center = np.argmin(mat.sum(axis=1))
32     return strings[center], center
33
34
35 def MSA(strings, f=delta):
36     Sc, center = center_star(strings)
37     alignment = [s for s in strings] # copy the input strings
38     for i in range(len(strings)):
39         if i == center: # skip center string
40             continue
41
42         V, P = needleman(Sc, strings[i], f)
43         s2, t2 = construct_alignment(P, Sc, strings[i])
44         alignment[i] = t2
45
46         add_spaces = compute_indx(Sc, s2)
47         for indx in add_spaces:
48             for j in range(i):
49                 alignment[j] = alignment[j][:indx] + '_' + alignment[j][indx:]
50     Sc = s2
51     alignment[center] = Sc
52
53     return alignment
54
55
56 def read_fasta(filepath):
57     """
58     Returns a dict key=sequence_name, value=sequence from the fasta file
59     :param filepath:
60     :return:
61     """
62     file = open(filepath)
63     print(f"Opening: {filepath}")
64     started = False
65     header = ''

```

```

66     text = ''
67     sequences = {}
68     for line in file: # reads in fasta file. Allows for sequence to continue on
multiple lines
69         if line.startswith('>'):
70             if started:
71                 sequences[header] = text.replace('\n', '')
72                 text = ''
73                 started = True
74                 header = line.replace('>', '').strip()
75             else:
76                 text += line.strip()
77     sequences[header] = text.replace('\n', '')
78     return sequences
79
80
81 def diff_generator(a, b):
82     def f(c1, c2):
83         # Helper function used for scoring
84         if c1 == c2:
85             return 0
86         elif c1 == '_' or c2 == '_':
87             return b
88         else:
89             return a
90     return f
91
92
93 if __name__ == "__main__":
94     if len(sys.argv) == 4:
95         alpha = float(sys.argv[1])
96         beta = float(sys.argv[2])
97         file = sys.argv[3]
98         seq = read_fasta(file)
99         f = diff_generator(alpha, beta)
100        ins = list(seq.values())
101
102        print(f'Center (string, index): {center_star(ins, f)}')
103        print("MSA:")
104        msa = MSA(ins, f)
105        for s in msa:
106            print(s)
107    else:
108        print(center_star(['ACTCTCGATC', 'ACTTCGATC', 'ACTCTCTATC', 'ACTCTCTAATC']))
109        print(center_star(['bananan', 'bann', 'aannnana', 'ananana', 'bananan']))
110        msa = MSA(['bananan', 'bann', 'aannnana', 'ananana', 'bananan'])
111        for s in msa:
112            print(s)

```
