

Homework 1

Elliott Pryor

7 September 2021

Problem 1

For the following mRNA sequence, can you extract its 5' UTR, 3' UTR and the protein sequence? ACTTGTCATGGTAACTCCGTCGTACCAGTAGGTCATG

So the above is actually DNA, since it has T's. So we assume that the T's are actually U's.

So in RNA we have: ACUUGUCAUGGUAACUCCGUCGUACCAGUAGGUCAUG

We then look for the start codon (Met) which is AUG. So we get:

ACUUGUC AUG GUA ACU CCG UCG UAC CAG UAG GUC AUG

The 5' UTR is everything left of the first AUG which is: ACUUGUC. The stop codons are either UAA, UAG, or UGA. We see UAG, so the 3' UTR is everything right of this: GUCAUG

Thus our coding region is:

AUG GUA ACU CCG UCG UAC CAG UAG → MET VAL THR PRO SER TYR GLN Stop

Problem 2 Implement the Z-algorithm in the language of your choice. Also, code up the Exact Pattern Matching algorithm that makes use of the Z algorithm. Demonstrate both algorithms on some test strings (use screen shots to show runs).

Below are some screenshots showing functionality. They were run in a python shell and imported the python file that I made for concision and clarity.

```
>>> import HW1.z_algorithm as z
>>> z.z_alg('aabbabaaa')
[0, 1, 0, 0, 1, 0, 2, 2, 1]
>>> z.z_alg('aaaaaaaa')
[0, 7, 6, 5, 4, 3, 2, 1]
>>> z.z_alg('ababbba')
[0, 0, 2, 0, 0, 0, 1]
>>> z.z_alg('aaaacaaaaa')
[0, 4, 3, 2, 1, 0, 5, 5, 4, 3, 2, 1]
```

Figure 1: Examples using base z algorithm

```
>>> import HW1.z_algorithm as z
>>> z.z_pattern_match("aaa", "abaaabaabaaa")
[2, 9]
>>> z.z_pattern_match('aaa', 'aaaaaaaaaaaa')
[0, 1, 2, 3, 4, 5, 6, 7, 8, 9]
>>> z.z_pattern_match('ab', 'hello world')
[]
```

Figure 2: Examples using pattern matching

I also did the bonus, which runs via the command line. Here is an example below showing functionality.

My example FASTA input:

```
>seq0, example file from
    http://prodata.swmed.edu/promals/info/fasta_format_file_example.htm
    with last sequence modified
FQTTWEEFSRAAEKLYLADPMKVRVVLKYRHVDGNLCIKVTDDLVLVYRTDQAQDVKKIEKF
>seq1
KYRTWEEFSTRAAEKLYQADPMKVRVVLKYRHCDGNLCIKVTDDVVCLLYRTDQAQDVKKIEKFHSQLMRLME
    LKVTDNKECLKFKTDQAQEAKKMEKLNNIFFTLM
>seq2
```

```

    EEYQTWEEFARAAEKLYLTDPMKVRVVLKYRHCDGNLCMKVTD DAVCLQYKTDQAQDVKKVEKLHGK
>seq3
MYQVWEEFSRAVEKLYLTDPMKVRVVLKYRHCDGNLCIKVTDNSVCLQYKTDQAQDVK
>seq4
EEFSRAVEKLYLTDPMKVRVVLKYRHCDGNLCIKVTDNSVVS YEMRLFGVQKDNFALEHSL
>seq5
SWEEFAKAAEVLYLEDPMKCRMCTKYRHVDHKL VVKLTDNHTVLKYVTDMAQDVKKIEKLTLLMR
>seq6
FTNWEEFAKAAERLHSANPEKCRFVTKYNHTKGELVLKLTDDVVCLQYSTNQLQDVKKLEKLSSTLLRSI
>seq7
SWEEFVERSVQLFRGDPNATRYVMKYRHCEGKLV LKVTD DRECLKFKTDQAQDAKKMEKLNNIFF
>seq8
SWDEFVDRSVQLFRADPESTRYVMKYRHCDGKLVLKVTDNKECLKFKTDQAQEAKKMEKLNNIFFTLM
>seq9
KNWEDFEIAAENMYMANPQNCRYTMKYVHSGHILLKMSDNVKCVQYRAENMPDLKK
>seq10
THISISABSOLUTGIBBRISHJUSTTOSHOWITSOUTPUTBUTNOMATCH

```

```

PS C:\Users\pryor\Documents\GithubProjects\CompBio\HW1> python .\z_algorithm.py E example.fasta
Opening: example.fasta
Matching pattern: E
Analyzing Sequence: seq0, example file from http://prodata.swmed.edu/promals/info/fasta\_format\_file\_example.htm
Found matches 4, 5, 11, 59
Analyzing Sequence: seq1
Found matches 5, 6, 12, 60, 71, 80, 91, 96
Analyzing Sequence: seq2
Found matches 0, 1, 6, 7, 13, 61
Analyzing Sequence: seq3
Found matches 5, 6, 12
Analyzing Sequence: seq4
Found matches 0, 1, 7, 43, 57
Analyzing Sequence: seq5
Found matches 2, 3, 9, 14, 57
Analyzing Sequence: seq6
Found matches 4, 5, 11, 19, 33, 59
Analyzing Sequence: seq7
Found matches 2, 3, 6, 29, 41, 57
Analyzing Sequence: seq8
Found matches 3, 17, 41, 52, 57
Analyzing Sequence: seq9
Found matches 3, 6, 10, 49
Analyzing Sequence: seq10
No match found

```

Figure 3: Program output on example Fasta input

```
1 """
2 Author - Elliott Pryor
3 Created on - 2 September 2021
4 """
5 import sys
6
7
8 def print_output(z):
9     if len(z) > 0:
10         z = [str(i) for i in z]
11         output = ', '.join(z)
12         print(f"Found matches {output}")
13     else:
14         print("No match found")
15
16 def find_next_match(s1, s2, S):
17     """
18     Finds the next match in S starting from s1 and s2
19     Compares S[s1, n] to S[s2, n]
20     We don't pass in slices because this is faster
21     :param s1: Starting of first string (base)
22     :param s2: Start of second string (prefix to match)
23     :param S: String to search
24     :return: q=s2 index of failed match in S, zk length of match
25     """
26     n = len(S)
27     zk = 0
28     while (s1 < n and s2 < n) and (S[s1] == S[s2]): # while they match and in bounds
29         s1 += 1
30         s2 += 1
31         zk += 1
32     return s2, zk
33
34
35 def z_alg(S):
36     """
37     Run the z prefix matching algorithm
38     :param S: String input
39     :return: Z score for each index (note first one is always 0)
40     """
41     l, r = 0, 0 # initialize
42     z = [0 for _ in S] # init z scores
43     for k in range(1, len(S)):
44         if k > r: # match them
45             _, zk = find_next_match(0, k, S)
46             if zk > 0:
47                 r = k + zk - 1
48                 l = k
49                 z[k] = zk
50         else:
51             kp = k - l
52             b = r - k + 1 # ||B|| from alg descrip
53             if z[kp] < b:
54                 z[k] = z[kp]
55             else:
```

```

56         q, zk = find_next_match(b, r + 1, S)
57         z[k] = q - k
58         l = k
59         r = q - 1
60     return z
61
62
63 def z_pattern_match(pattern, text, special_char=None):
64     """
65     Finds all occurrences of the pattern in the text
66     :param pattern: pattern to search for
67     :param text: text to search in
68     :param special_char: special division characters. Can specify a char, or list of
69                          chars.
70                          Defaults are $, %, ^, ~, \, #
71     :return: Index in text of start of pattern match
72     """
73     if special_char is None:
74         special_char = ['$ ', '% ', '^ ', '~ ', '\\ ', '# '] # our default options for
75         special chars
76
77     sep_char = ' '
78     if not isinstance(special_char, str): # if it isn't just a string
79         for c in special_char:
80             if c not in text and c not in pattern: # 2 * O(n) each time
81                 sep_char = c
82                 break
83
84     S = pattern + sep_char + text # put special char in the middle
85     z = z_alg(S)
86     m = len(pattern)
87     matches = []
88     for j in range(1, len(z)):
89         if z[j] == m:
90             matches.append(j - (m + 1))
91     return matches
92
93
94 def analyze_fasta(pattern, filepath):
95     file = open(filepath)
96     print(f"Opening: {filepath}")
97     print(f"Matching pattern: {pattern}")
98     started = False
99     text = ''
100     for line in file:
101         if line.startswith('>'):
102             if started:
103                 z = z_pattern_match(pattern, text.replace('\n', ' '))
104                 text = ''
105                 print_output(z)
106                 started = True
107                 print(f"Analyzing Sequence: {line.replace('>', ' ').strip()}")
108             else:
109                 text += line.strip()
110     z = z_pattern_match(pattern, text.replace('\n', ' '))
111     print_output(z)

```

```
109
110
111 if __name__ == '__main__':
112     if len(sys.argv) == 3:
113         # run fasta file
114         # expect arguments z_algorithm.py <pattern> <path to file>
115         pattern = sys.argv[1]
116         filepath = sys.argv[2]
117         analyze_fasta(pattern, filepath)
118     else:
119         print(z_alg('aabbabaaa'))
```