

Biologia Computacional 2019/2 - Lista de Exercícios 2 parte 2

Aluno: Leonardo Oliveira Wellausen

1a)

```
1 import numpy as np
2
3 word1 = list("MFVLKGSVVQAFVLLSIVCLEITTIADGVRVYVNAEWKRPEQSOEGRHSCTARRLEDNSEEVACSTEVKFRORAPAEYANKIKKAKDKLRRLESQFDDCCQOENDRKDRLLIQLOQANL")
4 word2 = list("MEKVPGEIMEIERERSEELSEAERKAVQAMWRLYANCEDVGVAIVLVRFFVNFPSAKOYFSQFKHMEOPLEMERSPOLRKHACRVNGALNTWENLHDPDKVSSVLALVGKAHA")
5
6 #word1 = 'aaaaabb'
7 #word2 = 'bbccccc'
8
9 size1 = len(word1)
10 size2 = len(word2)
11
12 gap = -2
13 match = 1
14 mismatch = -1
15
16 matrix = np.zeros((size1 + 1, size2 + 1))
17
18 def matching(a, b):
19     if a == b:
20         return match
21     else:
22         return mismatch
23
24 maxval = -1000
25
26 for i in range(1, size1 + 1):
27     for j in range(1, size2 + 1):
28         val = max(matching(word2[j-1], word1[i-1]) + matrix[i-1][j-1], gap + matrix[i-1][j], gap + matrix[i][j-1], 0)
29         if val >= maxval:
30             maxval = val
31         matrix[i][j] = val
32
```

Como nos exercícios da lista anterior, as sequências a serem alinhadas já estão inseridas no código, assim como os pesos do algoritmo. O algoritmo é implementado em um laço percorrendo e preenchendo a matriz.

```

34 def back_trace(starti, startj):
35     res1 = ''
36     res2 = ''
37     i = starti
38     j = startj
39     matches = 0
40     while matrix[i][j] != 0:
41         up = matrix[i-1][j]
42         left = matrix[i][j-1]
43         diag = matrix[i-1][j-1]
44
45         if diag >= up and diag >= left:
46             res1 += word1[i-1]
47             res2 += word2[j-1]
48             i -= 1
49             j -= 1
50             if word1[i] == word2[j]:
51                 matches += 1
52         elif up >= diag and up >= left:
53             res1 += word1[i-1]
54             res2 += '-'
55             i -= 1
56         else:
57             res1 += '-'
58             res2 += word2[j-1]
59             j -= 1
60     res1 = res1[::-1]
61     res2 = res2[::-1]
62     print('Cadeias alinhadas: Identidade do alinhamento: %2.2f' % (matches/len(res1)*100) + '%')
63     print(res1)
64     print(res2)
65     print("")

```

O back trace para recuperar os alinhamentos de sequência de interesse pode começar em qualquer ponto da matriz, sendo os pontos de interesse os de valor máximo, e termina ao encontrar um valor 0.

Para cálculo de identidade considerei o comprimento das subsequências alinhadas localmente, com relação ao número de matches dentro delas (valor alto, em geral).

```

68 for i in range(size1):
69     for j in range(size2):
70         if matrix[i][j] == maxval:
71             back_trace(i, j)
72
73 print('Tabela de alinhamento:')
74 print(matrix)
75 print('\nScore do alinhamento: %d' % maxval)
76

```

Como o valor máximo de score na matriz se apresenta diversas vezes, há mais de um alinhamento possível que contém esse score máximo. Todos estes alinhamentos são retornados. Ao final é retornada a tabela e o score máximo.

Resultado:

```
Cadeias alinhadas: Identidade do alinhamento: 100.00%
GALN
GALN

Cadeias alinhadas: Identidade do alinhamento: 83.33%
GALNKV
GALNTV

Cadeias alinhadas: Identidade do alinhamento: 80.00%
Q-VPN
QQVPN

Cadeias alinhadas: Identidade do alinhamento: 100.00%
SEAE
SEAE

Cadeias alinhadas: Identidade do alinhamento: 75.00%
LSEEDRKA
LSEAERKA

Tabela de alinhamento:
[[0. 0. 0. ... 0. 0. 0.]
 [0. 1. 0. ... 0. 0. 0.]
 [0. 0. 0. ... 0. 0. 0.]
 ...
 [0. 0. 0. ... 0. 0. 0.]
 [0. 0. 0. ... 0. 0. 0.]
 [0. 0. 0. ... 0. 0. 0.]]

Score do alinhamento: 4
```

1e) Observando que o valor obtido para score máximo foi 4, e as cadeias são extremamente pequenas em relação à sequência original posso concluir duas coisas: ou as sequências possuem baixa semelhança, possuindo muito poucas regiões bem alinhadas; ou os pesos utilizados para o preenchimento da matriz são inadequados para estas sequências.

2a)

```
1  import numpy as np
2
3  word1 = 'ATAGCAGAGTCTGGGAGGGCTGTCCGCTTTTCAGTTTGTCTAAAGGTTAGGGCAGTTGAT'
4  word2 = 'ATAGCAGAGTCTGGGAGGGCTGTCCGCTTTTCAGTTTGTCTAAAGGTTAGGGCAGTTGAT'
5
6  #word1 = 'abcdefghijklmnopqrstuvwxyz'
7  #word2 = 'abcdefghijklmnopqrstuvwxyz'
8
9  size1 = len(word1)
10 size2 = len(word2)
11
12 print(size2, size1)
13 matrix_points = np.zeros((size1, size2))
14
15
16
17 for i in range(size1):
18     for j in range(size2):
19         if word1[i] == word2[j]:
20             matrix_points[i][j] = 1
21         else:
22             matrix_points[i][j] = 0
23
24
```

Implementação da matriz de pontos para as sequências fornecidas. Pela tabela final, seria possível observar onde há mutações com base em valores 0 na diagonal principal. Porém não é possível observar gaps, pois não há como diferenciar uma mutação de um gap quando se observa valores 0 na diagonal.

Para as sequências fornecidas, como há gaps (a sequência mutada tem tamanho menor que a sequência original), caímos nesse caso de não saber onde ocorre mutações e onde ocorre gaps.