

PRACTICAL 4

Question 1. (a) Write a program to construct a dot plot for the alignment of **human and chicken hemoglobin beta chain**. Identify the segments, which are same in both sequences. (b) Construct the dot plot manually for the **residues 1-20** and verify with the plot obtained using program

Solution. (a) Below is the code for the program to construct a dot plot for the alignment.

```

1 import matplotlib.pyplot as plt
2 import seaborn as sns
3
4 # Function to create the dot plot for two sequences
5 def dot_plot(seq_a, seq_b):
6     dot_plot = [[0 for i in range(len(seq_b))] for j in range(len(seq_a))]
7     for i in range(len(seq_a)):
8         for j in range(len(seq_b)):
9             if seq_a[i] == seq_b[j]:
10                 dot_plot[i][j] = 1
11     return dot_plot
12
13 seq_human = """
14     MVHLTPEEKSAVTALWGKVNVDEVGGEALGRLLVVYPWTQRFFESFGDLSTPDAMGNPKVKAHGKKVLGAFS
15 DGLAHLNLDNLKGTATLSELHCDKLHVDPENFRLLGNVLVCVLAHHFGKEFTPPVQAAYQKVVAGVANALAHKYH"""
16 seq_chick = """
17     MVHWTAEKQLITGLWGKVNAECGAEALARLLIVYPWTQRFFASFGNLSPTAILGNPMVRAHGKKVLTSFG
18 DAVKNLDNIKNTFSQLSELHCDKLHVDPENFRLLGDILIIIVLAHFSKDFTEPCQAAWQKLVRVVAHALARKYH"""
19
20 matrix = dot_plot(seq_human, seq_chick)
21 sns.heatmap(matrix)
22 plt.show()

```

LISTING 1. Code to generate dot plot

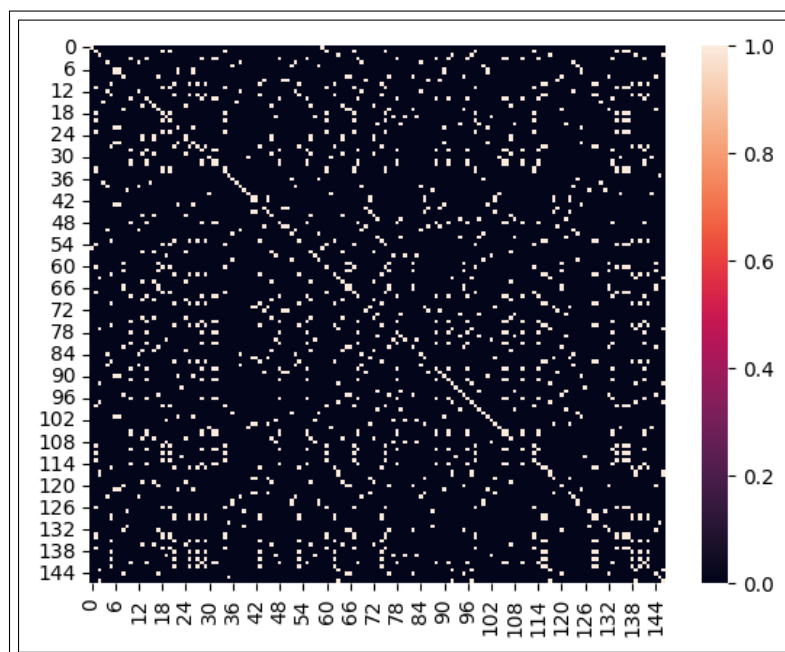


FIGURE 1. Dot plot for alignment of Human and Chicken hemoglobin chain

In the above figure, the numbers on the left and the bottom correspond to the **residue number** of the corresponding hemoglobin chain. The left denotes the **human hemoglobin beta chain**. The bottom denotes the **chicken hemoglobin beta chain**. **White** denotes a **match**, while the **dark blue** denotes a **mismatch**.

Observing the dot plot, we can see that a lot of **white points lie along the diagonal**. These correspond to the **same residue positions** on either sequence. This is indicative of the fact that the **hemoglobin beta chains** of humans and chickens have a **lot of residues in common** and in **corresponding positions**.

MVH-T-EEK—T-LWGKVN—E-G-EAL-RLL-VYPWTQRFF-SFG-LS-P-A—GNP-V-AHGKKVL—F-D—LDN-K-TF—LSELHCDKLHVDPENFRLLG—L—VLA-HF-K-FTP—QAA-QK-V—VA-ALA-KYH

Given above is the sequence where - indicates dissimilarity, while the rest positions are occupied by the respective amino acids.

(b) Considering **dot plots (both manually and program generated)** for residues **1–20** of the above **human and chicken hemoglobin beta chains**.

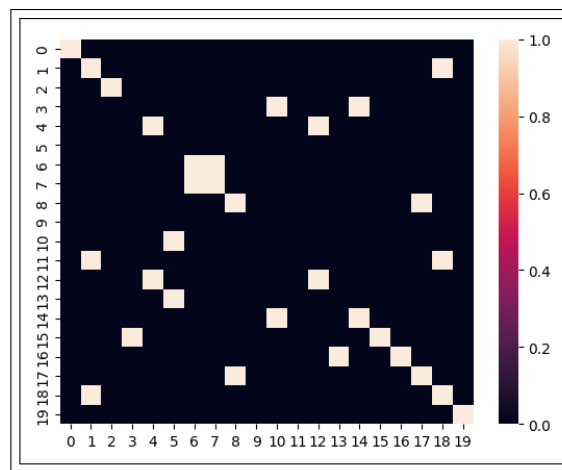


FIGURE 2. Dot plot code for alignment of first 20 residues

On comparing both plots, we can notice that both plots are **exactly the same**, thereby verifying the **correctness of the box plot code**.

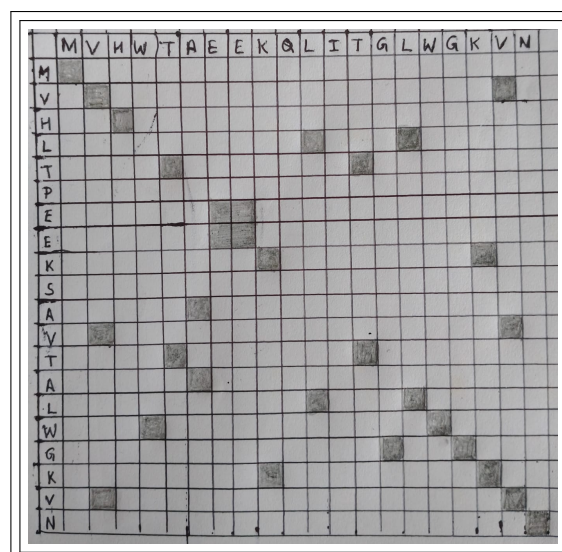


FIGURE 3. Dot plot manual for alignment of first 20 residues

Question 2. Calculate the score for the following alignments using code:

AATCTATA

AAG- -ATA

Assume that the **match score is 1**, the **mismatch score is 0**, the **origination penalty is -2**, and the **length penalty is -1**.

Solution. Below is the code for the program to compute the alignment score between two sequences. There are a set of assumptions here.

- The first sequence is always **complete without any gaps**, otherwise the concept of computing origination penalty will change relatively drastically.
- The second sequence is given **with the gaps enclosed** (if any).
- The **length of first sequence and second sequence are equal**.

```
1 def calculate_score(seq_1, seq_2, match_score, mismatch_score, origin_penalty
  , len_penalty):
2     # Considering length of seq_1 equal to seq_2 after insertion of gaps.
3     origin=True
4     score = 0
5     for i in range(len(seq_1)):
6         if (seq_2[i] == "-" and origin == True):
7             score += origin_penalty + len_penalty
8             origin=False
9         elif(seq_2[i] == "-" and origin == False):
10            score += len_penalty
11        elif(seq_2[i] == seq_1[i]):
12            score += match_score
13            origin = True
14        else:
15            score += mismatch_score
16            origin = True
17    return score
18
19 seq_1 = "AATCTATA"
20 seq_2 = "AAG--ATA"
21 match_score = +1
22 mismatch_score = 0
23 origin_penalty = -2
24 len_penalty = -1
25 score = calculate_score(seq_1, seq_2, match_score, mismatch_score,
  origin_penalty, len_penalty)
26 print(score)
```

LISTING 2. Code to compute alignment score

Alignment Score

AATCTATA

AAG- -ATA

The alignment score for above two sequences is +1

Apart from the above two sequences, I have tried computing the alignment scores between various other sequences, which were discussed in class. Below are the results of aligning those sequences.

Alignment Score

AATCTATA

AA-G-ATA

The alignment score for above two sequences is -1

Alignment Score

AATCTATA

AAG-AT-A

The alignment score for above two sequences is -3

Question 3. Verify the Q2 manually.

Solution. Below is an image where I have manually calculated the alignment score for the example given in Q2. The so manually obtained alignment score matches with the score obtained via code, which is +1.

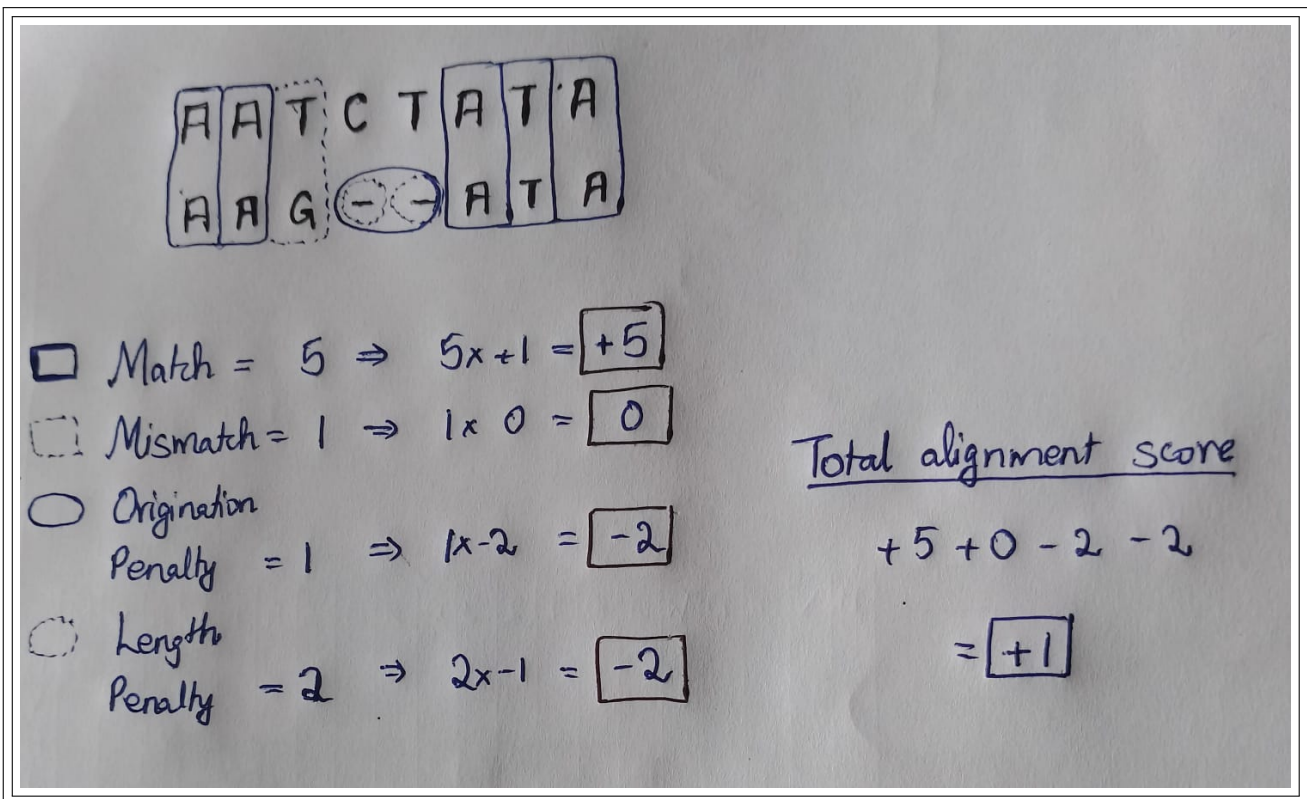


FIGURE 4. Manually calculated alignment score

Question 4. Using the Needleman and Wunsch dynamic programming method, construct the partial alignment score table and align the following two sequences (using code): ACAGTCGAACG and ACCGTCCG

use the scoring parameters: match score: +2; mismatch score: -1; and gap penalty: -2

Solution. The formula for computing the Needleman and Wunsch dynamic programming method is as follows:

$$s_{(i,j)} = \max \begin{cases} s_{(i-1,j)} + w(a_i, -) \\ s_{(i,j-1)} + w(-, b_j) \\ s_{(i-1,j-1)} + w(a_i, b_j) \end{cases} \quad w(a_i, b_j) = \begin{cases} match_score & \text{if } a_i = b_j \\ mismatch_score & \text{if } a_i \neq b_j \\ gap_penalty & \text{if } a_i = - \text{ or } b_j = - \end{cases}$$

This algorithm is widely used for optimal global alignment. The code for **Needleman and Wunsch algorithm** is given below:

```
1 # Calculate the matrix for Needleman-Wunsch algorithm
2 def needleman_wunsch_matrix(seq_1, seq_2, match_score, mismatch_score,
   gap_score):
3     matrix = [[0 for i in range(len(seq_2)+1)] for j in range(len(seq_1)+1)]
4     for i in range(len(seq_2)+1):
5         matrix[0][i] = gap_score*i
6     for i in range(len(seq_1)+1):
7         matrix[i][0] = gap_score*i
8     for i in range(1, len(seq_1)+1):
9         for j in range(1, len(seq_2)+1):
10            val_1 = matrix[i-1][j] + gap_score
11            val_2 = matrix[i][j-1] + gap_score
12            val_3 = matrix[i-1][j-1] + match_score if seq_1[i-1] ==
13                seq_2[j-1] else matrix[i-1][j-1] + mismatch_score
14            matrix[i][j] = max(val_1, val_2, val_3)
15     return matrix
16
17 # Backtrack to find the aligned sequences
18 def backtrack(matrix, seq_1, seq_2, match_score, mismatch_score, gap_score):
19     seq_align_1 = ""
20     seq_align_2 = ""
21     i = len(seq_1)
22     j = len(seq_2)
23     while(i>0 and j>0):
24         if(matrix[i-1][j] == matrix[i][j] - gap_score):
25             seq_align_1 += seq_1[i-1]
26             seq_align_2 += "-"
27             i -= 1
28         elif(matrix[i][j-1] == matrix[i][j] - gap_score):
29             seq_align_2 += seq_2[j-1]
30             seq_align_1 += "-"
31             j -= 1
32         else:
33             seq_align_1 += seq_1[i-1]
34             seq_align_2 += seq_2[j-1]
35             i -= 1
36             j -= 1
37     if(i == 0 and j != 0):
38         for k in range(j-1, -1, -1):
39             seq_align_2 += seq_2[k]
40             seq_align_1 += "-"
41     elif(j == 0 and i != 0):
42         for k in range(i-1, -1, -1):
43             seq_align_1 += seq_1[k]
44             seq_align_2 += "-"
45     seq_align_1 = seq_align_1[::-1]
46     seq_align_2 = seq_align_2[::-1]
47
48     return seq_align_1, seq_align_2
49
50 seq_1 = "ACAGTCGAACG"
51 seq_2 = "ACCGTCCG"
52 match_score = +2
```

```

53 mismatch_score = -1
54 gap_score = -2
55 matrix = needleman_wunsch_matrix(seq_1, seq_2, match_score, mismatch_score,
    gap_score)
56 seq1, seq2 = backtrack(matrix, seq_1, seq_2, match_score, mismatch_score,
    gap_score)
57 print(seq1)
58 print(seq2)
59
60 # Output results of seq1 and seq2
61 ACAGTCGAACG
62 ACCGTC---CG

```

LISTING 3. Code to compute sequences in Needleman Wunsch Algorithm

The partial alignment score table for the above two sequences in the question is given below. Also enclosed below is the code to generate the **partial alignment score table**. The **matrix** in the code is taken from the above code:

```

1 seq_1 = "ACAGTCGAACG"
2 seq_2 = "ACCGTCCG"
3 sns.heatmap(matrix, annot=True)
4 plt.xticks(ticks=np.arange(len(matrix[0])) + 0.5, labels=["-", "A", "C", "C",
    "G", "T", "C", "C", "G"])
5 plt.yticks(ticks=np.arange(len(matrix)) + 0.5, labels=["-", "A", "C", "A", "G",
    "T", "C", "G", "A", "A", "C", "G"])
6 plt.show()

```

LISTING 4. Code to generate alignment table in Needleman Wunsch Algorithm

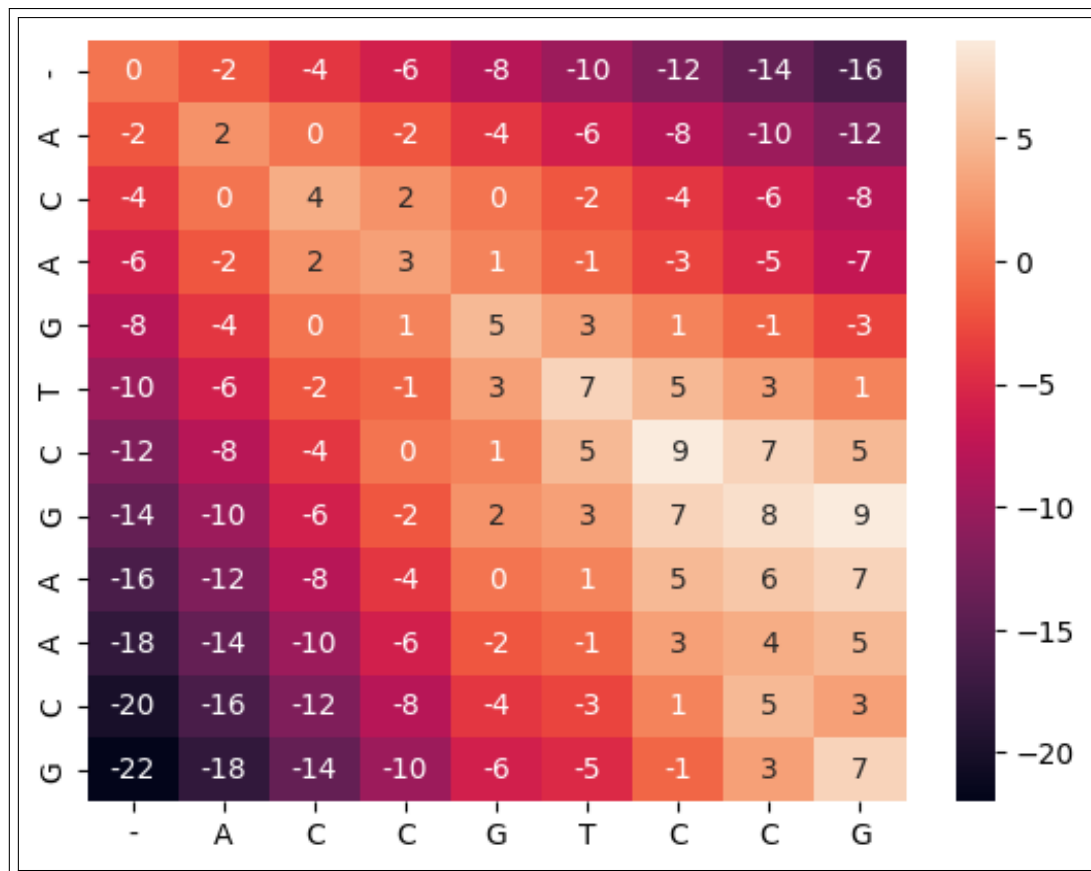


FIGURE 5. Partial alignment score table

Question 5. Verify Q4 manually

Solution. The partial alignment score table obtained in manually computing is shown in the image below. The table so obtained is found to be **exactly same** as the one computed by the code above, thereby verifying the correctness of the Needleman Wunsch algorithm code.

	-	A	C	C	G	T	C	C	G	
-	0	-2	-4	-6	-8	-10	-12	-14	-16	
A	-2	2	0	-2	-4	-6	-8	-10	-12	
C	-4	0	4	2	0	-2	-4	-6	-8	
A	-6	-2	2	3	1	-1	-3	-5	-7	
G	-8	-4	0	1	5	3	1	-1	-3	
T	-10	-6	-2	-1	3	7	5	3	1	
C	-12	-8	-4	0	1	5	9	7	5	
G	-14	-10	-6	-2	2	3	7	8	9	
A	-16	-12	-8	-4	0	1	5	6	7	
A	-18	-14	-10	-6	-2	-1	3	4	5	
C	-20	-16	-12	-8	-4	-3	1	5	3	
G	-22	-18	-14	-10	-6	-5	-1	3	7	

FIGURE 6. Manual Needleman Wunsch Algorithm

Question 6. Using the **Smith-Waterman method**, construct the partial alignment scoring table and align the following two sequences (using code):

ACGTATCGCGTATA and **GATGCGTATCG**

scoring parameters: **match score: +2; mismatch score: -1; and gap penalty: -2.**

Solution. The formula for computing the **Smith and Waterman dynamic programming method** is as follows:

$$s_{(i,j)} = \max \begin{cases} s_{(i-1,j)} + w(a_i, -) \\ s_{(i,j-1)} + w(-, b_j) \\ s_{(i-1,j-1)} + w(a_i, b_j) \\ 0 \end{cases} \quad w(a_i, b_j) = \begin{cases} match_score & \text{if } a_i = b_j \\ mismatch_score & \text{if } a_i \neq b_j \\ gap_penalty & \text{if } a_i = - \text{ or } b_j = - \end{cases}$$

This algorithm is widely used for optimal local alignment. The code for **Smith and Waterman algorithm** is given below:

```

1 # Calculate the matrix for Smith-Waterman algorithm
2 def smith_waterman_matrix(seq_1, seq_2, match_score, mismatch_score,
   gap_score):
3     matrix = [[0 for i in range(len(seq_2)+1)] for j in range(len(seq_1)+1)]
4     for i in range(len(seq_2)+1):
5         matrix[0][i] = gap_score*i
6     for i in range(len(seq_1)+1):
7         matrix[i][0] = gap_score*i
8     for i in range(1, len(seq_1)+1):
9         for j in range(1, len(seq_2)+1):
10            val_1 = matrix[i-1][j] + gap_score
11            val_2 = matrix[i][j-1] + gap_score
12            val_3 = matrix[i-1][j-1] + match_score if seq_1[i-1] ==
13                seq_2[j-1] else matrix[i-1][j-1] + mismatch_score
14            matrix[i][j] = max(val_1, val_2, val_3, 0)
15     return matrix
16
17 # Backtrack to find the local aligned sequence
18 def backtrack(matrix, seq_1, seq_2, match_score, mismatch_score, gap_score):
19     seq_align_1 = ""
20     seq_align_2 = ""
21     i = 0
22     j = 0
23     maximum = 0
24     for r in range(len(seq_1)+1):
25         for c in range(len(seq_2)+1):
26             if matrix[r][c] > maximum:
27                 maximum = matrix[r][c]
28                 i = r
29                 j = c
30
31     while(matrix[i][j] != 0):
32         if(matrix[i-1][j] == matrix[i][j] - gap_score):
33             seq_align_1 += seq_1[i-1]
34             seq_align_2 += "-"
35             i -= 1
36         elif(matrix[i][j-1] == matrix[i][j] - gap_score):
37             seq_align_2 += seq_2[j-1]
38             seq_align_1 += "-"
39             j -= 1
40     else:

```



```

41         seq_align_1 += seq_1[i-1]
42         seq_align_2 += seq_2[j-1]
43         i -= 1
44         j -= 1
45
46     seq_align_1 = seq_align_1[::-1]
47     seq_align_2 = seq_align_2[::-1]
48
49     index_start_1 = seq_1.index(seq_align_1)
50     index_start_2 = seq_2.index(seq_align_2)
51
52     index_end_1 = index_start_1 + len(seq_align_1)
53     index_end_2 = index_start_2 + len(seq_align_2)
54
55     prefix_1 = "-"*(max(index_start_1, index_start_2) - index_start_1) +
56               seq_1[:index_start_1]
57     prefix_2 = "-"*(max(index_start_1, index_start_2) - index_start_2) +
58               seq_2[:index_start_2]
59
60     suffix_1 = seq_1[index_end_1:] + "-"*(max(len(seq_1) - index_end_1,
61         len(seq_2) - index_start_2) - (len(seq_1) - index_end_1 + 1))
62     suffix_2 = seq_2[index_end_2:] + "-"*(max(len(seq_1) - index_end_1,
63         len(seq_2) - index_start_2) - (len(seq_2) - index_end_2 + 1))
64
65     seq_align_1 = prefix_1 + seq_align_1 + suffix_1
66     seq_align_2 = prefix_2 + seq_align_2 + suffix_2
67
68     return seq_align_1, seq_align_2
69
70 seq_1 = "ACGTATCGCGTATA"
71 arr_seq_1 = list("-"+seq_1)
72 seq_2 = "GATGCGTATCG"
73 arr_seq_2 = list("-"+seq_2)
74 match_score = +2
75 mismatch_score = -1
76 gap_score = -2
77 matrix = smith_waterman_matrix(seq_1, seq_2, match_score, mismatch_score,
78     gap_score)
79 seq1, seq2 = backtrack(matrix, seq_1, seq_2, match_score, mismatch_score,
80     gap_score)
81
82 # Output results of seq1 and seq2
83 ---ACGTATCGCGTATA
84 GATGCGTATCG-----
85 ....CGTATCG..... # The matching local alignment

```

LISTING 5. Code for Smith and Waterman Algorithm

The partial alignment score table for the above two sequences in the question is given below. Also enclosed below is the code to generate the **partial alignment score table**. The **matrix** in the code is taken from the above code:

```

1 seq_1 = "ACGTATCGCGTATA"
2 arr_seq_1 = list("-"+seq_1)
3 seq_2 = "GATGCGTATCG"

```

```

4 arr_seq_2 = list("-"+seq_2)
5
6 # Plot the heatmap
7 sns.heatmap(matrix, annot=True)
8 plt.xticks(ticks=np.arange(len(matrix[0])) + 0.5, labels=arr_seq_2)
9 plt.yticks(ticks=np.arange(len(matrix)) + 0.5, labels=arr_seq_1)
10 plt.show()

```

LISTING 6. Code to generate alignment table in Smith and Waterman algorithm

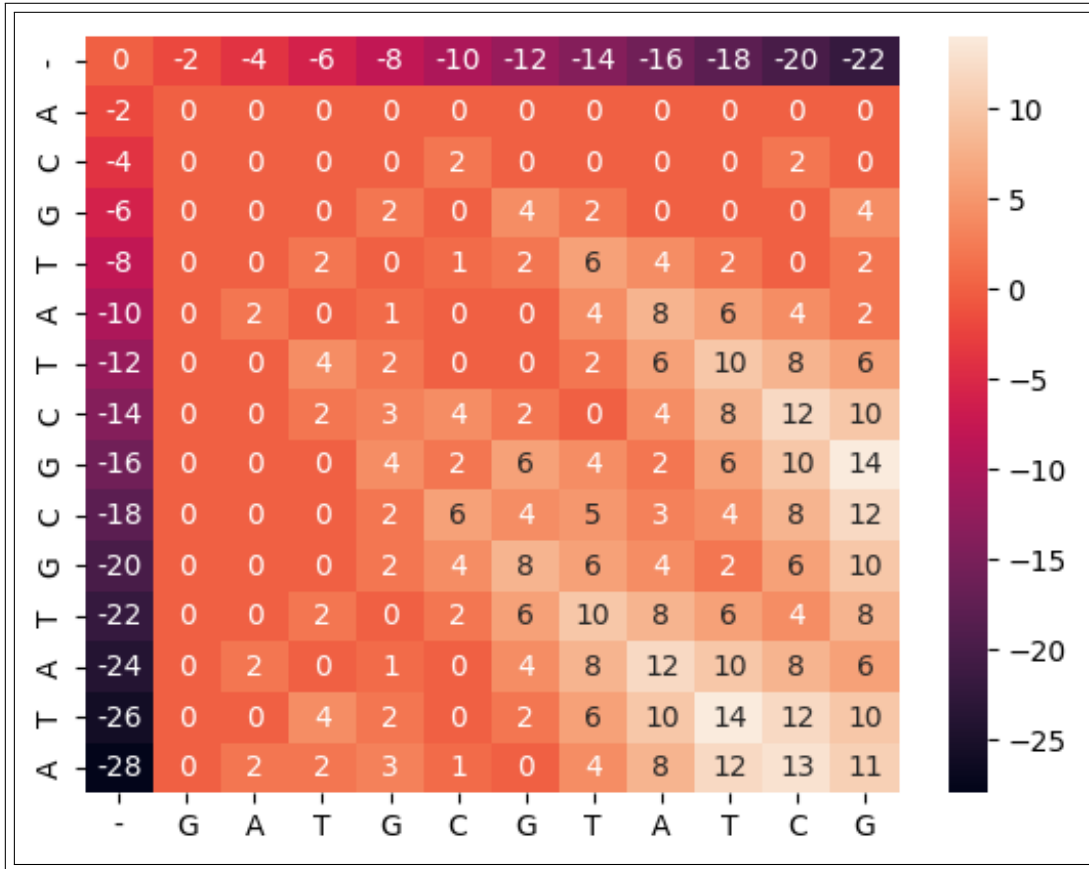


FIGURE 7. Partial alignment score table