Programming for Bioinformatics | BIOL7200

Dynamic Programming Algorithm

As mentioned in the class, there is a main graded assignment and a bonus, extra point assignment this week. The first few pages describe the algorithm at length. The text below summarizes the algorithm.

For this week as well, assume that the user gives you correct inputs all the time. **Your script will be graded on the output produced and not how all the errors are handled.**

The CI for this week will also not be public.

**Again, please do not use any modules other than "sys".** The CI will not install missing modules. Do not use input() for any input either, we do not handle this in the CI. The CI **will fail** if you do not follow these instructions.

## Instructions for submission
- **This assignment is due Monday, October 22, 2018 at 11:59pm. Late submissions will receive a 0**
- Your code must be available on GitLab at the above time to be graded
- Name the Needleman-Wunsch algorithm script as **nwAlign.py**
- Optional bonus assignment: Name the Smith-Waterman algorithm script as **swAlign.py**
- Your code should run as ./nwAlign.py <seq1.fa> <seq2.fa>  OR ./swAlign.py
- Both scripts should output their results to STDOUT
- DO NOT HARDCODE any file name!
- Please use the **#!/usr/bin/env python** as your shebang


Main assignment: Needleman-Wunsch (NW) algorithm

Max score: 100 points

This is an example of another complex problem that has a rather simple solution. NW algorithm is a classical bioinformatics algorithm designed to obtain optimal global alignment for a given pair of sequences. The algorithm falls under the class of dynamic programming which in simple language is the class of algorithm that work by breaking a problem into subproblems, solving each subproblem and joining the solutions to reach the global solution.

The algorithm can be divided into three steps:

1. *Initialization*: Construction of the matrix with the two sequences as each axis and selection of a suitable scoring system. For simplicity, let's have three types of scores:
    a. Match = +1
    b. Mismatch = -1
    c. Gap = -1

2. *Matrix filling*: Filling the matrix based on the scoring system. This occurs one row at a time, starting from the topmost row. Each cell in the matrix derives the value from the adjacent cells located to the left, top-left or on top of the current cell. Match score is added or gap/mismatch penalty is

subtracted from these adjacent cells and the maximum value is carried over to the current cell (Figure 1).

3. *Backtracking*: Once the matrix has been filled up, backtracking is done to compute the optimal alignment(s). The backtracking step starts from the very last cell filled in the matrix (the bottom-right cell) and proceeds to the first cell filled in matrix (the cell with 0 in the upper left corner of the matrices in Figure 1). This backtrack path is computed by moving through the adjacent cells (cells to the left, top-left and on top of the current cell) with the maximum score such that the path has the maximum total score (Figure 2). If multiple paths exist, then all of them are considered to be the optimal paths. This path is converted to an alignment by the following rule: the path moves diagonally to the left if there is a match or if the maximum score of the adjacent cells is present in the diagonal left cell. If either of these are true, the two corresponding characters from each sequence are aligned together. When the maximum score is obtained by moving horizontally, then a gap is introduced in the sequence on the vertical axis, and if the path moves vertically, then a gap is introduced in the sequence on the horizontal axis.

Backtracking rules:

1. Always take the diagonal when the diagonal is either (1) the highest score or (2) tied for highest score
2. If the diagonal is not the highest score, take the "Up" if it is either (1) the highest score or (2) tied for highest score.
3. Take the "Left" if the diagonal and "Up" are not the highest

For the purpose of this assignment, you will only observe a single optimal path with the above rules in the test sequences we use. You will not have to worry about multiple, optimal paths.

*As before, do not think too much or you will spend most of the week doing this. Think simple, use arrays and loops. If you are not too confident with coding, start by writing pseudocode. When you have it right, then code it up. Ideally, designing the pseudocode should take 80% of your time, translating it into code should take 10% and testing should take 10%. If you aren't sincere with the pseudocode, you are going to have a bad time with this assignment.*

<u>Syntax</u>: `./nwAlign.py <input FASTA file 1> <input FASTA file 2>`

<u>Example usage</u>: `./nwAlign.py seq1.fa seq2.fa`

seq1.fa contains:

>seq1.fa
AATTCCTT

seq2.fa contains:

>seq2.fa
AACTCTT


<u>Output format</u>:

```
AATTCCTT
||| | |||
AACT-CTT
Alignment score: 4
```

Match = +1  Mismatch = -1  Gap = -1

D is the cell to be filled, it takes the maximum of (A +/- Match/Mismatch), (B - Gap) & (C - Gap)

|   |   | G | C | A | T | G | C |
|---|---|---|---|---|---|---|---|
|   | 0 |   |   |   |   |   |   |
| G |   |   |   |   |   |   |   |
| A |   |   |   |   |   |   |   |
| T |   |   |   |   |   |   |   |
| T |   |   |   |   |   |   |   |
| A |   |   |   |   |   |   |   |
| C |   |   |   |   |   |   |   |

|   |   | G | C | A | T | G | C |
|---|---|---|---|---|---|---|---|
|   | 0 | -1 | -2 | -3 | -4 |   |   |
| G |   |   |   |   |   |   |   |
| A |   |   |   |   |   |   |   |
| T |   |   |   |   |   |   |   |
| T |   |   |   |   |   |   |   |
| A |   |   |   |   |   |   |   |
| C |   |   |   |   |   |   |   |

|   |   | G | C | A | T | G | C |
|---|---|---|---|---|---|---|---|
|   | 0 | -1 | -2 | -3 | -4 | -5 | -6 |
| G | -1 | +1 |   |   |   |   |   |
| A |   |   |   |   |   |   |   |
| T |   |   |   |   |   |   |   |
| T |   |   |   |   |   |   |   |
| A |   |   |   |   |   |   |   |
| C |   |   |   |   |   |   |   |

|   |   | G | C | A | T | G | C |
|---|---|---|---|---|---|---|---|
|   | 0 | -1 | -2 | -3 | -4 | -5 | -6 |
| G | -1 | +1 | 0 | -1 | -2 | -3 | -4 |
| A |   |   |   |   |   |   |   |
| T |   |   |   |   |   |   |   |
| T |   |   |   |   |   |   |   |
| A |   |   |   |   |   |   |   |
| C |   |   |   |   |   |   |   |

|   |   | G | C | A | T | G | C |
|---|---|---|---|---|---|---|---|
|   | 0 | -1 | -2 | -3 | -4 | -5 | -6 |
| G | -1 | +1 | 0 | -1 | -2 | -3 | -4 |
| A | -2 | 0 | 0 | 1 | 0 | -1 | -2 |
| T | -3 | -1 | -1 | 0 | 2 | 1 | 0 |
| T | -4 | -2 | -2 | -1 | 1 | 1 | 0 |
| A | -5 | -3 | -3 | -1 | 0 | 0 |   |
| C |   |   |   |   |   |   |   |

|   |   | G | C | A | T | G | C |
|---|---|---|---|---|---|---|---|
|   | 0 | -1 | -2 | -3 | -4 | -5 | -6 |
| G | -1 | +1 | 0 | -1 | -2 | -3 | -4 |
| A | -2 | 0 | 0 | 1 | 0 | -1 | -2 |
| T | -3 | -1 | -1 | 0 | 2 | 1 | 0 |
| T | -4 | -2 | -2 | -1 | 1 | 1 | 0 |
| A | -5 | -3 | -3 | -1 | 0 | 0 | 0 |
| C | -6 | -4 | -2 | -2 | -1 | -1 | 1 |

Figure 1. Matrix filling step of dynamic programming

|     | G  | C  | A  | T  | G  | C  |
|-----|----|----|----|----|----|----|
| 0   | -1 | -2 | -3 | -4 | -5 | -6 |
| G -1 | +1 | 0  | -1 | -2 | -3 | -4 |
| A -2 | 0  | 0  | 1  | 0  | -1 | -2 |
| T -3 | -1 | -1 | 0  | 2  | 1  | 0  |
| T -4 | -2 | -2 | -1 | 1  | 1  | 0  |
| A -5 | -3 | -3 | -1 | 0  | 0  | 0  |
| C -6 | -4 | -2 | -2 | -1 | -1 | 1  |

Starting point of the path

|     | G  | C  | A  | T  | G  | C  |
|-----|----|----|----|----|----|----|
| 0   | -1 | -2 | -3 | -4 | -5 | -6 |
| G -1 | +1 | 0  | -1 | -2 | -3 | -4 |
| A -2 | 0  | 0  | 1  | 0  | -1 | -2 |
| T -3 | -1 | -1 | 0  | 2  | 1  | 0  |
| T -4 | -2 | -2 | -1 | 1  | 1  | 0  |
| A -5 | -3 | -3 | -1 | 0  | 0  | 0  |
| C -6 | -4 | -2 | -2 | -1 | -1 | 1  |

Optimal alignment:

```
C
C
```

|     | G  | C  | A  | T  | G  | C  |
|-----|----|----|----|----|----|----|
| 0   | -1 | -2 | -3 | -4 | -5 | -6 |
| G -1 | +1 | 0  | -1 | -2 | -3 | -4 |
| A -2 | 0  | 0  | 1  | 0  | -1 | -2 |
| T -3 | -1 | -1 | 0  | 2  | 1  | 0  |
| T -4 | -2 | -2 | -1 | 1  | 1  | 0  |
| A -5 | -3 | -3 | -1 | 0  | 0  | 0  |
| C -6 | -4 | -2 | -2 | -1 | -1 | 1  |

Optimal alignments:

```
-C        GC
AC        AC
```

|     | G  | C  | A  | T  | G  | C  |
|-----|----|----|----|----|----|----|
| 0   | -1 | -2 | -3 | -4 | -5 | -6 |
| G -1 | +1 | 0  | -1 | -2 | -3 | -4 |
| A -2 | 0  | 0  | 1  | 0  | -1 | -2 |
| T -3 | -1 | -1 | 0  | 2  | 1  | 0  |
| T -4 | -2 | -2 | -1 | 1  | 1  | 0  |
| A -5 | -3 | -3 | -1 | 0  | 0  | 0  |
| C -6 | -4 | -2 | -2 | -1 | -1 | 1  |

Optimal alignments:

```
G-C       -GC       TGC
TAC       TAC       TAC
```

|     | G  | C  | A  | T  | G  | C  |
|-----|----|----|----|----|----|----|
| 0   | -1 | -2 | -3 | -4 | -5 | -6 |
| G -1 | +1 | 0  | -1 | -2 | -3 | -4 |
| A -2 | 0  | 0  | 1  | 0  | -1 | -2 |
| T -3 | -1 | -1 | 0  | 2  | 1  | 0  |
| T -4 | -2 | -2 | -1 | 1  | 1  | 0  |
| A -5 | -3 | -3 | -1 | 0  | 0  | 0  |
| C -6 | -4 | -2 | -2 | -1 | -1 | 1  |

Optimal alignments:

```
ATG-C     AT-GC     A-TGC
ATTAC     ATTAC     ATTAC
```

|     | G  | C  | A  | T  | G  | C  |
|-----|----|----|----|----|----|----|
| 0   | -1 | -2 | -3 | -4 | -5 | -6 |
| G -1 | +1 | 0  | -1 | -2 | -3 | -4 |
| A -2 | 0  | 0  | 1  | 0  | -1 | -2 |
| T -3 | -1 | -1 | 0  | 2  | 1  | 0  |
| T -4 | -2 | -2 | -1 | 1  | 1  | 0  |
| A -5 | -3 | -3 | -1 | 0  | 0  | 0  |
| C -6 | -4 | -2 | -2 | -1 | -1 | 1  |

Optimal alignments:

```
GCATG-C   GCAT-GC   GCA-TGC
G-ATTAC   G-ATTAC   G-ATTAC
```
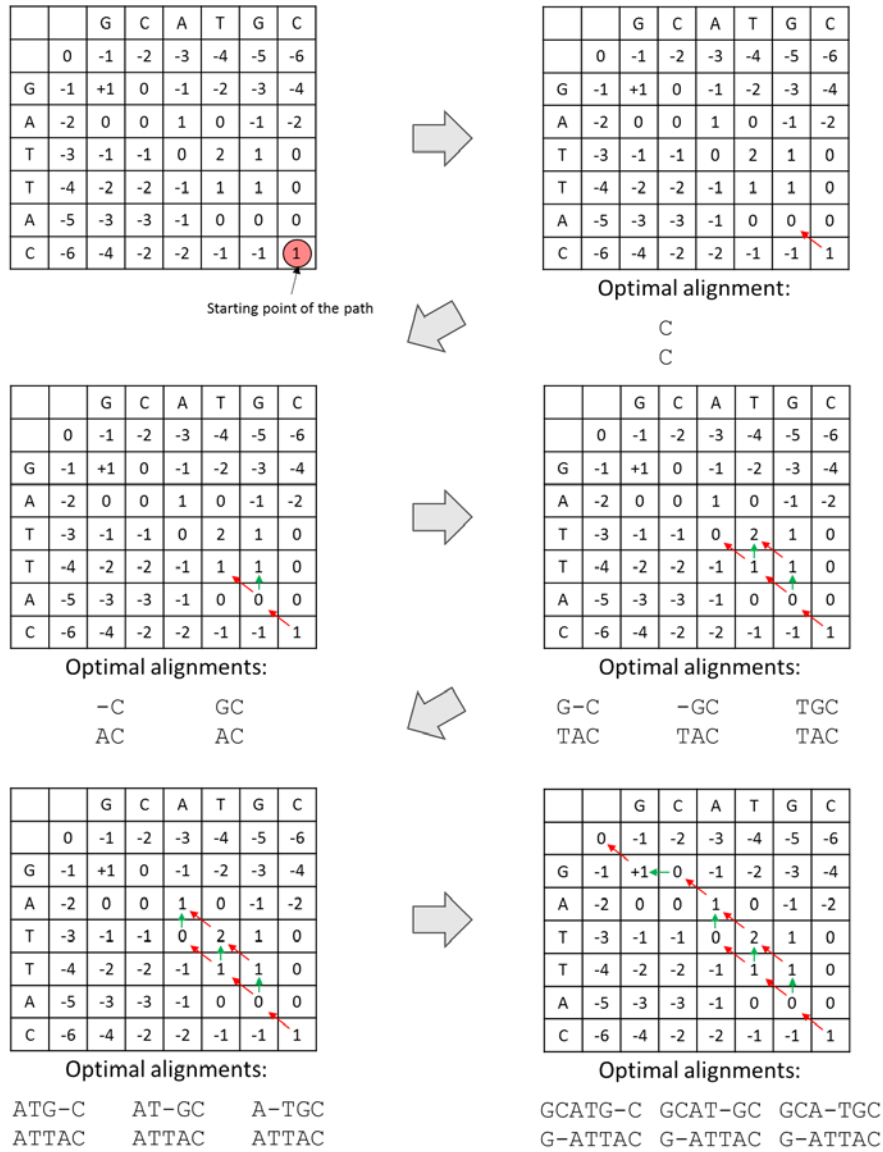
Figure 2. Matrix backtracking step and the generation of optimal alignments

Bonus assignment: Smith-Waterman (SW) algorithm

Bonus score: 100 points