

---

### SMITH-WATERMAN ALGORITHM (15%)

---

Bioinformatics is a field that uses techniques from computer science, statistics, and engineering to help study large amounts of biological data. Analysis of protein structure and sequence is very important in bioinformatics, especially for understanding how cells work, which makes it easier to make drugs for metabolic pathways. Protein sequence alignment is a method for figuring out the relationships between different protein structures by finding the similarities between their structures.

The Smith-Waterman algorithm uses local sequence alignment to find similar regions between two strings of nucleic acid or protein sequences. Rather than comparing the entire sequence, the Smith-Waterman algorithm compares segments of varying lengths and optimizes the similarity measure.

---

### MOTIVATION

---

Genome projects on various organisms have generated massive amounts of sequence data for genes and proteins in recent years, necessitating computational analysis. Sequence alignment reveals the relationships between genes or proteins, allowing us to understand their homology and functionality better. The alignment of sequences can also reveal conserved domains and motifs.

---

### ALGORITHM DESIGN

---

This algorithm aims to compare the structures of two amino acids and is proposed for use in local sequence alignment in bioinformatics. The first row and column are both initialized to a value of zero. Then, **Equation 1** is applied to each cell in the matrix. Set the Gap to -2. Score 1 when the two strings being compared are identical and score 0 otherwise. This procedure must be repeated until the entire matrix has been filled in. When this happens, the matrix's last cell incurs the dissimilarity cost. **Figure 1** shows the steps for algorithm written in pseudocode.

$$H_{i,j} = \text{MAX} \begin{cases} \text{MAX}(H_{i-1,j-1} + \text{Score}, 0) \\ \text{MAX}(H_{i-1,j} + \text{Gap}, \text{Gap}) \\ \text{MAX}(H_{i,j-1} + \text{Gap}, \text{Gap}) \end{cases} \quad \text{Equation 1}$$

---

#### Algorithm 1 Sequential Implementation of Smith-Waterman

---

```

1: procedure SW(Str1, Str2, N)
2:   Initial first row and first column from 1 to N
3:   for <Row = 0 to N -1> do
4:     for <Column = 0 to N -1> do
5:       if Str1[Row - 1] == Str2[Column - 1] then
6:         Score = 0
7:       else
8:         Score = 1
9:       end if
10:      H[Row][Column] = Calculate Distance using
          Equation (1)
11:     end for
12:   end for
13: end procedure

```

---

**Figure 1:** pseudocode of Smith-Waterman sequence alignment

## Exercise 02

Due Data: 14-Sep-2022, 13:30 PM

Table 1: deoxyribonucleic acid structure

Nucleotide Name	Code
Adenine	A
Guanine	G
Cytosine	C
Thymine	T

Random DNA Sequence results
>random sequence 1 consisting of 10 bases. GTCGATTGA
>random sequence 2 consisting of 10 bases. ACGAAAGAGG
>random sequence 3 consisting of 10 bases. TTCTGTCTGA
>random sequence 4 consisting of 10 bases. ACTGGCATTG
>random sequence 5 consisting of 10 bases. CGTACCGTCG
>random sequence 6 consisting of 10 bases. GGCCCAAGGA
>random sequence 7 consisting of 10 bases. CAATAGTCCG
>random sequence 8 consisting of 10 bases. TTACAGGTGT
>random sequence 9 consisting of 10 bases. ACTTGTCGTC
>random sequence 10 consisting of 10 bases. CTACGCCTAC

Figure 2 displays an example to calculate the similarity between sequence 1 and sequence 2 using online tool<sup>1</sup>. You have to implement Java code for this algorithm to measure the similarity between any DNA sequences. You can use the previous 10 sequences as test cases to your implementation.

Grading schema:

- Read from text file (5 Pts): **Note: you need to handle I/O exceptions.**
- Java implementation for Smith-Waterman algorithm (5 Pts).
- Presentation (5 Pts)

---

<sup>1</sup> <https://rna.informatik.uni-freiburg.de/Teaching/index.jsp?toolName=Smith-Waterman>

## Exercise 02

Due Data: 14-Sep-2022, 13:30 PM

Input:

Sequence *a*: 
  
Sequence *b*: 
  
Scoring in *s*: Match  Mismatch  Gap 
  
Hint:
 For similarity maximization,
 match scores should be positive and all other scores lower.

Recursion:
 
$$S_{i,j} = \max \begin{cases} S_{i-1,j-1} + s(a_i, b_j) \\ S_{i-1,j} + s(a_i, -) \\ S_{i,j-1} + s(-, b_j) \\ 0 \end{cases} = \max \begin{cases} S_{i-1,j-1} + 1 & a_i = b_j \\ S_{i-1,j-1} + 0 & a_i \neq b_j \\ S_{i-1,j} + -2 & b_j = - \\ S_{i,j-1} + -2 & a_i = - \\ 0 \end{cases}$$

Output:

<i>S</i>		A <sub>1</sub>	C <sub>2</sub>	G <sub>3</sub>	A <sub>4</sub>	A <sub>5</sub>	A <sub>6</sub>	G <sub>7</sub>	A <sub>8</sub>	G <sub>9</sub>	G <sub>10</sub>
	0	0	0	0	0	0	0	0	0	0	0
G <sub>1</sub>	0	0	0	1	0	0	0	1	0	1	1
T <sub>2</sub>	0	0	0	0	1	0	0	0	1	0	1
C <sub>3</sub>	0	0	1	0	0	1	0	0	0	1	0
G <sub>4</sub>	0	0	0	2	0	0	1	1	0	1	2
A <sub>5</sub>	0	1	0	0	3	1	1	1	2	0	1
T <sub>6</sub>	0	0	1	0	1	3	1	1	1	2	0
T <sub>7</sub>	0	0	0	1	0	1	3	1	1	1	2
T <sub>8</sub>	0	0	0	0	1	0	1	3	1	1	1
G <sub>9</sub>	0	0	0	1	0	1	0	2	3	2	2
A <sub>10</sub>	0	1	0	0	2	1	2	0	3	3	2

Score: 3

Results

You can select a result to get the related traceback.

CGA T T T GA

\*\*\* | | \*\*

CGA \_ AAGA

CGATTTGA

\*\*\* | | \*\*

CGAA \_ AGA

CGATTTGA

\*\*\* | | \*\*

CGAAA \_ GA

CGATTTGA

\*\*\* | | | |

CGAAAGAG

Figure 2: Shows an example for Smith-Waterman Algorithm.