

Sequence Alignment

It is a way of arranging two (or) more "sequences" of characters to identify regions of similarity.

Sequence:- can be taken as order of string of letters.

↳ These sequences can be DNA, RNA, Protein

sequences.

Dynamic programming: The process of breaking down the

problem into subproblem. Stores the result of the

subproblem - which will help to avoid the

computing same results again and again.

↳ Dynamic programming will not be useful when

there are no common (overlapping) subproblems.

Needleman - Wunsch Algorithm → This algorithm will

introduce the dynamic algorithm paradigm.

which will provide the optimal solution - will

the corresponding score.

Sequence Alignment: The process of alignment of nucleotide

/ character which will maximize the number of matches

and minimize the number of mismatches and gap.

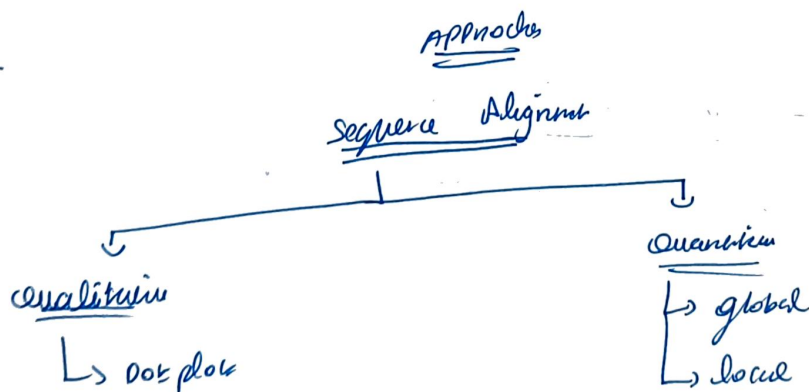
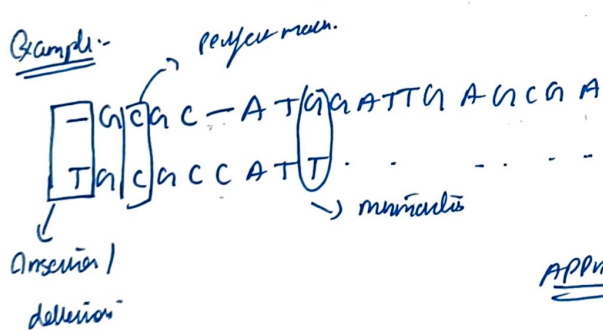
The algorithm which defines a recurrence relation which will provide an reward for a match / character match, and penalty for a mismatch.

Results of Sequence Alignment:

when we are giving two sequences. The only the sequence alignment can be. consisting

- (i) Insertion / deletion (indel)
- (ii) mismatches
- (iii) Perfect match.

Example:-



→ pictorial representation + relationship between two sequences.

→ uses a table.

→ does not guarantee

Quantitative:

⇒ Construction of the best alignment between the sequences

Global sequence: The best alignment over the entire length of two sequences.



when similar: if two sequences are of same length, with a significant degree of similarity throughout.

Local sequence: Comparison should be proportional of sequences

↳ when similar: when short length of similarity and difference in string length.

How to quantify?

⇒ Introduced a scoring schema. → all cells have an

alignment score. → goodness of alignment.

Substitution Matrix:

20 nucleic acids contain A, C, G, T. (which also contain a scoring function σ) → which will define a score to

various substitution matrix.

↳ Match → 2

mis-match → -1

| | A | C | G | T |
|---|----|----|----|----|
| A | 2 | -1 | -1 | -1 |
| C | -1 | 2 | -1 | -1 |
| G | -1 | -1 | 2 | -1 |
| T | -1 | -1 | -1 | 2 |

hap: \rightarrow consecutive run of spaces is an alignment
(insertion or deletion).

Gap penalty: \rightarrow defines a score to give insertion or deletion

$$\hookrightarrow -1 // (a) 0$$

for any dynamic programming, we need to define
a recurrence relation.

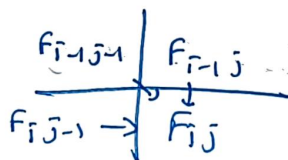
$$T(i, j) = \max \begin{cases} T(i-1, j-1) + \sigma(s_1(i), s_2(j)) & \nearrow \text{match / mismatch} \\ T(i-1, j) + \text{gap penalty} \\ T(i, j-1) + \text{gap penalty} \end{cases}$$

Initialization:

$$F(0, 0) = 0$$

$$T(i, 0) = F(i-1, 0) - \text{gap penalty}$$

$$T(0, j) = T(0, j-1) - \text{gap penalty}.$$



Termination: Bottom right.

Steps involved:

\hookrightarrow Step 1: Fill up the matrix table (T) using the recurrence relation.

\rightarrow Step 2: Traceback the matrix to work out for a best alignment.

Sequences:

$$S_1 = \blacktriangle T A G T A$$

$$S_2 = A T C G T$$

Table Initialization:

| | $\bar{i}=0$ | $\bar{i}=1$ | $\bar{i}=2$ | $\bar{i}=3$ | $\bar{i}=4$ | $\bar{i}=5$ |
|-------------|-------------|-------------|-------------|-------------|-------------|-------------|
| $\bar{j}=0$ | | T | G | G | T | G |
| $\bar{j}=1$ | A | 0 | -1 | -1 | -1 | -1 |
| $\bar{j}=2$ | T | 0 | 2 | 1 | 0 | 1 |
| $\bar{j}=3$ | C | 0 | 1 | 1 | 0 | 0 |
| $\bar{j}=4$ | G | 0 | 0 | 3 | 3 | 2 |
| $\bar{j}=5$ | T | 0 | 2 | 2 | 2 | 5 |

$$A \begin{bmatrix} & A & C & G & T \\ \begin{bmatrix} 2 & -1 & -1 & -1 \\ -1 & 2 & -1 & -1 \\ -1 & -1 & 2 & -1 \\ -1 & -1 & -1 & 2 \end{bmatrix} \end{bmatrix}$$

$$\text{Alignment} = 5 //$$

→ step-1 ⇒ The value of $T(0,0)$ is zero → Starting condition

→ step-2 ⇒ Gap penalty → 0

→ step-3 ⇒ Match → +2, mismatch = -1

→ Initialization.

→ Look for strand changes //

$$S_1 = T A G T A$$

$$S_2 = A T C G T$$

$$\begin{array}{cccccc} & T & G & G & T & G \\ & 1 & & 1 & 1 & \\ A & T & C & G & T & - \end{array}$$

$$\text{Score} = \text{Number of matches} + \text{number of mismatches} + \text{gap in}$$

$$= 3(2) + -1 + 0$$

$$= 5 //$$

Performance:

→ N-w → approach → n^2 //

→ accessing all possible combinations → $2^n C_n$

$$N^2 < 2^n C_n //$$

N-w faster than access all possible elements one by

one.