



## 1 Sequence alignment

- Basics
- Importance
- Varieties of sequence alignment

## 2 Global alignment

- Finding the best global alignment score
- Deriving the optimal global alignment
- Alignment with gap penalty functions
- Alignment using hidden Markov models

## 3 Semi-global alignment

## 4 Local alignment

- Finding the best local alignment score
- Deriving the optimal local alignment

## 5 Multiple alignment

- Scoring the alignment of multiple sequences
- Deriving the best multiple alignment

## 6 Hands-on





# Varieties of sequence alignment

We are often interested in finding the best alignments between two sequences. This can be categorized into different types of problems as listed below.

- **Global alignment:** It refers to the alignment of the entire sequence pairs.
- **Semi-global alignment:** It refers to the alignment of prefixes and suffixes of the given pair of sequences not any arbitrary substrings.
- **Local alignment:** It refers to the alignment of just the substrings of a pair of sequences.
- **Multiple alignment:** It refers to the alignment of multiple (more than a pair) sequences.





# Deriving similarity of sequences

We can store the different possible alignment scores in an array 'a' of dimension  $(m + 1) \times (n + 1)$  by aligning the sequence prefixes.

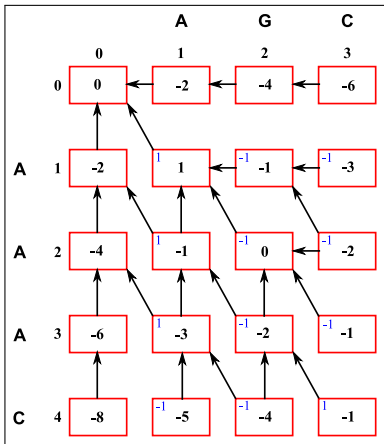
The value for an entry  $(i, j)$ , represented as  $a[i, j]$ , can be computed by looking at just three previous entries: those for  $(i - 1, j)$ ,  $(i - 1, j - 1)$ , and  $(i, j - 1)$ . The reason is that there are just three ways of obtaining an alignment between  $s[1 \dots i]$  and  $t[1 \dots j]$ , and each one uses one of these previous values.

In fact, to get an alignment for  $s[1 \dots i]$  and  $t[1 \dots j]$ , we have the following three choices:

- 1 Align  $s[1 \dots i]$  with  $t[1 \dots j - 1]$  and match a gap with  $t[j]$ .
- 2 Align  $s[1 \dots i - 1]$  with  $t[1 \dots j - 1]$  and match  $s[i]$  with  $t[j]$ .
- 3 Align  $s[1 \dots i - 1]$  with  $t[1 \dots j]$  and match  $s[i]$  with a gap.



# Deriving similarity of sequences



Deriving optimal global alignment between AAAC and AGC (the corner of a cell  $(i,j)$  reflects whether  $s[i] = t[j]$ )





# Improving the space complexity

It is possible to improve the space complexity from quadratic to linear and keep the same generality as follows.

**Input:** Sequences  $s$  and  $t$ .

**Output:** Matrix  $a$  containing the similarities between  $s$  and  $t$ .

```

1:  $m \leftarrow |s|$  // Length of  $s$ 
2:  $n \leftarrow |t|$  // Length of  $t$ 
3: for  $j \leftarrow 0$  to  $n$  do
4:    $a[j] \leftarrow j \times g$  // Filling up the  $j^{\text{th}}$  row
5: end for
6: for  $i \leftarrow 1$  to  $m$  do
7:    $old \leftarrow a[0]$ 
8:    $a[0] \leftarrow i \times g$ 
9:   for  $j \leftarrow 1$  to  $n$  do
10:     $temp \leftarrow a[j]$ 
11:     $a[j] \leftarrow \max(a[j] + g, old + p(i, j), a[j - 1] + g)$ 
12:     $old \leftarrow temp$ 
13:  end for
14: end for
15: return  $a[m, n]$ 

```

# Deriving the optimal global alignment

**Input:** Indices  $i, j$ , and the array  $a$  given by the previous algorithm.

**Output:** Alignments in  $align-s$ ,  $align-t$ , and length in  $len$ .

```

1: if  $i = 0$  and  $j = 0$  then
2:    $len \leftarrow 0$ 
3: else
4:   if  $i > 0$  and  $a[i, j] = a[i - 1, j] + g$  then
5:     Recursive-call( $i - 1, j, len$ )
6:      $len \leftarrow len + 1$ 
7:     Set  $align-s[len] \leftarrow s[i]$  and  $align-t[len] \leftarrow -$ 
8:   else
9:     if  $i > 0$  and  $j > 0$  and  $a[i, j] = a[i - 1, j - 1] + p(i, j)$  then
10:      Recursive-call( $i - 1, j - 1, len$ )
11:       $len \leftarrow len + 1$ 
12:      Set  $align-s[len] \leftarrow s[i]$  and  $align-t[len] \leftarrow t[j]$ 
13:    else
14:      Recursive-call( $i, j - 1, len$ )
15:       $len \leftarrow len + 1$ 
16:      Set  $align-s[len] \leftarrow -$  and  $align-t[len] \leftarrow t[j]$ 
17:    end if
18:  end if
19: end if

```

# Deriving the optimal global alignment

The optimal global alignment for the example shown earlier can be derived using this algorithm as follows.

**Step 1:** Start from  $a[m = 4, n = 3]$  and align C with C.

**Step 2:** Move diagonally to  $a[m = 3, n = 2]$  and align A with G.

**Step 3:** Move up to  $a[m = 2, n = 2]$  and align A with a gap (-).

**Step 4:** Move diagonally to  $a[m = 1, n = 1]$  and align A with A.

**Step 5:** Move diagonally to  $a[m = 0, n = 0]$  and stop.

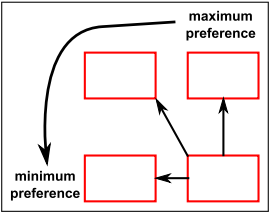
Thus, the final alignment becomes the following:

**AAAC**

**A-GC**

# Deriving the optimal global alignment

Note that, many optimal alignments may exist for a given pair of sequences. The presented algorithm returns just one of them, giving preference to the edges leaving  $(i, j)$  in counterclockwise order as shown below.



So, when there is choice, a column with a gap in  $t$  has precedence over a column with two symbols, which in turn has precedence over a column with a gap in  $s$ . This can be changed by changing the order of the *if-else* blocks in the algorithm.





# Alignment with gap penalty functions

In this algorithm, we cannot break an alignment in two parts and expect the total score to be the sum of the partial scores. However, score additivity is still valid if we break the alignment in block boundaries.

Every alignment can be uniquely decomposed into a number of consecutive blocks. There can be three kinds of blocks as listed below.

- 1 Two aligned characters from the alphabet set.
- 2 A maximal series of consecutive characters in  $t$  aligned with spaces in  $s$ .
- 3 A maximal series of consecutive characters in  $s$  aligned with spaces in  $t$ .



# Alignment with gap penalty functions

To compare sequence  $s$  of length  $m$  to sequence  $t$  of length  $n$ , we use three arrays of size  $(m + 1) \times (n + 1)$ , one for each type of ending block. Array  $a$  is used for alignments ending in character-character blocks;  $b$  is used for alignments ending in spaces in  $s$  and  $c$  is used for alignments ending with spaces in  $t$ .

The initialization is done as follows:

$$a[0, 0] = 0$$

$$b[0, j] = f(j)$$

$$c[i, 0] = f(i)$$

# Alignment with gap penalty functions

The following recurrence relations are used for updating the cells:

$$a[i, j] = p(i, j) + \max \begin{cases} a[i - 1, j - 1] \\ b[i - 1, j - 1] \\ c[i - 1, j - 1] \end{cases}$$

$$b[i, j] = \max \begin{cases} a[i, j - k] + f(k), 1 \leq k \leq j \\ c[i, j - k] + f(k), 1 \leq k \leq j \end{cases}$$

$$c[i, j] = \max \begin{cases} a[i - k, j] + f(k), 1 \leq k \leq i \\ b[i - k, j] + f(k), 1 \leq k \leq i \end{cases}$$



# Introduction to Hidden Markov models

Processes are of two types – deterministic (e.g., rolling the SHOLAY coin) and stochastic (e.g., rolling a dice).

## Definition

A stochastic process  $X = (X_t : t \in I)$  on a probability space  $(\Omega, \mathcal{F}, \mathbb{P})$  is said to process the Markov property if,  $\forall A \in \mathcal{F}$  and  $s, t \in I, s < t$ , we have

$$\mathbb{P}(X_t \in A | \mathcal{F}_s) = \mathbb{P}(X_t \in A | \sigma(X_s)),$$

where  $\{\mathcal{F}_t\}_{t \in I}$  is the natural filtration.

If the process takes discrete values and is indexed by a discrete time, this can be reformulated as follows.

$$\mathbb{P}(X_n = x_n | X_{n-1} = x_{n-1} \cdots X_0 = x_0) = \mathbb{P}(X_n = x_n | X_{n-1} = x_{n-1}).$$





# Basics

Using the already adopted scoring scheme, the first alignment turns out to be better, however, the second one appears to be more appropriate due to its *continuity*. This continuity can be effectively quantified based on the count of end gaps.

To obtain the score of the optimal alignment between  $s$  and  $t$  without penalizing the gaps after the end of  $s$  (final gaps), all we need to do is to find the best similarity between  $s$  and a prefix of  $t$ . In the previous algorithms, the entry  $(i, j)$  of matrix  $a$  contains the similarity between  $s[1 \dots i]$  and  $t[1 \dots j]$ . Therefore, it suffices to take the maximum value in the last row (or column) of the array, i.e.,

$$\mathcal{S}(s, t) = \max_{j=1}^n a[m, j].$$

**Note:** Here  $\mathcal{S}(s, t)$  indicates the similarity score ignoring the end gaps.







## Basic constructions

Following initialization, the array can be filled in the usual way, with  $a[i, j]$  depending on the value of three previously computed entries as shown below.

$$a[i, j] = \max \begin{cases} a[i, j - 1] + g \\ a[i - 1, j - 1] + p(i, j) \\ a[i - 1, j] + g \\ 0 \end{cases}$$

For any entry  $(i, j)$ , there is always the alignment between the empty suffixes of  $s[1 \dots i]$  and  $t[1 \dots j]$ , which has the score zero. Therefore, the array will have non-negative entries only.

























# Pairwise projection of multiple alignments

Consider the following multiple sequence alignment.

**ADNMQPHLLL-**

**ADNMLR-LL-Y**

**ADNMK--LLLY**

**-DNMPPVLHLY**

Suppose we take only the second and third one out of the above.

**ADNMLR-LL-Y**

**ADNMK--LLLY**

By removing additional spaces, the induced pairwise alignment (projection) is obtained as follows.

**ADNMLRLL-Y**

**ADNMK-LLLY**







