

Computational Molecular Biology and Bioinformatics

Sequence Alignment

Malay Bhattacharyya

Assistant Professor

Machine Intelligence Unit
Indian Statistical Institute, Kolkata

October, 2021

Basics

Consider the two DNA sequences given by GACGGATTAG and GATCGGAATAG. It is obvious that we can align them one above the other as follows.

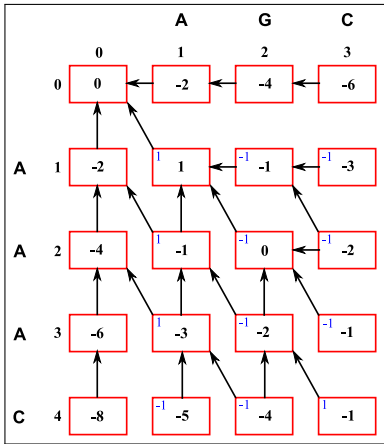
GA-CGGATTAG

GATCGGAATAG

Note that the only two differences that are distinguishable in the above alignment are given below.

- 1 There appears an extra T in the second sequence (gap), and
- 2 There is a change from A to T in the fourth position from right to left (mismatch).

Deriving similarity of sequences



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

Finding the best global alignment score

The following dynamic programming approach (known as Needleman-Wunsch algorithm) can be used to compute the best alignment score (which is -1 here) for the previous example.

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  $i \leftarrow 0$  to  $m$  do
4:    $a[i, 0] \leftarrow i \times g$  // Filling up the first column
5: end for
6: for  $j \leftarrow 0$  to  $n$  do
7:    $a[0, j] \leftarrow j \times g$  // Filling up the first row
8: end for
9: for  $i \leftarrow 1$  to  $m$  do
10:  for  $j \leftarrow 1$  to  $n$  do
11:     $a[i, j] \leftarrow \max(a[i - 1, j] + g, a[i - 1, j - 1] + p(i, j), a[i, j - 1] + g)$ 
12:  end for
13: end for
14: return  $a[m, n]$ 

```

Complexity analysis

Time complexity:

This algorithm consists of four *for* blocks. The first two blocks (initialization steps) consume time $O(m)$ and $O(n)$, respectively. The last two nested *for* blocks are used to fill the rest of the matrix. The number of operations performed depends essentially on the number of entries that must be computed, that is, the size of the matrix. Thus, we spend time $O(mn)$ in this part and the time complexity becomes

$$O(m) + O(n) + O(mn) = O(mn).$$

Space complexity:

As we need to fill the entries of the matrix 'a', the space complexity becomes $O(mn)$.

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:    $\text{len} \leftarrow 0$ 
3: else
4:   if  $i > 0$  and  $a[i, j] = a[i - 1, j] + g$  then
5:     Recursive-call( $i - 1, j, \text{len}$ )
6:      $\text{len} \leftarrow \text{len} + 1$ 
7:     Set  $\text{align-s}[\text{len}] \leftarrow s[i]$  and  $\text{align-t}[\text{len}] \leftarrow -$ 
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( $i - 1, j - 1, \text{len}$ )
11:       $\text{len} \leftarrow \text{len} + 1$ 
12:      Set  $\text{align-s}[\text{len}] \leftarrow s[i]$  and  $\text{align-t}[\text{len}] \leftarrow t[j]$ 
13:    else
14:      Recursive-call( $i, j - 1, \text{len}$ )
15:       $\text{len} \leftarrow \text{len} + 1$ 
16:      Set  $\text{align-s}[\text{len}] \leftarrow -$  and  $\text{align-t}[\text{len}] \leftarrow t[j]$ 
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.

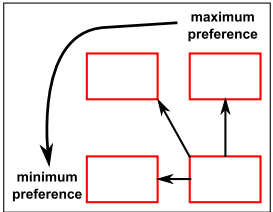
Thus, the final alignment becomes the following:

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

Let us redefine a *gap* as a consecutive number of $k > 1$ spaces. The formation of such (consecutive) gaps with k spaces is more probable than k isolated spaces during mutations.

As of now, no distinction has been made between the consecutive and isolated gaps. The gaps were penalized in the previous cases through a linear function given by

$$f(k) = kg,$$

where g is the score associated with a single space and k is the number of spaces.

We introduce an algorithm that computes similarities with respect to general gap penalty functions that consider non-additive scores.

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:

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

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

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

Complexity analysis

Time complexity:

For computing $a[i, j]$, $b[i, j]$, and $c[i, j]$, we need to perform $(3 + 2j + 2i)$ accesses. So, the total worst case time complexity becomes

$$\begin{aligned}
 & \sum_{i=1}^m \sum_{j=1}^n (3 + 2j + 2i) \\
 = & \sum_{i=1}^m (2ni + n^2 + 4n) \\
 = & O(mn^2 + m^2n).
 \end{aligned}$$

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

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

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

Alignment using hidden Markov models

We can prepare a matrix of transitional probabilities (where an entry at cell (i, j) denotes the probability of occurrence of j immediately after i), from the given sequences. E.g., consider the sequence **AAGGAATTAGC** and the corresponding transitional probabilities as shown below.

	-	A	T	C	G
-	-	1	0	0	0
A	0	2/5	1/5	0	2/5
T	0	1/2	1/2	0	0
C	1	0	0	0	0
G	0	1/3	0	1/3	1/3

The alignment can be generated from these transition probabilities. We use the relation $\arg \max_{i \in \{A, T, C, G\}} P(X_n = i | X_{n-1} = x_{n-1})$ to derive the aligned sequence.

Basics

In a semi-global comparison, we score alignments ignoring some of the end gaps (that appear before the first or after the last character) in the sequences. With a slight modification to the already presented algorithms, we can control the penalty associated with end gaps.

The end gaps are welcome because they might provide more acceptable alignments. E.g., consider the sequences AGCACTTGGATTCTCGG and CAGCGTGG, and their following two possible alignments.

-AGCACTTGGATTCTCGG
CAGC-----G-T-----GG

[Match = 7, Mismatch = 0, Gap = 11]

AGCA-CTTGGATTCTCGG
---CAGCGTGG-----

[Match = 6, Mismatch = 1, Gap = 11]

Basics

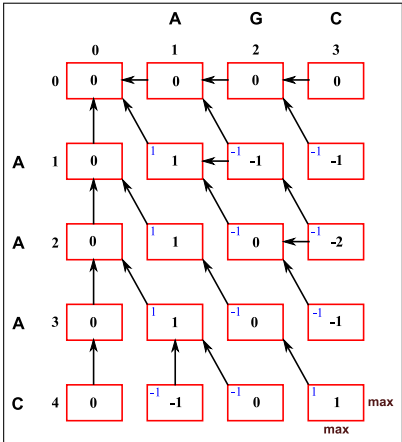
Using the already adopted scoring scheme, the first alignment turns out to be better, however, the second one appears to be more appropriate due to its *continuity*. This continuity can be effectively quantified based on the count of end gaps.

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.

Deriving similarity of sequences



Deriving optimal semi-global alignment between AAAC and AGC by ignoring the end spaces at the beginning and end

Deriving the best semi-global alignment

The maximum similarity scores over the rows (or columns) basically gives the score of the best alignment. To recover the alignment itself, we proceed just as in the previous algorithms, but starting at (k, m) (or (m, k) for columns) where k is such that $S(s, t) = a[k, m]$.

However, the initializations will be different based on the different versions of the same problem as follows.

Where gaps are not charged	Action
Beginning of first sequence	Initialize first row with zeros
End of first sequence	Look for maximum in last row
Beginning of second sequence	Initialize first column with zeros
End of second sequence	Look for maximum in last column

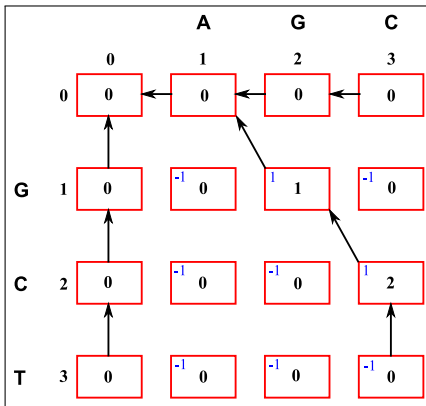
Basic constructions

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.

Deriving similarity of sequences



Deriving optimal local alignment between the sequences
AGC and GCT

Finding the best local alignment score

The following dynamic programming approach (known as Smith-Waterman algorithm) can be used for the previous example.

Input: Sequences s and t .

Output: Similarity between s and t .

```

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

```


Complexity analysis

Time complexity:

The algorithm for finding the best local alignment score consumes time $O(mn)$ and the algorithm for deriving the optimal local alignment consumes $O(m + n)$.

Space complexity:

For both the above cases, space complexity becomes $O(mn)$.

Basics

We are often required to align more than two sequences simultaneously in the best possible way. This refers to the problem of multiple sequence alignment.

Given a set of sequences over the same alphabet, a multiple alignment is obtained by inserting gaps in these sequences such that their sizes become the same. We generally place the extended sequences in a vertical list so that characters or gaps in the corresponding positions occupy the same column.

Such an example with protein sequences is shown below.

```

ADNMQPHLLL-
ADNMLR-LL-Y
ADNMK---LLLY
-DNMPPVLHLY

```

Scoring the alignment of multiple sequences

As because scoring a multiple alignment is more complex than its pairwise counterpart, we therefore restrict ourselves to purely additive functions here, i.e., the alignment score is the sum of column scores.

We note the following two important requirements for deriving such a score.

- 1 The function must be independent of the order of arguments.
- 2 The function should reward the presence of many equal or strongly related characters and penalizes unrelated residues and gaps.

We use the sum-of-pairs (SP) function for this purpose.

Deriving the best multiple alignment

Suppose, for simplicity, that we have k sequences, all of the same length n . We use a k -dimensional array a of length $n + 1$ in each dimension to hold the optimal scores for multiple alignments of prefixes of the sequences.

After initializing with $a[0, \dots, 0] \leftarrow 0$, we fill in this entire array by computing $a[i] \leftarrow \max_{b \neq 0} (a[i - b] + \text{SP-score}(\text{Column}(s, i, b)))$, where b ranges over all nonzero binary vectors of k elements. Here, $\text{Column}(s, i, b) = (c_j)_{1 \leq j \leq k}$ with $c_j = s_j[i_j]$, if $b_j = 1$, and $c_j = -$, otherwise.

Note: The cell $a[i_1, \dots, i_k]$ holds the score of the optimal alignment involving $s_1[1 \dots i_1], \dots, s_k[1 \dots i_k]$.

Complexity analysis

Time complexity:

This algorithm works on every cell of the array a to compute the values that consumes $O(n^k)$ time. Again, for each entry this computation depends on $2^k - 1$ entries, thereby requiring $O(2^k)$ time. Additionally, the algorithm uses SP-score for scoring the alignments by computing pairwise alignments. This consumes $O(k^2)$ time. Therefore, the total worst case complexity becomes $O(k^2 2^k n^k)$, where k is the number of sequences with length n .

Space complexity:

As we need to fill the entries of the matrix a , the space complexity becomes $O(n^k)$.

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

```

Relevance test

Theorem

A cell is relevant when each of its pairwise projections is part of an optimal alignment of the two sequences corresponding to the projection. Let α be an optimal alignment involving s_1, \dots, s_k . If $SP\text{-score}(\alpha) \geq L$, then

$$\text{score}(\alpha_{ij}) \geq L_{ij},$$

where $L_{ij} = L - \sum_{x < y, (x,y) \neq (i,j)} (\mathcal{S}(s_x, s_y))$.

Relevance test

Proof.

The relation $\text{SP-score}(\alpha) \geq L$ can be written as

$$\sum_{x < y} (\text{score}(\alpha_{xy})) \geq L$$

$$\Rightarrow \sum_{x < y, (x,y) \neq (i,j)} (\text{score}(\alpha_{xy})) \geq L - \text{score}(\alpha_{ij})$$

$$\Rightarrow \sum_{x < y, (x,y) \neq (i,j)} (\mathcal{S}(s_x, s_y)) \geq L - \text{score}(\alpha_{ij})$$

$$\Rightarrow \text{score}(\alpha_{ij}) \geq L - \sum_{x < y, (x,y) \neq (i,j)} (\mathcal{S}(s_x, s_y)).$$



Hands-on

- 1 BLASTn (Nucleotide BLAST) the reference genome sequence of SARS-CoV-2 using the NCBI database. BLAST is available at: <https://blast.ncbi.nlm.nih.gov/Blast.cgi>
- 2 Download the following paper and do the following:
Li, H., Minimap2: pairwise alignment for nucleotide sequences. *Bioinformatics*, 34(18), pp.3094-3100, 2018.
 - i) Get the implementation from:
<https://github.com/lh3/minimap2>
 - ii) Apply it on some supplementary available at:
<https://academic.oup.com/bioinformatics/article/34/18/3094/4994778>
 - iii) Could you identify some limitations of this implementation? How can you overcome that? Any suggestions?