

Lecture 2: December 13, 1998

*Lecturer: Ron Shamir**Scribe: Eti Ezra, Guy Gelles¹*

2.1 Pairwise Alignment

The lecture introduces the problem of comparing two strings while allowing certain mismatches between the two. The problem will be referred to as *(pairwise) sequence alignment* or *inexact matching*.

We first present the problem and provide biological motivation. We then define similarity and difference between strings and present algorithms for computing them, with analysis of the complexity of these algorithms.

The algorithms use the *dynamic programming* technique. For each algorithm the following information will be given:

- Intuitive explanation of the recursive process.
- Formal definition of the recursive process.
- Discussion of the complexity.

2.1.1 Problem Definition and Biological Motivation

Motivation

A large variety of the biologically motivated problems in computer science primarily involve sequences or strings. For instance:

- Reconstructing long sequences of DNA from overlapping string fragments.
- Determining physical and genetic maps from probe data under various experiments protocols.
- Storing, retrieving and comparing DNA strings.
- Comparing two or more strings for similarities.
- Searching databases for related strings and substrings.

¹Partly based on scribe by Amichay Oren November 6th, 1995

- Exploring frequently occurring patterns of nucleotides.
- Finding informative elements in protein and DNA sequences.

Many of these research problems aim at learning about functionality or the structure of protein without performing any experiments and actually without having to physically construct the protein itself. The basic idea is that similar sequences produce similar proteins. Thus, in order to predict the characteristics of a protein using only its sequence data, we can use the structure/function information on known protein with similar sequences available in databases.

For instance, when considering protein folding, it usually suffices that two protein sequences are identical at 25% of their positions for their three dimensional structures to be almost identical. Classical example is the establishment of an association between cancer and uncontrolled cells growth [1]. This discovery was enabled by comparing the sequence of a cancer associated gene against the sequence of protein which had already been known to be influed cell growth. The correlation between these two sequences was very high, proving the connection between cancer and cellular growth.

2.1.2 Similarity and Difference

The resemblance of two DNA sequences taken from different organisms can be explained by the theory that all contemporary genetic material has one ancestral ancient DNA. According to this theory, during the course of evolution mutations occurred, creating differences between families of contemporary species. Most of these changes are due to local mutations, each modifying the DNA sequence at a specific manner. These local modifications between nucleotide sequences, or more generally, between strings over an arbitrary alphabet can be either:

- *Insertion* - an insertion of a letter or several letters to the sequence.
- *Deletion* - deleting a letter (or more) from the sequence.
- *Substitution* - replacing a sequence letter by another.

Insertion and deletion are the reverse of one another: given two sequences, if the insertion of a character (or more) into one yields the other, then equivalently its deletion from the latter sequence transforms it to the first one. Due to this reciprocity between insertion and deletion, they are usually called *indel* for short.

The notion of *distance* derives its definition from the concept of mutations: by assigning weights to each mutation. Given two strings, the *distance* between them is the minimal sum of weights for a set of mutations transforming one into the other.

The notion of *similarity* derives its definition from the concept of one ancestral ancient DNA: by assigning weights corresponding for resemblance. Given two strings the *similarity* between them is the maximal sum of such weights.

Models for Inexact Matching

In this lecture we consider four variants of biologically motivated inexact matching problems:

Problem 2.1 Global Alignment

INPUT: *Two strings S and T of roughly the same length.*

QUESTION: *What is the difference (or similarity) between the two?*

Problem 2.2 Local Alignment

INPUT: *Two strings S and T .*

QUESTION: *What is the maximum similarity (minimum difference) between a substring of S and a substring of T ? What are these most similar substrings?*

Problem 2.3 Ends free-space alignment

INPUT: *Two strings S and T of different length.*

QUESTION: *What is the maximum similarity between substrings of S and T , respectively? where at least one of these substrings must be a prefix of the original string and one (not necessary by the other) must be a suffix?*

Problem 2.4 Gap penalty

INPUT: *Two strings S and T of different length.*

QUESTION: *Defining a gap as any maximal, consecutive run of spaces in a single string of a given alignment, and the length of a gap as the number of indel operations on it², What is the similarity between the two strings, given a weight function for gaps (depending on their lengths)?*

2.1.3 Global Alignment

Definition Informally, an *alignment* of two strings S and T is obtained by first inserting chosen spaces, either into or at the ends of S and T so the length of the strings will match, and then placing the two resulting strings one above the other so that every character or space in one of the strings is matched to a unique character or a unique space in the other string.

²Gap Penalty

Example 2.5 Given a string "ACBCDDDB" and a string "CADBDAD", one possible alignment will be:

```

A C - - B C D D D B
  |       | | | |
- C A D B - D A D -

```

A more useful than the general case is the following problem:

Problem 2.6 INPUT: Two strings S and T $|T| = m$, $|S| = n$ (n and m are of roughly the same magnitude)

QUESTION: Establish the optimal alignment according to the alignment quality (or scoring) which will be defined next.

Notation Let $\sigma(a, b)$ be the score (weight) of the alignment of character a with character b . (including spaces) **Notation** Let $V(i, j)$ be the optimal score of the alignment of $S_1 \dots S_i$ and $T_1 \dots T_j$ ($0 \leq i \leq n, 0 \leq j \leq m$).

Lemma 2.7 $V(A, B)$ has the following properties:

$$\begin{aligned}
 \text{Base conditions :} \quad & V(i, 0) = \sum_{k=0}^i \sigma(S_k, -) \\
 & V(0, j) = \sum_{k=0}^j \sigma(-, T_k) \\
 \text{Recurrence relation :} \quad & \text{for } 1 \leq i \leq n, 1 \leq j \leq m : \\
 & V(i, j) = \max \begin{cases} V(i-1, j-1) + \sigma(S_i, T_j) \\ V(i-1, j) + \sigma(S_i, -) \\ V(i, j-1) + \sigma(-, T_j) \end{cases}
 \end{aligned}$$

Proof: **Base condition:**

The only way to align the first i elements of the string S with zero elements of the string T is to align each of the elements with a space in the string T . The score for that operation is by definition $\sigma(S_i, -)$ for each of the i elements and $V(i, 0) = \sum_{k=0}^i \sigma(S_k, -)$ for the total sum.

Similarly, the expression $V(0, j) = \sum_{k=0}^j \sigma(-, T_k)$ follows from matching the first j elements of T with i blanks in string S . ■

Proof: **Recurrence relation:**

Let us consider an optimal alignment of $S_1 \dots S_i$ and $T_1 \dots T_j$. We shall distinguish between three cases according to the three possible scoring for the three operations are:

- **Aligning S_i with T_j :** The score in this case is the score of the operation $\sigma(S_i, T_j)$ plus the score of aligning $i - 1$ elements of S with $j - 1$ elements of T , namely, $V(i - 1, j - 1) + \sigma(S_i, T_j)$
- **Aligning S_i with a space character in string T :** The score in this case is the score of the indel operation $\sigma(S_i, -)$ plus the score of aligning the previous $i - 1$ elements of S with j elements of T (Since the space is not an original character of T), $V(i - 1, j) + \sigma(S_i, -)$
- **Aligning T_j with a space character in string S :** Similar to the previous case, the score will be $V(i, j - 1) + \sigma(-, T_j)$

■

Tabular computation of optimal alignment

The problem can be evaluated systematically using a tabular computation. In this approach, we compute $V(i, j)$ for the all possible values for i and j starting from smaller such values and increasing them in a row-wise manner.

We use a table of size $(n + 1) \times (m + 1)$ in which we store the values of $V(i, j)$ for all choices of i and j . Finally $V(n, m)$ is the required alignment score.

The following pseudo code describes the algorithm:

```

for i=1 to n do
begin
  for j=1 to m do
  begin
    Calculate V(i,j) using V(i-1,j-1), V(i,j-1), V(i-1,j)
  end
end
end

```

Example 2.8 *Figure 2.1 illustrates a snapshot at some point during the computation. In this example and the following one the value of σ is -1 for a mismatch and 2 for a match.*

The Traceback

One way to traceback the alignments is to establish pointers in the cells of the table as the values are computed. The direction of the pointers in cell (i, j) indicates which cell contributed the most to $V(i, j)$.

Theorem 2.9 *The time complexity of the algorithm is $O(nm)$. Space complexity is $O(n + m)$, if only $V(S, T)$ is required and $O(mn)$ for the reconstruction of the alignment.*

		j						
		0	1	2	3	4	5	
i			c	a	d	b	d	←T
	0	0	-1	-2	-3	-4	-5	
1	a	-1	-1	1				
2	c	-2						
3	b	-3						
4	c	-4						
5	d	-5						
6	b	-6						

↑
S

Figure 2.1: Snapshot of computing the table

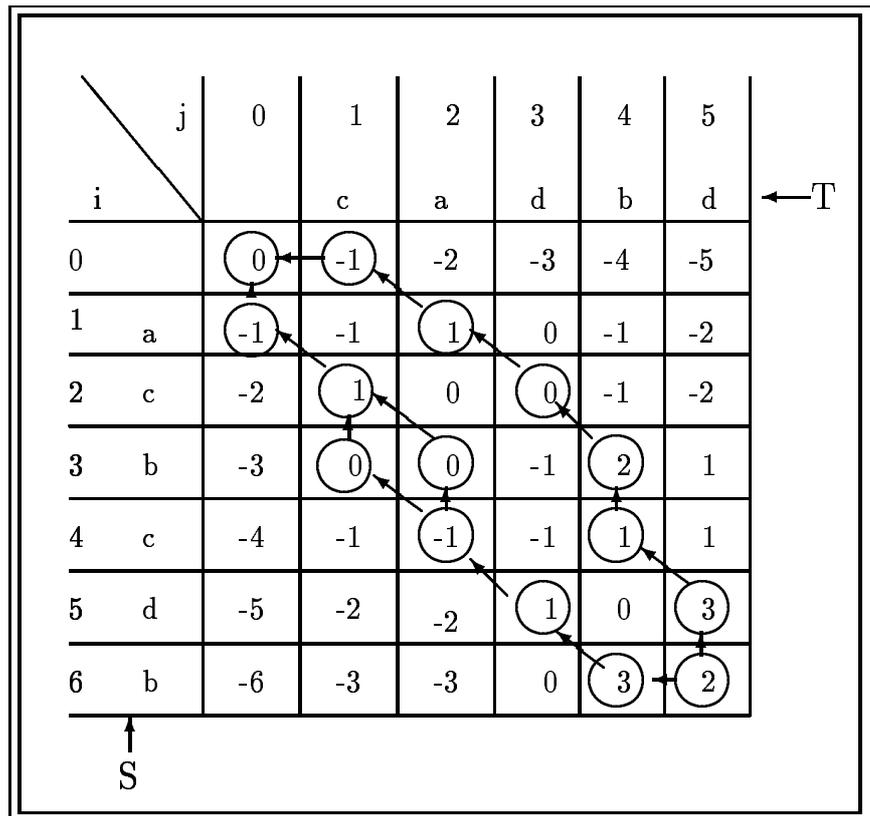


Figure 2.2: Backtracing the alignment

Proof:

- Time complexity - When computing the value for a specific cell (i, j) only cells $(i - 1, j - 1)$, $(i, j - 1)$ and $(i - 1, j)$ are examined, along with two characters S_i and T_j . Hence, to fill in one cell takes a constant number of cell examinations and comparisons. There are $n \times m$ cells in the table. So the time complexity is also $O(nm)$.
- Space complexity - Using the algorithm, computing the value of cell (i, j) involves one cell in row j $(i - 1, j)$ and two cells in the previous row $((i - 1, j - 1)$ and $(i, j - 1))$. Since the computation is performed one row at a time, when computing the values in row k only row $k - 1$ has to be stored, using, $O(n + m)$ space. In order to reconstruct the alignment from the recursion, pointers must be set for allowing the back-tracing, therefore, the space complexity is $O(nm)$. ■

2.1.4 Global Alignment in linear space

The backtracing algorithm requires the entire matrix to be saved in memory. The space complexity consequently increase to $O(nm)$. Hirschberg [4] developed a more practical space-reduction method for solving dynamic programming problems that reduces the required space from $O(nm)$ to $O(m)$ (for $m < n$).

Notation $V^r(i, j)$ will denote an optimal alignment value of last i characters in string S against last j characters in string T .

Lemma 2.10

$$V(n, m) = \max_{0 \leq k \leq m} \{V(\frac{n}{2}, k) + V^r(\frac{n}{2}, m - k)\}$$

Proof: [3] (chapter 12) For any fixed position k' in T , there is an alignment of S and T consisting of an alignment of $S_1 \dots S_{\frac{n}{2}}$ and $T_1 \dots T_{k'}$ followed by a disjoint alignment of $S_{\frac{n}{2}+1} \dots S_n$ and $T_{k'+1} \dots T_m$. By definition of V and V^r , the best alignment of the first type has value $V(\frac{n}{2}, k')$ and the best alignment of the second type has value $V^r(\frac{n}{2}, m - k')$, so the combined alignment has value $V(\frac{n}{2}, k') + V^r(\frac{n}{2}, m - k') \leq V(n, m)$. since this argument holds for any k' , it follows that $\max_k [V(\frac{n}{2}, k) + V^r(\frac{n}{2}, m - k)] \leq V(n, m)$.

Conversely, for an optimal alignment of S and T , let k' be the right-most position in T that is aligned with a character at or before position $\frac{n}{2}$ in S . Then the optimal alignment of S and T consists of an alignment of $S_1 \dots S_{\frac{n}{2}}$ and $T_1 \dots T_{k'}$ followed by an alignment of $S_{\frac{n}{2}+1} \dots T_n$ and $T_{k'+1} \dots T_m$. Let the value of the first alignment be denoted p and the value of the second alignment be denoted q . Then p must be equal to $V(\frac{n}{2}, k')$, for if $p < V(\frac{n}{2}, k')$ we could replace the alignment of $S_1 \dots S_{\frac{n}{2}}$ and $T_1 \dots T_{k'}$ with the alignment of $S_1 \dots S_{\frac{n}{2}}$

and $T_1 \dots T_{k'}$ that allows value $V(\frac{n}{2}, k')$. That would create an alignment of S and T whose value is larger than the claimed optimal. Hence, $p = V(\frac{n}{2}, k')$. By similar reasoning, $q = V^r(\frac{n}{2}, m - k')$. So $V(n, m) = V(\frac{n}{2}, k') + V^r(\frac{n}{2}, m - k') \leq \max_k [V(\frac{n}{2}, k) + V^r(\frac{n}{2}, m - k)]$. Having shown both sides of the inequality, we conclude that $V(n, m) = \max_k [V(\frac{n}{2}, k) + V^r(\frac{n}{2}, m - k)]$. ■

The Algorithm:

1. Compute $V(A, B)$ while saving the values of the $\frac{n}{2}$ -th row. Denote $D(A, B)$ as the *Forward Matrix* F.
2. Compute $V(A^r, B^r)$ while saving the $\frac{n}{2}$ -th row. Denote $D(A^r, B^r)$ as the *Backward Matrix* B.
3. Find the column k^* so that the crossing point $(\frac{n}{2}, k^*)$ satisfies:

$$F(\frac{n}{2}, k^*) + B(\frac{n}{2}, m - k^*) = F(n, m)$$

4. Now that k^* is found, recursively partition the problem to two sub problems:
 - (i) Find the path from $(0,0)$ to $(\frac{n}{2}, k^*)$.
 - (ii) Find the path from (n, m) to $(\frac{n}{2}, m - k^*)$.

Lemma 2.11 Time complexity of Hirschberg's algorithm is $O(nm)$.

Proof: Time complexity

Let $T^*(n, m)$ = Time to find the value of a $n \times m$ problem. Let $T(n, m)$ = Time to find the path (solution) of a $n \times m$ problem.

$$T^*(n, m) = 2T(n, m) + T^*(\frac{n}{2}, k^*) + T^*(\frac{n}{2}, m - k^*) \quad (2.1)$$

$T(n, m) = cnm$, therefore $T^*(n, m) = 4T(n, m) = 4cnm$. according to 2.1: $4cnm = 2cnm + 4c(\frac{n}{2} \times k^*) + 4c(\frac{n}{2} \times (m - k^*))$. Therefore, the time complexity will remain $O(nm)$. ■

Lemma 2.12 Space complexity of Hirschberg's algorithm is $O(m)$

Proof: Space Complexity

Each Dynamic Programming computation requires storing one additional row (middle one) which can be discarded once the middle point is found. Therefore the space complexity will be $O(m)$. ■

2.1.5 Alignment Graph

It is often useful to represent dynamic programming solutions of string problems in terms of a weighted graph.

Definition Given two strings S and T of lengths n and m respectively. An *alignment graph* is a directed graph $G = (V, E)$ on $(n + 1) \times (m + 1)$ nodes, each labeled with a distinct pair (i, j) ($0 \leq i \leq n, 0 \leq j \leq m$), with the following weighted edges:

1. $((i, j), (i + 1, j))$ with weight $\sigma(S_{i+1}, -)$
2. $((i, j), (i, j + 1))$ with weight $\sigma(-, T_{j+1})$
3. $((i, j), (i + 1, j + 1))$ with weight $\sigma(S_{i+1}, T_{j+1})$

A path from node $(0, 0)$ to node (n, m) in the alignment graph corresponds to an alignment and its total weight is the alignment score. Our goal is to find the heaviest path from node $(0, 0)$ to node (n, m) .

This *alignment graph* is used to map the problem of optimal alignment into the world of graphs, opening the door for new and exciting algorithms.

2.1.6 Edit Distance

Definition The *edit distance* between two strings is defined as the minimum number of edit operations (insertions, deletions and substitutions) needed to transform the first string into the other.

Each operation is given a score (weight), usually, insert and delete (indel) operations are given the same score, and using alignment algorithms, we search for the minimal scoring (or the maximum negative scoring), representing the minimal number of edit operations.

Since most of the changes to a DNA during evolution are due to the three common local mutations: insertion, deletion and substitution, the edit distance can be used as a way to roughly measure the number of DNA replications that occurred between two DNA sequences.

2.1.7 Local Alignment

In many applications two strings may not be highly similar as a whole, but many contain regions that are highly similar. The task is to find and extract a pair of regions, one from each of the two given strings, that exhibit high similarity. This is called the *local alignment* or *local similarity* problem and is defined formally below:

Definition Given two strings S and T , *local alignment* problem is defined as the problem of finding the substrings α and β of S and T respectively, whose *similarity* (optimal global

alignment) is maximum over all such pairs of substrings.

Example 2.13 Consider the two strings:

$S = a b c x d e x$

$T = x x x c d e$

If we give each match a value of 2 and each mismatch a value of -1, then the two substrings: $\alpha = cxde$ and $\beta = c-de$ of S and T respectively have the optimal alignment.

Motivation

In many biological applications local similarity is far more meaningful than global similarity. This is particularly true when long stretches of non-coding DNA are compared, since only small regions within those strings may be related. When comparing protein sequences, local alignment is also critical because proteins from very different families often share the same structural or functional subunits, and local alignment is an appropriate tool for searching such moduls.

Computing local alignment

Given a pair of indices $i \leq n$ and $j \leq m$, the local suffix alignment problem is finding a (possibly empty) suffix α of $S_{1..i}$ and a (possibly empty) suffix β of $T_{1..j}$ such that the value of their alignment is the maximum over all values of alignments of suffixes of $S_{1..i}$ and $T_{1..j}$. We use $V(i, j)$ to denote the value of the optimal local suffix alignment for a given pair i, j of indices.

We choose the weights of the editing operations as:

$$\sigma(x, y) = \begin{cases} \geq 0 & \text{if } x, y \text{ match} \\ \leq 0 & \text{if } x, y \text{ do not match or one of them is -} \end{cases}$$

The algorithm needs to:

1. Find maximum similarity between suffixes of $S_{1..i}$ and $T_{1..j}$.
2. Discard the prefixes $S_{1..i}, T_{1..j}$ whose similarity is ≤ 0 , and therefore decreases the overall similarity.
3. Find the best indices i^*, j^* of S and T respectively after which the similarity only decreases.

Note that any extension of the optimal solution either to the right or to the left decreases the overall similarity.

Recursive definition: The base condition will be: $V(i, 0) = 0$ and $V(0, j) = 0 \forall i, j$ since we can always choose an empty suffix.

For $i > 0$ and $j > 0$ the proper recurrence for $V(i, j)$ is

$$V(i, j) = \max\{0, V(i-1, j-1) + \sigma(s_i, T_j), V(i, j-1) + \sigma(-, T_j), V(i-1, j) + \sigma(s_i, -)\}$$

Compute i^*, j^* so that:

$$V(i^*, j^*) = \max_{1 \leq i \leq n, 1 \leq j \leq m} V(i, j)$$

Observe that the recurrence for computing local suffix alignment is almost identical to the one used for computing global alignment. The only difference is the inclusion of zero in the case of local suffix alignment. In both global alignment and local suffix alignment of prefixes $S_{1\dots i}$ and $T_{1\dots j}$, the terminating characters of any alignment are specified, but in the case of local suffix alignment, any number of initial characters can be ignored.

The zero in the recurrence implements this, 'restarting' the recurrence. Adding 0 to the maximization makes sure that negative prefixes are discarded from the computation.

Adding the '0' to the constraint only handles mismatched prefixes, there's still a need to determine, when should a computation of a transformation be stopped, so that the similarity value will not decrease. Therefore, after computing the table of $V(i, j)$ values, and there's a need to search for a cell with the maximal value and ignore all table entries from that point on.

Example 2.14 *Figure 2.3 illustrates the calculation of the $n \times m$ entries table for the two strings taken σ as 2 for match and -1 for mismatch.*

As usual, pointers are created while filling in the values of the table. After cell (i^*, j^*) is found, the substrings α and β giving the optimal local alignment of S and T are found by tracing back the pointers from cell (i^*, j^*) until reaching an entry (i', j') that has value zero. Then the optimal local alignment substrings are $\alpha = S_{i' \dots i^*}$ and $\beta = T_{j' \dots j^*}$. As it seems from here, space complexity will be $O(mn)$, we will show that only $O(m)$ space is needed:

Lemma 2.15 *The optimal local alignment of two strings S and T can be computed in linear space.*

Proof: The optimal local alignment of S and T identifies substrings α and β whose global alignment has maximum value over all pairs of substrings. Hence, if α and β can be found using only linear space, then their actual alignment can be found in linear space, using Hirschberg's method for global alignment. The value of the optimal local alignment is found in cell i^*, j^* . Those indices specify the terminating points of the strings α and β . The values of each row can be computed in a row wise fasion and the algorithm must store values for only two rows at a time. Hence, the end positions (i^*, j^*) can be computed in linear space. To find the starting position of the two substrings, the algorithm can execute the reverse dynamic programing using linear space (the details are left as an exercise). ■

i \ j		0	1	2	3	4	5	6	T
			x	x	x	c	d	e	
0		0	0	0	0	0	0	0	
1	a	0	0	0	0	0	0	0	
2	b	0	0	0	0	0	0	0	
3	c	0	0	0	0	2	1	0	
4	x	0	2	2	2	1	1	0	
5	d	0	1	1	1	1	3	2	
6	e	0	0	0	0	0	2	5	
7	x		2	2	2	1	1	4	
S									

Figure 2.3: finding local alignment

Complexity :

- Time Complexity - since it takes constant number of operation per cell to compute $V(i, j)$, it takes only $O(mn)$ time to fill in the entire table. The search for $V(i^*, j^*)$ requires only $O(nm)$ time as well. Hence the total time complexity is $O(nm)$.
- Space Complexity - As shown in lemma local alignment, the space complexity is $O(m)$.

2.1.8 End free-space alignment

In this variant, any indel operations at the end or the beginning of the alignment contribute a weight of zero, no matter what weight other spaces contribute.

Example 2.16 Consider the alignment

```
S = - - c a c - d b d v l
T = l t c a b d d b - - -
```

the two leading spaces at the left end of the alignment are free, as well as the three trailing spaces at the right end.

Motivation

One example where end-spaces should be free is in the shotgun sequence assembly procedure. In this problem, one has a large set of partially overlapping substrings that come from many copies of one original but unknown DNA sequences. The problem is to use comparisons of pairs of substrings to infer the correct original string. Two random substrings from the set are unlikely to have nearby starting positions in the original string, and this is reflected by a low end-space free alignment score for those two substrings. But if two substrings do overlap in the original string, then an alignment may be between suffix of one to a prefix of the other with only a small number of spaces and mismatches. This overlap is detected by an end-space free weighted alignment with high score. Similarly, the case when one substring contains another can be detected in this way. See figure 2.4 for illustration.

When comparing two strings, it is not obvious how to place the two strings, so that the *similarity* between the two will be maximal. One possibility, denoted by the *ends free problem* is to disregard leading and trailing indel operations (in the usual similarity strategy, all indel operations reduce the similarity).

To implement this we will change the algorithm presented for the *global alignment* problem, as follows:

- Set initial conditions:

$$\begin{aligned} V(i, 0) &= 0 && \text{for } 1 \leq i \leq n \\ V(0, j) &= 0 && \text{for } 1 \leq j \leq m \end{aligned}$$

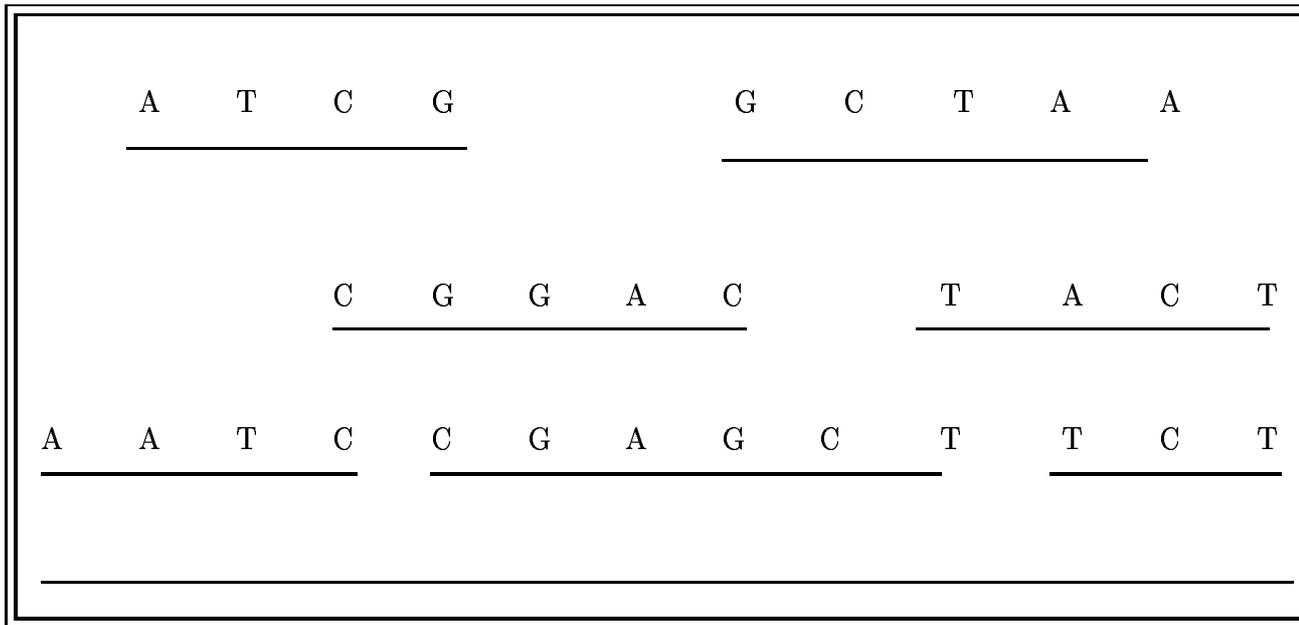


Figure 2.4: sequence assembly

- Use the same recurrence for $1 \leq i \leq n, 1 \leq j \leq m$

$$V(i, j) = \max \begin{cases} V(i-1, j-1) + \sigma(S_i, T_j) \\ V(i-1, j) + \sigma(S_i, -) \\ V(i, j-1) + \sigma(-, T_j) \end{cases}$$

- Instead for looking at $V(n, m)$ the algorithm will search for i^* and j^* so that:

$$V(n, i^*) = \max_i V(n, i)$$

$$V(j^*, m) = \max_j V(j, m)$$

- The similarity will be defined as:

$$V(S, T) = \max\{V(n, i^*), V(j^*, m)\}$$

Looking for i^* means searching for a cell in the last row of the table, produced while computing $V(n, m)$. Looking for j^* means searching for a cell in the last column of the same table. This eliminates trailing indel operations. Leading indel operations will not be taken into account due to the changes in the initial conditions.

Complexity :

- Time complexity - Computing the matrix takes $O(nm)$. Finding j^* and i^* takes $O(n + m)$. Therefore the total complexity remains $O(nm)$.
- Space complexity - Computing the matrix takes $O(n + m)$ space. Computing the maximizing values i^* , j^* requires the last row and column to be saved, which is also $O(n + m)$. Therefore the total complexity remains $O(n + m)$.

2.1.9 Gap Penalty

Until now the central constructs used to measure the value of an alignment have been matches, mismatches and spaces. Now we introduce another important construct, *gaps*. Gaps help create alignments that better conform to underlying biological models and more closely fit patterns that one expects to find in meaningful alignment. The idea is to take in account the number of continuous gaps and not only the number of spaces when calculating an alignment mark. This section presents a gap penalty model for evaluating the weight of a sequence of consecutive indel operations. The model states that consecutive indel operations have different total weight than simply the sum of their weights.

Definition A *gap* is any maximal, consecutive run of spaces in a single string of a given alignment.

Example 2.17 Consider the alignment:

```
S = a t t c - - g a - t g g a c c
T = a - - c g t g a t t - - - c c
```

which has four gaps containing a total of eight spaces. That alignment would be described as having seven matches, no mismatch, four gaps and eight spaces.

Definition The *length of the gap* will be the number of indel operations in it. The number of gaps in the alignment will be denoted as #gaps.

Motivation

The concept of a gap in an alignment is important in many biological application, because the insertion or deletion of an entire substring often occurs as single mutational event. Moreover, many of these single mutational events can create gaps of quite varying sizes. At the protein level, two protein sequences might be relatively similar over several intervals but differ in intervals where one contains a protein subunit that the other does not.

One concrete illustration of the use of gaps in the alignment model comes from the problem of cDNA matching [2] (chapter 11). In this problem, one string is much longer than the other, and the alignment best reflecting their relationship should consist of a few regions of very high similarity interspersed with 'long' gaps in the shorter string. Note that

the matching regions can have mismatches and spaces, but these should amount only to a small fraction of the region.

An RNA molecule is transcribed from DNA of the gene. That RNA transcript is a complement of the DNA in the gene in that each A in the gene is replaced by U in the RNA, each T is replaced by A, each C by G, and each G by C. Moreover, the RNA transcript covers the entire gene, introns as well as exons. Then in a process that is not completely understood, each introns-exon boundary in the transcript is located, the RNA corresponding to the introns is spliced out, and the RNA regions corresponding to exons are concatenated. Additional processing occurs. The resulting RNA molecule is called the *messenger* RNA (*mRNA*): it leaves the cell nucleus and is used to create the protein it encodes.

Each cell (usually) contains a copy of all the chromosomes and hence, of all the genes of the entire individual, yet in each specialized cell (a liver cell for example) only a small fraction of the genes are expressed. That is, only a small fraction of the proteins encoded in the genome are actually produced in that specialized cell. A standard method to determine which proteins are expressed in the specialized cell line, and to hunt for the location of the encoding genes, involves capturing the mRNA in that cell after it leaves the cell nucleus. That mRNA is then used to create a DNA sequence complementary to it. This sequence is called *cDNA* (complementary DNA). Compared to the original gene, the cDNA sequence consists only of the concatenation of exons in the gene. After cDNA is obtained, the problem is to determine where the gene associates with that cDNA resides, and it becomes one of aligning the cDNA sequence against the longer DNA sequence in a way that reveals the exons.

gap penalties types

- Constant solved in $O(nm)$ time. (see below)
- Affine solved in $O(nm)$ time. (see below)
- Convex solved in $O(nm \log m)$ time. ([2])
- Arbitrary solved in $O(nm^2 + n^2m)$ time. (left as an exercise)

Constant gap weight model

The simplest choice is the *constant* gap weight, where each individual space is free, and each gap is given a weight of W_s independent of the number of spaces in the gap. Letting σ denote the weights of match and mismatch only ($\sigma(x, -) = \sigma(-, x) = 0$ for every character x). Thus we have to find an alignment that maximizes:

$$\sum \sigma(S'_i, T'_i) + W_g \times \#gaps$$

where S' and T' represent S and T after inserting space. A generalization of the constant gap weight model is to add a weight W_s for each space in the gap. In this case, W_g represents the cost of starting a gap, and W_s represents the cost of extending the gap by one space. This leads us to the *affine gap weight model*. This is called affine gap weight model because the weight contributed by a single gap of length q is given by the affine function $W_g + qW_s$. The constant gap weight model is simply the affine model with $W_s = 0$. Thus we have to find an alignment that maximizes:

$$\Sigma\sigma(S'_i, T'_i) + W_g \times \#gaps + W_s \times \#spaces$$

while S' and T' represent S and T after inserting space and $\sigma(x, -) = \sigma(-, x) = 0$ for every character x .

It has been suggested that some biological phenomena are better modeled by a gap weight function where each additional space in a gap contributes less to the gap weight than the preceding space. In other words, a gap weight that is a *convex*, but not affine function of its length. An example is the function $W_g + \log q$, where q is the length of the gap. Finally, the most general gap weight that might be considered is the *arbitrary gap weight*, where the weight of a gap is an arbitrary function $\omega(q)$ of its length q . The constant, affine and convex weight models are ofcourse restricted cases of the arbitrary weight model.

Affine gaps penalty

To align strings S, T , consider as usual the prefixes $S_{1..i}$ of S and $T_{1..j}$ of T . Any alignment of these two prefixes is one of the following three types:

1. S ————i
T ————j
alignment of $S_{1..i}$ and $T_{1..j}$ where characters $S(i)$ and $T(j)$ are aligned opposite each other. This includes both the case that $S_i = T_j$ and that $S_i \neq T_j$.

2. S ————i_ ————
T —————j

alignment of $S_{1..i}$ and $T_{1..j}$ where character S_i is aligned to a character strictly to the left of character T_j . Therefore, the alignment ends with a gap in S.

3. S —————i
T ————j_ ————

alignment of $S_{1..i}$ and $T_{1..j}$ where character S_i is aligned to a character strickly to the right of character T_j . Therefore, the alignment ends with a gap in T.

Notation We will use $G(i, j)$ to denote the maximum value of any alignment of type 1, $E(i, j)$ as the maximum value of any alignment of type 2 and $F(i, j)$ as the maximum value of any alignment of type 3. We finally define $V(i, j)$ as the maximum value of the three terms $E(i, j)$, $F(i, j)$, $G(i, j)$.

Hence the base conditions are:

$$\begin{aligned} V(i, 0) &= E(i, 0) = W_g + iW_s \\ V(0, j) &= F(0, j) = W_g + jW_s \end{aligned}$$

and the recursive computation of $V(i, j)$ will be:

$$V(i, j) = \max\{E(i, j), F(i, j), G(i, j)\}$$

while

$$\begin{aligned} G(i, j) &= V(i - 1, j - 1) + \sigma(S_i, T_j) \\ E(i, j) &= \max\{E(i, j - 1) + W_s, V(i, j - 1) + W_g + W_s\} \\ F(i, j) &= \max\{F(i - 1, j) + W_s, V(i - 1, j) + W_g + W_s\} \end{aligned}$$

The optimal value alignment is the maximum value in the n th row or m column.

Complexity

- Time complexity - As before $O(nm)$, as we only compute four matrices instead of one.
- Space complexity - There's a need to save four matrices (for E, F, G, and V respectively) during the computation. Hence, $O(nm)$ space is needed, for the trivial implementation.

2.1.10 Longest Common Subsequence

Definition *Subsequence* is defined as a subset of the characters of string S arranged in their original "relative" order. Formally: a subsequence of string S of length n is specified by a list of indices $i_1 < i_2 < i_3 < \dots < i_k$, for some $k \leq n$. The subsequence specified by this list of indices is the string $S_{i_1} S_{i_2} \dots S_{i_k}$.

Definition Given two strings S and T , a **common subsequence** is a subsequence that appears both in S and T . The **longest common subsequence problem** is to find a longest subsequence common to both S and T .

The problem of Longest Common Subsequence can be modeled and solved using the optimal alignment algorithm with the following scoring:

$$\sigma(x, y) = \begin{cases} 1 & x = y \\ 0 & x \neq y \end{cases}$$

$$\sigma(x, -) = \sigma(-, y) = 0$$

Or, directly compute $V(n, m)$ with:

$$V(i, 0) = V(0, j) = 0$$

$$V(i, j) = \max \begin{cases} V(i-1, j-1) + \sigma(S_i, T_j) \\ V(i-1, j) \\ V(i, j-1) \end{cases}$$

Each character in the string S can be align with the same character in the string T or with a space in T (in this case no substitution is done). Since the goal is to find maximum length, character matched are valued as '1', while a space match is valued '0'.

2.1.11 Polymerase Chain Reaction

It is possible to replicate substrings of a DNA sequence starting at almost any point as long as we know a small number of the nucleotides, appearing just before that point. This replication is done using a technology called *polymerase chain reaction*, which has had a tremendous impact on experimental molecular biology.

Reaction of replication in a cell starts as a consequence of initiation a replication and then by an enzyme which continues that replication (according to the reading direction of the DNA). When considering *polymerase chain reaction* we perform the initiation and let the enzyme do the rest.

Knowing few nucleotides allows one to synthesize a string that is complementary to those few nucleotides. This complementary string can be used to create a 'primer', which finds its way to the point in the long DNA sequence containing the complement of the primer. It then hybridizes (bonds) with the longer string at that point. This creates the conditions that allow the replication of part of the original string: the replication is done by exposing the DNA to an enzyme which has the ability to duplicate strings called *DNA polymerase*. After a while we get two substrings. Now we can warm what we got for a short time and the bonds will separate - and again we can apply the enzyme and do the same to the separated substrings and get more and more copies.

Bibliography

- [1] F. Alizadeh, R. Karp, L. Newberg, and D. Weisser. Simian sarcoma virus onc gene v-sis, is derived from the gene encoding a platelet-derived growth factor. *Science*, 221:275–277, 1983.
- [2] Gusfield Dan. *Algorithms on Strings, Trees, and Sequences*. Cambridge University Press, 1997.
- [3] Dan Gusfield. *Algorithms on Strings, Trees, and Sequences*. Cambridge University Press, 1997.
- [4] D. S. Hirschberg. Algorithms for the longest common subsequence problem. *J.ACM*, 24:664–675, 1977.