# Sequence Alignment with Gap Opening and Extension

In the last lecture, we learned about a dynamic programming algorithm for sequence alignment. This algorithm uses a single fixed gap penalty when aligning amino acids against gaps. In this lecture, we introduce a variant of the dynamic programming algorithm that takes into account the gaps' tendency to cluster together.

Rather than using a single score for all gaps, our algorithm will use two gap scores: the *gap opening* score $g_o$, and the *gap extension* score $g_e$. The gap opening score $g_o$ will be used for the first gap in each gap cluster; that is, the gap that immediately follows an amino acid in the aligned sequence. The gap extension score $g_e$ will be used for each gap following the opening gap. Since gaps tend to bundle together, the presence of an opening gap makes it more likely that the next character in the alignment is also a gap. Therefore, $g_o$ should be *higher* than $g_e$.

$$\downarrow g_o \quad \downarrow g_e$$

|       |   |   |   |   |   |   |   |   |
|-------|---|---|---|---|---|---|---|---|
| **A:** | A | C | – | – | R | T | W | Y |
| **B:** | A | C | A | I | – | S | W | Y |

$$\uparrow g_o$$

Figure 1: An example of gap opening vs. gap extension scores

How do we implement the two different gap scores into our sequence alignment algorithm? The idea of simply using $g_o$ or $g_e$ instead of the fixed gap score $g$ does not work. Consider, for instance the following alignment.

$$\downarrow$$

|       |   |   |   |   |   |   |
|-------|---|---|---|---|---|---|
| **A:** | . | . | . | . | . | – |
| **B:** | . | . | . | . | . | $b_m$ |

Here, the characters prior to $b_m$ have been aligned, and we want to align $b_m$ against a gap. In order to determine whether to use the score $g_o$ or $g_e$ in aligning $b_m$ against a gap in sequence A, we need to know whether the character preceding the gap in A was an amino acid or a gap. However, this information is not stored in the score matrix generated by our dynamic

programming algorithm. Therefore, we will need to adopt a slightly modified approach to implement gap extensions and openings.

The key idea in our new approach is that at any point in the alignment process, we can extend our existing alignment in three possible ways: by aligning $a_n$ against $b_m$, by aligning $a_n$ against a gap, or by aligning $b_m$ against a gap. Therefore, we will have three dynamic score matrices rather than just one. The matrices will be defined as follows:

$T(n, m)$: The matrix of optimal scores so that $a_n$ is aligned against $b_m$.

$T_-(n, m)$: The matrix of optimal scores so that $a_n$ is aligned against a gap, and $b_m$ precedes that gap.

$T^-(n, m)$: The matrix of optimal scores so that $b_m$ is aligned against a gap, and $a_n$ precedes that gap.

$$
\begin{array}{cc}
T(n, m) & \left|\begin{array}{c} \cdots\cdots \\ \cdots\cdots \end{array}\right.\left|\begin{array}{c} a_n \\ b_m \end{array}\right. \\[2ex]
T_-(n, m) & \left|\begin{array}{c} \cdots\cdots \\ \cdots b_m \end{array}\right.\left|\begin{array}{c} a_n \\ - \end{array}\right. \\[2ex]
T^-(n, m) & \left|\begin{array}{c} \cdots a_n \\ \cdots\cdots \end{array}\right.\left|\begin{array}{c} - \\ b_m \end{array}\right.
\end{array}
$$

Figure 2: The three dynamic score matrices

In figure 3, we illustrate how to set up the recurrence relations for use with our new dynamic matrices. Each equation is shown, along with the picture of the alignment corresponding to that equation. These recurrence relations follow naturally from the definitions of our dynamic matrices. The initialization for each of the matrices is identical to the initialization of the dynamic matrix for sequence alignment with a single fixed gap penalty.

$$T(n+1, m+1) = \max \begin{cases} T(n, m) + S(a_{n+1}, b_{m+1}) \\ T_-(n, m) + S(a_{n+1}, b_{m+1}) \\ T^-(n, m) + S(a_{n+1}, b_{m+1}) \end{cases} \qquad \begin{array}{ccc} \cdots & a_n & a_{n+1} \\ \cdots & b_m & b_{m+1} \\ \hline \cdots\cdots & a_n & a_{n+1} \\ \cdots b_m & - & b_{m+1} \\ \hline \cdots a_n & - & a_{n+1} \\ \cdots\cdots & b_m & b_{m+1} \end{array}$$

$$T_-(n+1, m+1) = \max \begin{cases} T(n, m+1) + g_o \\ T_-(n, m+1) + g_e \\ T^-(n, m+1) + g_o \quad (*) \end{cases} \qquad \begin{array}{ccc} \cdots & a_n & a_{n+1} \\ \cdots & b_{m+1} & - \\ \hline \cdots\cdots & a_n & a_{n+1} \\ \cdots b_{m+1} & - & - \\ \hline \cdots a_n & - & a_{n+1} \\ \cdots\cdots & b_{m+1} & - \end{array}$$

$$T^-(n+1, m+1) = \max \begin{cases} T(n+1, m) + g_o \\ T_-(n+1, m) + g_o \quad (*) \\ T^-(n+1, m) + g_e \end{cases} \qquad \begin{array}{ccc} \cdots & a_{n+1} & - \\ \cdots & b_m & b_{m+1} \\ \hline \cdots\cdots & a_{n+1} & - \\ \cdots b_m & - & b_{m+1} \\ \hline \cdots a_{n+1} & - & - \\ \cdots\cdots & b_m & b_{m+1} \end{array}$$

Figure 3: Recurrence relations for the dynamic score matrices

Note that in the equations marked with (*), the alignment segments $\left(\begin{smallmatrix} - & a_{n+1} \\ b_{m+1} & - \end{smallmatrix}\right)$ and $\left(\begin{smallmatrix} a_{n+1} & - \\ - & b_{m+1} \end{smallmatrix}\right)$ both contribute a score of $2g_o$, whereas aligning $\left(\begin{smallmatrix} a_{n+1} \\ b_{m+1} \end{smallmatrix}\right)$ would contribute the score $S(a_n, b_m)$. Any reasonable value of $g_o$ will be quite small, so that for any pair of amino acids, $S(a_n, b_m) > 2g_o$. Therefore, both equations marked by (*) are redundant, and the resulting recurrence relations are:

$$T(n+1, m+1) = \max \begin{cases} T(n, m) & + & S(a_{n+1}, b_{m+1}) \\ T_-(n, m) & + & S(a_{n+1}, b_{m+1}) \\ T^-(n, m) & + & S(a_{n+1}, b_{m+1}) \end{cases} \qquad (1)$$

3

$$T_-(n+1, m+1) \;=\; \max \begin{cases} T(n, m+1) & + & g_o \\ T_-(n, m+1) & + & g_e \end{cases} \qquad (2)$$

$$T^-(n+1, m+1) \;=\; \max \begin{cases} T(n+1, m) & + & g_o \\ T^-(n+1, m) & + & g_e \end{cases} \qquad (3)$$

Finally, to recover the optimal alignment, we would need to backtrack through our scoring matrices. To do that, we must first determine which of the three matrices has the maximum value at position $(N, M)$, where $N$ is the length of sequence A and $M$ is the length of sequence B. Afterwards, we proceed according to the recurrence relation, identically to the single fixed gap sequence alignment algorithm. In the figure below, for example, we start off at $T(M, N)$ and backtrack to $T^-(M-1, N-1)$. This corresponds to the last two characters in the alignment being $\left( \begin{smallmatrix} - & A_N \\ B_{M-1} & B_M \end{smallmatrix} \right)$.
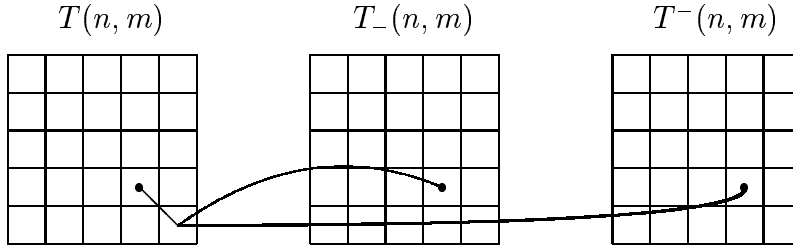
$T(n, m)$         $T_-(n, m)$         $T^-(n, m)$



Figure 4: Back-tracking through the matrices