

1

2 Implementing My Cost Function

I now show how to implement this cost function with dynamic programming (with space and time complexity $\mathcal{O}(nm)$, m and n are the lengths of the two sequences). I assume the sequences are from an alphabet of "ABC", and I will use the following single letter substitution score matrix:

$$\begin{matrix} & A & B & C & - \\ \begin{matrix} A \\ B \\ C \\ - \end{matrix} & \begin{pmatrix} 2 & -2 & -2 & -3 \\ -2 & 2 & -2 & -3 \\ -2 & -2 & 2 & -3 \\ -3 & -3 & -3 & 0 \end{pmatrix} \end{matrix}$$

A better substitution matrix could be designed for use comparing DNA sequences that penalises the less likely pairs more.

2.1 Algorithm

The algorithm works similarly to the standard quadratic space and time DP algorithm. The change is that, at each step, it considers the indel options (up and left) and checks if they continued beforehand, ie. if it was checking the up indel, we follow the pointers as long as they point in the up direction, and report how many there were. This is the job of the traceback function. This allows us to give different scores depending on how many indels there are from the max alignment back from that point.

Algorithm 1 Dynamic programming implementation of my cost function

Data: Two sequences S_1, S_2 only containing letters "ABC"

Result: Score, Path₁, Path₂

pointers $\leftarrow [n \times m \text{ matrix of 2's}]$, sc_matrix $\leftarrow [n \times m \text{ matrix of 0's}]$

```

for  $j = 1:n$  do
  for  $i = 1:m$  do
    upcount = traceback(i-1, j, 'up')
    case0 = sc_matrix[i-1, j] + (if upcount == {0: delete( $S_1[i-1]$ ), 1: -2, 2: -1, else: 0})
    leftcount = traceback(i, j-1, 'left')
    case1 = sc_matrix[i, j-1] + (if leftcount == {0: insert( $S_2[j-1]$ ), 1: -2, 2: -1, else: 0})

    local_options = [case0, case1, 0, sc_matrix[i-1, j-1] + match( $S_1[i-1]$ ,  $S_2[j-1]$ )]
    pointers[i, j] = argmax(local_options)
    matrix[i, j] = local_options[pointers[i, j]]
  end
  mx  $\leftarrow \max(\text{sc\_matrix}[:, j])$ 
  if local_max < mx then
    local_max  $\leftarrow$  mx
    max_pos  $\leftarrow (\text{argmax}(\text{sc\_matrix}[:, j]), j)$ 
  end
end
return fulltraceback(max_position)

```

I have defined the pointers to have codes for each direction: 0 = up, 1 = left, 2 points to itself (an alignment starts here), 3 = diagonal. This is so that the argmax for local_options will select any traces of previous indels with equal score to either of the other two options (starting afresh and matching (diagonally)). The functions: delete(a), insert(a) and match(a, b) correspond to the substitutions scores for score(a, '-'), score('-', a) and score(a, b).