

## ALIGNMENT

The sequence alignment problem:

- given:
  - 2 sequences ( $X$  and  $Y$ )
  - scoring matrix ( $S$ )

- compute: the pairwise alignment of  $X$  and  $Y$  of MAXIMUM score

**Global alignment**: ~optimal alignment along the entirety of both sequences

For example: given:

$$\begin{aligned} X &= \text{ACAAAT} \\ Y &= \text{TCAAGAT} \end{aligned}$$

with scoring scheme:

- +0 for gap
- +0 for mismatch
- +1 for match

we could get the alignment:

$$\begin{array}{ccccccc} T & C & A & G & A & T \\ A & C & A & - & A & T \end{array}$$

$$\text{score: } 0+1+1+0+1+1 = 4$$

Thus, there are 3 possible alignments for a letter in a sequence:

- MATCH**: align letter w/ same letter in other sequence ( $A$ )
- MISMATCH**: align letter not w/ same letter ( $T$ )
- GAP/INDEL**: align letter w/ gap (-)

\* biological application of indels: an insertion/deletion mutation @ some point in evolutionary history

\* There is a bijection (1:1 correspondence) between alignments of  $X$  and  $Y$  and directed paths from the top left cell (beginning) to bottom right cell (end) of edit graph



- the edit graph is a directed graph with edge weights
- max alignment score = max directed path from beginning  $\rightarrow$  end

Suppose sequence  $X$  is of size  $m$  and  $Y$  is of size  $n$ :

$\rightarrow$  # of alignments b/w  $X$  and  $Y$  is exponential

However, this algorithm will find the optimal alignment in quadratic ( $O(mn)$ ) time!

Hooray! what a beautiful algorithm!  
We love dynamic programming!

Now, for the algorithm:

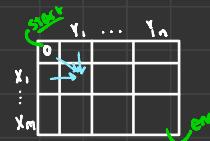
$$X = x_1, x_2, x_3, \dots, x_m ; Y = y_1, y_2, y_3, \dots, y_n$$

• Edit graph:

- dimensions:  $(m+1) \times (n+1)$

- entries have form  $(i, j)$

$$1 \leq i \leq m ; 1 \leq j \leq n$$



• edges: 3 types : horizontal, vertical, diagonal

- vertical : gap in Y  $(i-1, j) \rightarrow (i, j)$  ( $\delta(x_i, -)$ )

- horizontal : gap in X  $(i, j-1) \rightarrow (i, j)$  ( $\delta(-, y_j)$ )

- diagonal: alignment (match/mismatch)  $(i-1, j-1) \rightarrow (i, j)$  ( $\delta(x_i, y_j)$ )

•  $S(i, j)$  = score of the max score path from start to  $i, j$

ex]

$i-1, j-1$	$i-1, j$
$i, j-1$	$i, j$

→ any optimal path from start  $\rightarrow (i, j)$  must use one of the 3 green edges

scoring scheme:

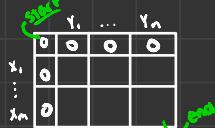
$$S(i, j) = \max \begin{cases} S(i-1, j) + \delta(x_i, -), \\ S(i, j-1) + \delta(-, y_j), \\ S(i-1, j-1) + \delta(x_i, y_j) \end{cases}$$

lost of Y gap

lost of X gap

cost of aligning  $x_i$  and  $y_j$   
(either match or mismatch)

First, you must initialize the edit graph (depending on the scoring scheme - this one has <sup>assume</sup> +0 gap penalty)



Then, you can go cell by cell, calculating  $S(i, j)$  based on the 3 surrounding cells.

