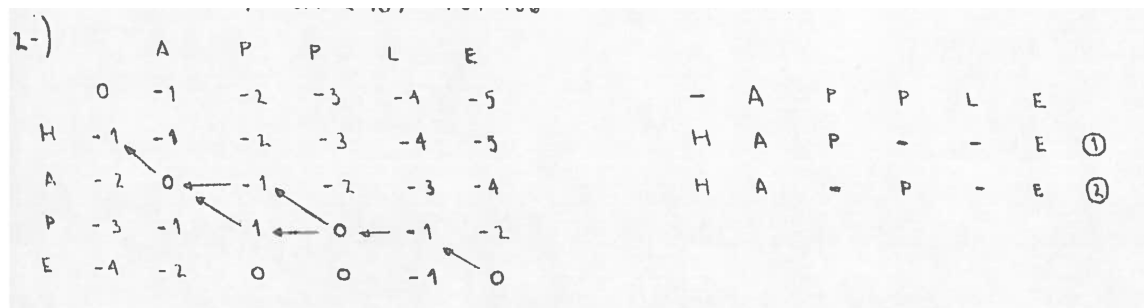


Question 1:



Question 2:

S = CCCC

T = CACACAC

Best alignment with linear gap penalty:

C - C - C - C - C

C A C A C A C A C

Score with linear gap penalty = +1

Score with affine gap penalty = -5

Best alignment with affine gap penalty:

C C C C C - - - -

C A C A C A C A C

Score with linear gap penalty = -3

Score with affine gap penalty = -3

Question 3.

- a) The two sequences cannot have lengths that are too different, as this would force us to use two consecutive gaps somewhere. In fact, a solution is only achievable if $|m-n| \leq \min(m,n)+1$. That's because $|m-n|$ is the number of gaps that will need to be inserted in the shorter sequence. Those gaps have to be interleaved with nucleotides. Starting and ending with a gap, we get $\min(m,n)+1$.

Example for $m=7$, $n=3$

AAAAAAA

-A-A-A-

- b) We first observe that the no-multi-gap alignment is a special case of the pairwise alignment problem with affine gap penalty, when the gap extension penalty is infinite: $\text{score}(L) = -d - e*(L-1)$, where $e = +\infty$. (Note: this is slightly different from the affine penalty scoring scheme presented in class, which was $\text{cost}(L) = a + b*L$, but it doesn't make a big difference).

Following the algorithm presented in the Durbin et al. book (Equation 2.16), we introduce three dynamic programming tables: M , I_x , and I_y . M is computed with no change. Since $e = +\infty$, $I_x(i,j)$ and $I_y(i,j)$ reduce to just $M(i-1,j)-d$ and $M(i,j-1)-d$ (respectively). Initialization is done as follows (not detailed in the book):

$M(0,0) = 0$

$M(i,0) = M(0,j) = +\infty$, for $i,j > 0$

$I_x(0,0) = 0$

$I_x(1,0) = -d$

$I_x(i,0) = -\infty$ for $i > 1$

$I_y(0,0)=0$
 $I_y(0,1) = -d$
 $I_y(0,j) = -\text{inf}$ for $j>1$

The rest of the algorithm proceeds like described in the book. Trace-back is performed from $\max(M(m,n), I_x(m,n), I_y(m,n))$.

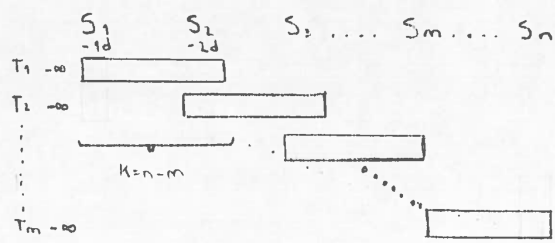
c) I'm too old for this...

Question 4:

4-) Simply not allow for mismatches by making the penalty of a mismatch $-\infty$ or
 Change the recurrence

$$A_{i,j} = \max \begin{cases} A_{i-1,j-1} + 1 & \text{if } S_i = T_j \\ A_{i,j-1} \\ A_{i-1,j} \end{cases} \quad 0 \text{ if } i=0 \text{ or } j=0$$

Question 5:

5-) 

$$F(i,j) \forall i,j (i-j \geq K) = \max \begin{cases} F(i-1,j-1) + s(x_i, y_j) \\ F(i-1,j) - d \end{cases}$$

Basis $F_{0,j} = d * j$
 $F_{i,0} = -\infty$

Question 6:

#1) $S_1: TGC$
 $S_2: TCA$
 $S_3: T GCA$

First align S_1 and S_2

the best alignment is $\begin{array}{ccc} T & G & C \\ & T & C & A \end{array}$ score: -1

then align $\begin{array}{|c|c|c|} \hline T & G & C \\ \hline T & C & A \\ \hline \end{array}$ with $TGCA$

the best alignment is $\begin{array}{|c|c|c|c|} \hline T & - & G & C \\ \hline T & - & C & A \\ \hline T & G & C & A \\ \hline \end{array}$ score = $2 - 4 + 0 + 0 = -2$

or $\begin{array}{|c|c|c|c|} \hline T & G & C & - \\ \hline T & C & A & - \\ \hline T & G & C & A \\ \hline \end{array}$ score = $2 + 0 + 0 - 4 = -2$

and the final score of the alignment is:

$\begin{array}{|c|c|c|} \hline T & - & G & C \\ \hline T & - & C & A \\ \hline T & G & C & A \\ \hline \end{array}$ score: $3 - 4 - 1 - 1 = -3$

But there exist a better alignment

$\begin{array}{|c|c|c|} \hline T & G & C & - \\ \hline T & - & C & A \\ \hline T & G & C & A \\ \hline \end{array}$ score $3 - 3 + 3 - 3 = 0$

Question 7:

$S = A A A A A A A A A A A A A A A A A A$
 $T = A A A A C A A A A C A A A A C A A A A C$

These two sequences are 80% identical but contain no exact match of size 5.

Bonus Question:

This can be done using the Hirschberg algorithm, which is described here.
https://en.wikipedia.org/wiki/Hirschberg%27s_algorithm