

COMP 561 Ass 1:

1) match = +1, mismatch = -1, indel = -1 * L

		0	1	2	3	4
0		0	-1	-2	-3	-4
A	1	-1	-1	0	-1	-2
P	2	-2	-2	-1	+1	0
P	3	-3	-3	-2	0	0
L	4	-4	-4	-3	-1	-1
E	5	-5	-5	-4	-2	0

opt sol #1 (pink)

- A P P L E
H A P - - E } cost = 0

opt sol #2 (blue)

- A P P L E
H A - P - E } cost = 0

these are the optimal global alignments.

2) let $S = C$, $T = CCC$

using linear gap penalty (with indel = -1 * L):

3 alignments to consider

① S' - C -
T' CCC

$$\text{score} = -1 + 1 - 1 = -1$$

② S' C - -
T' CCC

$$\text{score} = +1 - 1(2) = -1$$

③ S' - - C
T' CCC

$$\text{score} = -1(2) + 1 = -1$$

so all 3 alignments are optimal

using affine gap penalty (with indel = -2 - 0.5L)

considering same 3 alignments

① S' - C -
T' CCC

$$\begin{aligned} \text{score} &= -2 - 0.5(1) + 1 - 2 - 0.5(1) \\ &= -4 \end{aligned}$$

② S' C - -
T' CCC

$$\begin{aligned} \text{score} &= 1 - 2 - 0.5(2) \\ &= -2 \end{aligned}$$

③ S' - - C
T' CCC

$$\begin{aligned} \text{score} &= -2 - 0.5(2) + 1 \\ &= -2 \end{aligned}$$

so opt global alignments are only 2 & 3 →

different from result using linear gap penalty (1 was also opt).

3) a) you can only have multi-gap if it is possible to have 2+ consecutive gaps

ie CANT have if $S = A, T = AA$ (can have max 1 gap)

so constraint is $|m - n| \leq \min(m, n) + 1$

b) Algo: // assumes input M, b, S, T ($b < 0$)

initializing {
for $i = 1 \rightarrow m$:
 $X_{i,0} = b \times i$
 $P_{i,0} = \uparrow$ // let P = pointer matrix
for $j = 1 \rightarrow n$:
 $X_{0,j} = j \times b$
 $P_{0,j} = \leftarrow$
 $P_{0,0} = 0$ // know when to stop

calculation {
for $i = 1 \rightarrow m$:
 for $j = 1 \rightarrow n$:
 $\text{maxval} = \max(X_{i,j} + M(S_i, T_j), X_{i-1,j} + b, X_{i,j-1} + b)$
 if $\text{maxval} = X_{i,j} + M(S_i, T_j)$:
 $P_{i,j} = \nwarrow$
 if $\text{maxval} = X_{i-1,j} + b$:
 if $X_{i-1,j-1}$ doesn't contain \uparrow : // prevents having $\uparrow\uparrow$ as only option
 $P_{i,j} += \uparrow$
 if $\text{maxval} = X_{i,j-1} + b$: // prevents having $\leftarrow\leftarrow$ as only option
 if $X_{i,j-1}$ doesn't have \leftarrow :
 $P_{i,j} += \leftarrow$
 if $P_{i,j}$ is empty:
 $\text{maxval} = X_{i,j} + M(S_i, T_j)$
 $P_{i,j} = \nwarrow$
 $X_{i,j} = \text{maxval}$

3b) continued:

$S' = ''$, $T' = ''$, $p_copy = p.copy()$

$num = 0$, $iters = 1$, $prev = []$, $branch = []$

while $iters > 0$: // continue until $iters = 0$ OR $p(m,n) = ''$

if $num \neq 0$:

$S'.append('')$

$T'.append('')$

$i = m$

$j = n$

$lastS = ''$

$lastT = ''$

$others = ''$

while $(i > 0 \text{ or } j > 0)$:

if $(len(P_{i,j}) = 0)$:

if not others: // options left in previous step &
current step is empty

$S'[num] = ''$, $T'[num] = ''$ //empty S' & T'

$iters += 1$

$num -= 1$

else: // no options left

$S'[num] = ''$

$T'[num] = ''$

break

if $(len(P_{i,j}) > 1)$: // multiple options (ie. \uparrow or \leftarrow)

$iters += len(P_{i,j}) - 1$

$branch.append((i,j))$

$others = P_{i,j}[:-1]$

$move = P_{i,j}[-1]$

else:

// only one option

$move = P_{i,j}$

$others = move$

if move = ↖

$s'[num] = s[i-1] + s'[num]$

$T'[num] = T[j-1] + T'[num]$

$i -= 1$

$j -= 1$

$lasts = s[i-1]$

$lastT = T[j-1]$

elif move = ←

if lasts != '-':

$s'[num] = "-" + s'[num]$

$T'[num] = T[j-1] + T'[num]$

$lasts = '-', lastT = T[j-1]$

$j -= 1$

else:

$i = 0, j = 0$

$s'[num] = "", T'[num] = ""$

elif move = ↑

if lastT != '-':

$s'[num] = s[i-1] + s'[num]$

$T'[num] = '-' + T'[num]$

$lasts = s[i-1], lastT = "-"$

$i = -1$

else:

$i = 0, j = 0$

$s'[num] = "", T'[num] = ""$

iters -= 1

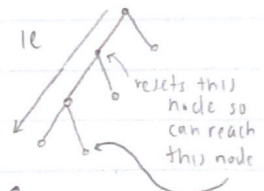
num += 1

if branch:

if prev:

if $prev[-1] < branch[-1]$ then $p[prev[-1] \text{ coords}] = p_{\text{copy}}[prev[-1]]$,
 $prev = prev[:-1]$

$p[branch[-1] \text{ coords}] = p[branch[-1] \text{ coords}][:-1]$ //removes last operation



```
prev.append(branch[-1])  
branch = branch[:-1]  
if pmin == '': break
```

```
// print alignments & scores  
for k = 0 → len(s'):  
    print s'[k]  
    print T'[k]  
print("SCORE", x_min).
```

3c) see Q3c.py for code, Q3_results.txt for results

4) let $X_{i,j}$ be the length of the longest common subsequence of $S = s_1 \dots s_i$ and $T = t_1 \dots t_j$

- we know that smith-waterman (s-w) algorithm computes

$$X_{i,j} = \begin{cases} 0 \\ X_{i-1,j-1} + M(s_i, t_j) \\ X_{i-1,j} + c \\ X_{i,j-1} + c \end{cases}$$

- If want longest common subsequence, want \nwarrow if is a match, don't allow \nwarrow if is a mismatch (meaning only consider \leftarrow or \uparrow).

CHANGES:

① set $c=0 \rightarrow$ subsequence allows for gaps, so don't want to penalize \rightarrow also, LCS algo from source says, if $s_i \neq t_j$, $X_{i,j} = \max(X_{i-1,j}, X_{i,j-1})$

② set $M(s_i, t_j) = \begin{cases} 1 & \text{if } s_i = t_j \\ -\min(\text{len}(S), \text{len}(T)) \end{cases}$

\hookrightarrow is because, if mismatch, don't want to consider \nwarrow , so want

$$X_{i-1,j-1} + M(s_i, t_j) \leq 0$$

$$\text{so } M(s_i, t_j) \leq -(X_{i-1,j-1})$$

\hookrightarrow worst case = only matches before

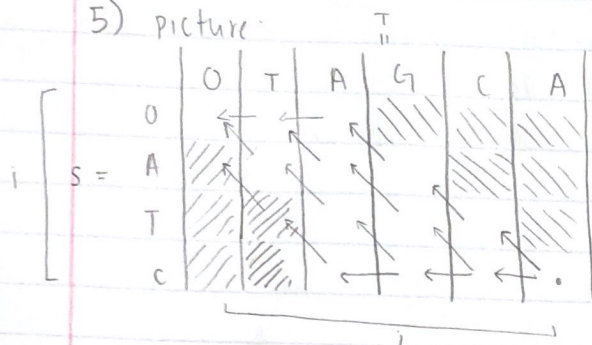
		O	A	A	A
O		0	0	0	0
A		0	1	1	1
T		0	1	1	1

\hookrightarrow DON'T want: $1-2 = -1 < 0$

So, run S-W algo, with above changes.

- Then, during traceback, start NOT at $X_{i,j}$ but rather at $\max(X_{i,j})$. Additionally, count s_i or t_j in X if \nwarrow (don't count if \leftarrow or \uparrow because is nucleotide only present in one of the sequences).

5) picture:



deletion from S NOT allowed \rightarrow
so "T" not allowed

/// = can't go because would require \uparrow

let $S = S_1 \dots S_m$, $T = T_1 \dots T_n$, with $m \leq n$, $k = m - n$,

$b = \text{linear gap penalty } (b < 0)$

algo: for $j = 0 \rightarrow k$:

$$X_{0,j} = b * j$$

// DONT need to initialize $X_{i,0}$ b/c can never reach

for $i = 1 \rightarrow m$: // because can reach all rows

for $j = 1 \rightarrow k+1$:

$$X_{i,j} = \max \begin{cases} X_{i-1,j-1} + M(S_i, T_j) \\ X_{i,j-1} + b \end{cases}$$

// give matrix " \leftarrow " or " \uparrow "

// do same trace back as in N-W algo

\hookrightarrow this is $O(m * k)$

6) given $c = -1$, mismatch = -1 , match = 0

assume

```

      ACT
     /  \
    /    \
   /      \
  /        \
 /          \
/            \
ACGT
  
```

algo gives:

AC-T	or	A-CT
AG-T		A-GT
ACGT		ACGT
ACGT		ACGT

$$\text{score} = (-1) + (-1) + (-1) + (-1) + (-1) + (-1) = -7$$

$$\text{score} = (-1) + (-1) + (-1) + (-1) + (-1) + (-1) = -7$$

optimal solution is:

```

AC-T
A-GT
ACGT
ACGT
  
```

$$\text{score} = (-1) + (-1) + (-1) + (-1) + (-1) = -6$$

↳ score of $-6 > -7$ → so algo DOESN'T give optimal result.

7) S = AAAAC AAAAC AAAAC AAAAC

T = AAAAA AAAAA AAAAA AAAAA

↳ never get 5-mers that match, so fail to align