

COMP561 Assignment 1

by Jingyan Wang
ID: 260860682

<Question 1>

Substitution matrix

M: A P L E H

insertion/deletion
penalty = $-1 \times L$

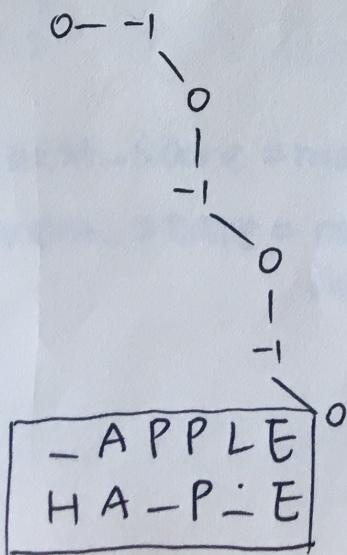
A	+1	-1	-1	-1	-1
P	-1	+1	-1	-1	-1
L	-1	-1	+1	-1	-1
E	-1	-1	-1	+1	-1
H	-1	-1	-1	-1	+1

Perform Needleman-Wunsch alg.

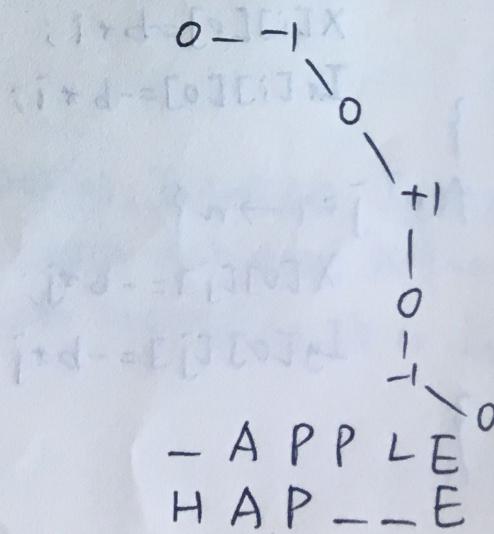
H A P E

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

∴ Solution 1



→ Solution 2



<Question 2>

Suppose pair of sequences : S: A , T: AAA.

there are 3 good alignments . which are :

① A -- ② - A - ③ - - A
AAA AAA AAA.

a) If we use linear gap penalty $\text{cost}(L) = -1 * L$,

all 3 alignments have score of $\text{score}(\text{match}) - 2$,
i.e. ~~we will possibly choose all three~~ alignments are optimal.

b) If we use affine gap penalty $\text{cost}(L) = -2 - 0.5 * L$.
alignment ① & ③ will achieve score of $\text{score}(\text{match}) - 3$
whereas alignment ② has score of $\text{score}(\text{match}) - 5$.
That is alignment ① & ③ are optimal.

<Question 3 >

a) constraint : $|m-n| \leq \min(m,n) + 1$

b) Input: M, b, S, T.

Optimal global alignment {M, b, S, T}:

State X, State Ix, State Iy [m+1][n+1] // to store FSA state
for i = 1 → m { with each condition

$$X[i][0] = b * i;$$

$$I_x[i][0] = -b * i;$$

for j = 1 → n {

$$X[0][j] = -b * j;$$

$$I_y[0][j] = -b * j$$

}

for $i = 1 \rightarrow m\{$

for $j = 1 \rightarrow n\{$

// X is highest score given s_i aligned to t_j
 $X[i][j] = \max(X[i-1][j-1] + M[i][j],$

$I_x[i-1][j-1] + M[i][j],$

$I_y[i-1][j-1] + M[i][j]);$

if ($X[i][j] = X[i-1][j-1] + M[i][j]$) store ~~$X \rightarrow X$~~ ;

if ($X[i][j] = I_x[i-1][j-1] + M[i][j]$) store ~~$X \rightarrow I_x$~~ ;

if ($X[i][j] = I_y[i-1][j-1] + M[i][j]$) store $X \rightarrow I_y$;
// store $X \rightarrow X$ means State $X[i][j] = X$

// store $X \rightarrow I_x$ means State $X[i][j] = I_x$

// store $X \rightarrow I_y$ means State $X[i][j] = I_y$.

// I_x is highest score given s_i aligned to '_'

$I_x[i][j] = \max(X[i-1][j] - b, I_y[i-1][j] - b);$

if ($I_x[i][j] = X[i-1][j] - b$) State $I_x[i][j] = X$;

if ($I_x[i][j] = I_y[i-1][j] - b$) State $I_x[i][j] = I_y$;

// I_y is highest score given t_j aligned to '_'

$I_y[i][j] = \max(X[i][j-1] - b, I_x[i][j-1] - b);$

if ($I_y[i][j] = X[i][j-1] - b$) State $I_y[i][j] = X$;

if ($I_y[i][j] = I_x[i][j-1] - b$) State $I_y[i][j] = I_x$;

}

highest-score = $\max(X[m][n], I_x[m][n], I_y[m][n])$;

current-state = matrix which achieved highest-score
in its last item; // i.e. X or I_x . or I_y .

// trace back with states (similar to FSA)

i=m; j=n;

while i>0 & j>0 {

if (current-state = X){

current-state = State[i][j];

align s_i with t_j ;

i-=1; j-=1;

if (current-state = I_x) {

current-state = State[I_x][i][j];

align s_i with '-' ;

i-=1;

if (current-state = I_y) {

current-state = State[I_y][i][j];

align '-' with t_j

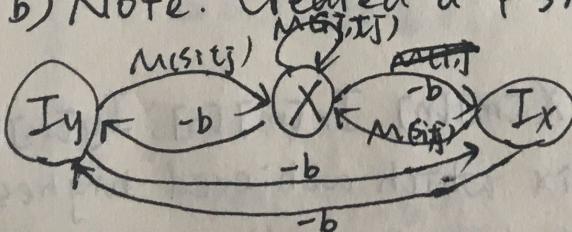
j-=1;

if (i=0) { while(j!=0) { align t_j with '-' ; j-=1; } }

if (j=0) { while(i!=0) { align s_i with '-' ; i-=1; } }

print (alignment, highestScore).

Q3. b) Note. Created a FSA like:



Q3. c) Please see align.py for code.

and optimal-alignment.txt for optimal alignment and scores.

<Question 4 >

Let $X_{i,j}$ be the length of longest common sequence of $s_i \dots s_i$ and $t_i \dots t_j$.

a) Use Smith Waterman alg to calculate $X_{i,j}$ each time

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

where we input parameters of this formula.

$$\textcircled{1} \quad M(s_i, t_j) = \begin{cases} 1, & \text{where } s_i = t_j \\ -\min(\text{len}(S), \text{len}(T)), & \text{where } s_i \neq t_j. \end{cases}$$

(assign a constant value to $M(s_i, t_j)$ where $s_i \neq t_j$,
the value is less or equal to $\frac{-2}{\min(\text{len}(S), \text{len}(T))}$)

to ensure $X_{i-1,j-1} + M(s_i, t_j)$ is not larger than 0
when $s_i \neq t_j$)

$$\textcircled{2} \quad c = 0$$

b) When trace back, start with $\arg \max_{i,j} \{X_{i,j}\}$
count s_i or t_j whenever there is a " Δ " as part of the common sequence.

c) Other steps are identical with Smith Waterman alg.

<Question 5>

~~Similar~~ to N-W alg, use $X_{i,j}$ to denote best score of alignment for $s_1 \dots s_i$ and $t_1 \dots t_j$.

for $j=0 \rightarrow k$ // initialization.

$X_{0,j} = -b*j$ // b is linear gap penalty.

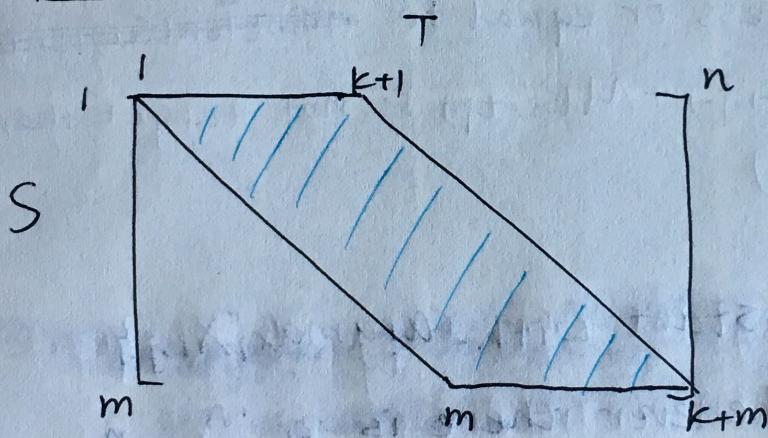
for $i=1 \rightarrow m$

for $j=i \rightarrow k+i$

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

// and give " \uparrow " or " \leftarrow " to the ~~matrix~~ matrix.

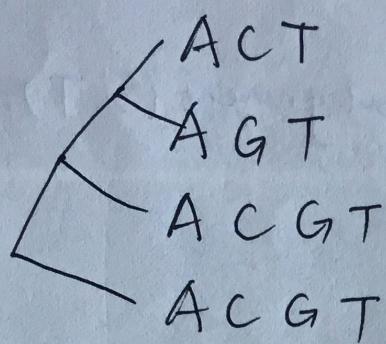
// traceback same with N-W.



After performing algorithm above, we are recording $X_{i,j}$ inside this parallelogram, which gives base k , height m , and proves the alg is in time $O(m*k)$

<Question 6 >

⇒ Set of DNA sequences:



⇒ Alignment by the alg

A C - T	6R	A - C T
A G - T		A - G T
A C G T		A C G T
A C G T		A C G T

Score: both -7

⇒ optimal alignment

A C - T
A - G T
A C G T
A C G T

Score: -6.

<Question 7 >

S: AAAAAA CCCCCC GGGGGG TTTTTT

T: AAAAAT CCCCAGGGGCTTTTG

< Bonus Question >

To create such an algorithm, we can use Divide & Conquer scheme to find one element of alignment at one time, with following function.

func divide-and-conquer-alignment (S, T):

$m = \text{len}(S);$

$n = \text{len}(T);$

if ($m \leq 1$ or $n \leq 2$) {

use Needleman Wunsch alg,

return alignment

}

Perform ForgetPast-NW on $(S, T[1 : \frac{n}{2}])$;

Perform Backward-ForgetPast-NW on $(S, T[\frac{n}{2} + 1 : n])$;

Suppose the result of last 2 steps is f and g ,

find q that ~~maximize~~ $f(i, j) + g(i, j)$.

add q into a global alignment table P .

// begin divide & conquer

divide-and-conquer-alignment ($S[1 : q], T[1 : \frac{n}{2}]$)

divide-and-conquer-alignment ($S[q : m], T[\frac{n}{2} : n]$)

return P

}

Notes: ① Backward-ForgetPast-NW is to consider the $[m, n]$ be the start point, and perform the same alg in the opposite direction

② each time, q is the optimal alignment point we need

③ the size of P is at most $m+n$, which

takes up space usage of this alg (when no match in nucleotide)