

Lab 9

For this lab, you will implement the Needleman-Wunsch algorithm for sequence alignment. The main recursion is:

$$\text{score}_{i,j} = \max \begin{cases} \text{score}_{i-1,j-1} + \begin{cases} M & \text{if } v_i = w_j \\ N & \text{if } v_i \neq w_j \end{cases} \\ \text{score}_{i-1,j} + I \\ \text{score}_{i,j-1} + I \end{cases}$$

Here is the pseudocode (where strings are 1-indexed):

Provide strings v and w and output an alignment as two strings

Also provide I , the insert cost, M the match gain, and N , the non-match cost.

NeedlemanWunsch(v, w, I, M, N)

Build a 2D array, called `score`, of size $(\text{length}(v) + 1) \times (\text{length}(w) + 1)$

for $i=0$ to $\text{length}(v)$

$\text{score}_{i,0} \leftarrow I \times i$

for $j=0$ to $\text{length}(w)$

$\text{score}_{0,j} \leftarrow I \times j$

for $i=1$ to $\text{length}(v)$

 for $j=1$ to $\text{length}(w)$

 if $(v_i = w_j)$

$\text{match} = M$

 else

$\text{match} = N$

$\text{score}_{i,j} \leftarrow \max(\text{match} + \text{score}_{i-1,j-1}, \text{score}_{i-1,j} + I, \text{score}_{i,j-1} + I)$

$\text{Alignment}_v \leftarrow ""$

$\text{Alignment}_w \leftarrow ""$

$i \leftarrow \text{length}(v)$

$j \leftarrow \text{length}(w)$

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

 if $(i > 0 \text{ and } j > 0)$

 if $(v_i = w_j)$

$\text{match} = M$

 else

$\text{match} = N$

 if $\text{score}_{i,j} = \text{score}_{i-1,j-1} + \text{match}$

$\text{Alignment}_v \leftarrow v_i + \text{Alignment}_v$

$\text{Alignment}_w \leftarrow w_j + \text{Alignment}_w$

$i \leftarrow i - 1$

$j \leftarrow j - 1$

```

else if scorei,j = scorei-1,j + l
    Alignment_v ← vi + Alignment_v
    Alignment_w ← "-" + Alignment_w
    i ← i - 1
else
    Alignment_v ← "-" + Alignment_v
    Alignment_w ← wj + Alignment_w
    j ← j - 1
else if (i > 0)
    Alignment_v ← vi + Alignment_v
    Alignment_w ← "-" + Alignment_w
    i ← i - 1
else
    Alignment_v ← "-" + Alignment_v
    Alignment_w ← wj + Alignment_w
    j ← j - 1

return Alignment_v, Alignment_w

```

One tricky aspect is setting up the 2D score array. In python, this is a list of lists. Here is python code for making the list of lists:

```

#create a list of lists that is indexed for v and then w
#this 0-indexed list of lists has 0 for the nucleotide before the
start
infinity = 10000
score=[]
for nucleotide in v:
    row = [infinity]*(len(w)+1)
    score.append(row)

row = [infinity]*(len(w)+1)
score.append(row)

```

Run the algorithm with the sequences:

v=

ACUUAGCUAAAACGUUUGGUUCAAACAUUUGCUUGCUGUCUUGGCAUAACAUCAAUAAAGGCAUAAACAU
CGCAAAACAAUGGUUAUAUAUAAAUGGCUAUGAGGAUGGUUUUAGUACGUAGGCGUUGCGGAACUUCGGU
UCAGAUAGAGCAAUGAAUCGUGCAUGCUAGGAAAACUGACCACACGCAGUUGGCAGCCCUAGUAUCUUUCG
AUAGAUUUCCAUACCUCCGCGAUC

w=

ACUUAGCCUAUACACUAUGUUGGAGAGAGACGCUUGCACCUAGGCAUAAUGUGAAUUAGGUUAUAAACAUC

GUGGUUGUAAACUUGAGUGGGUUUUAGUACGGUAUGCGUGAUUACUUCGUAAUCAUGAAUCGUGCAUGC
UAGUGGGGUUUGGCCUCCACUAGUAUCUUUGAAGAUUUUCCUCCUCAGCGAUC

And with $M = 1$, $N = -0.5$, $I = -1$