
NLP Homework 2, Part 2

Vahid Kharazi

November 27, 2015

1 LOCAL ALIGNMENT

According to wikipedia and course slides, we have to implement Smith-Waterman version of edit distance to solve local alignment problem.

$$H(i, 0) = 0, 0 \leq i \leq m$$

$$H(0, j) = 0, 0 \leq j \leq n$$

$$H(i, j) = \max \left\{ \begin{array}{ll} 0 & \\ H(i-1, j-1) + s(a_i, b_j) & \text{Match/Mismatch} \\ \max_{k \geq 1} \{H(i-k, j) + W_k\} & \text{Deletion} \\ \max_{l \geq 1} \{H(i, j-l) + W_l\} & \text{Insertion} \end{array} \right\}, 1 \leq i \leq m, 1 \leq j \leq n$$

Where:

- a, b = Strings over the **Alphabet** Σ
- $m = \text{length}(a)$
- $n = \text{length}(b)$
- $s(a, b)$ is a similarity function on the alphabet
- $H(i, j)$ - is the maximum Similarity-Score between a suffix of $a[1...i]$ and a suffix of $b[1...j]$
- W_i is the **gap-scoring** scheme

Test: There is an example in wikipedia to compute edit distance between ACACACTA and AGCACACA:

- $s(a, b) = +2$ if $a = b$ (match), -1 if $a \neq b$ (mismatch)
- $W_i = -1$

$$H = \begin{pmatrix} - & A & C & A & C & A & C & T & A \\ - & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ A & 0 & 2 & 1 & 2 & 1 & 2 & 1 & 2 \\ G & 0 & 1 & 1 & 1 & 1 & 1 & 0 & 1 \\ C & 0 & 1 & 3 & 2 & 3 & 2 & 3 & 2 \\ A & 0 & 2 & 2 & 5 & 4 & 5 & 4 & 4 \\ C & 0 & 1 & 4 & 4 & 7 & 6 & 7 & 6 \\ A & 0 & 2 & 3 & 6 & 6 & 9 & 8 & 8 \\ C & 0 & 1 & 4 & 5 & 8 & 8 & 11 & 10 \\ A & 0 & 2 & 3 & 6 & 7 & 10 & 10 & 12 \end{pmatrix} \quad T = \begin{pmatrix} - & A & C & A & C & A & C & T & A \\ A & 0 & \nearrow & \leftarrow & \nwarrow & \leftarrow & \nwarrow & \leftarrow & \nwarrow \\ G & 0 & \uparrow & \nwarrow & \uparrow & \nwarrow & \uparrow & \nwarrow & \uparrow \\ C & 0 & \uparrow & \nwarrow & \leftarrow & \nwarrow & \leftarrow & \nwarrow & \leftarrow \\ A & 0 & \nwarrow & \uparrow & \nwarrow & \leftarrow & \nwarrow & \leftarrow & \nwarrow \\ C & 0 & \uparrow & \nwarrow & \uparrow & \nwarrow & \leftarrow & \nwarrow & \leftarrow \\ A & 0 & \nwarrow & \uparrow & \nwarrow & \uparrow & \nwarrow & \leftarrow & \nwarrow \\ C & 0 & \uparrow & \nwarrow & \uparrow & \nwarrow & \uparrow & \nwarrow & \leftarrow \\ A & 0 & \nwarrow & \uparrow & \nwarrow & \uparrow & \nwarrow & \uparrow & \nwarrow \end{pmatrix}$$

Similar to wikipedia example, this is my result for same:

```
(venv)vahid@kharazi:~/dev/nlp/two$ python two.py
0 0 0 0 0 0 0 0 0
0 2 1 2 1 2 1 1 2
0 1 1 1 1 1 1 0 1
0 1 3 2 3 2 3 2 2
0 2 2 5 4 5 4 4 4
0 1 4 4 7 6 7 6 6
0 2 3 6 6 9 8 8 8
0 1 4 5 8 8 11 10 10
0 2 3 6 7 10 10 10 12
D D D D D D D D
D T D T D T D DTL T
D T D L D L D L L
D D T D L D L L DL
D T D T D L D L L DL
D T D T D T D L L
D D T D T D T D D
break
D , L D L D L L D
D , D T D T D DTL T
D T , L D L D L L
D D T , L D L L DL
D T D T , L D L L DL
D D T D T , L L
D T D T D T , L
D D T D T D T D ,
['A', 'C', 'A', 'C', 'A', 'C', 'T', 'A']
['-', 'i', '-', '-', '-', '-', '-', 'd', '-']
['A', 'G', 'C', 'A', 'C', 'A', 'C', 'A']
```

The first matrix is H, Second matrix is T(T = Top, L = Left, D = Diag). Finally you can find the path of minimum edit distance:

```
['A', 'C', 'A', 'C', 'A', 'C', 'T', 'A']
['-', 'i', '-', '-', '-', '-', '-', 'd', '-']
['A', 'G', 'C', 'A', 'C', 'A', 'C', 'A']
```