

CS612 – Algorithms in Bioinformatics

Fall 2017 – Homework Assignment 1
Due Wed. Oct. 11 2017 in class

Objectives

- Learn how to perform sequence alignment and analyze the results
- Understand and apply the theory of sequence alignment

Part 1 – Practice

1. **Sequence alignment hands-on exercise.** You are given the following protein sequence:

TCPFADPAALYSRQDTTSGQSPLAAYEVDDSTGYLTSVGGPIQDQTSLKAGIRGPTLLEDFMFRQKIQHFDHERVPERAV

- (a) Go to the Blast website at <http://blast.ncbi.nlm.nih.gov>. Under “basic Blast” select “Protein Blast” to get to the BlastP website. Paste the above sequence to the top window. Use the default parameters – nr database (non-redundant protein sequence database). Hit “BLAST” – the search may take a few seconds. Save the search results using the “download” button and the “txt” option above. Submit a printout of the first 2 pages of the saved search as part of your homework. **DO NOT print out the entire search result – it’s very large.**
- (b) Repeat the search above with the SwissProt database and submit its first 2 pages as part of your homework.
- (c) Repeat search a. above with PAM30 as a substitution matrix. This can be done in the blastP homepage by opening “algorithm parameters” at the bottom of the page. Observe the changes between the results of a and c due to the change in the substitution matrix: Look at the first entry that differs between a and c. What is its rank in a and c? What is the name of the protein sequence in this entry?

2. **DNA sequence alignment:** The following sequence was constructed by NCBI scientist Mark Boguski for Michael Chrichton’s “The Lost World” of the Jurassic Park series:

>DinoDNA from THE LOST WORLD p. 135

```
GAATTCCGGAAGCGAGCAAGAGATAAGTCCTGGCATCAGATACAGTTGGAGATAAGGACG
GACGTGTGGCAGCTCCCGCAGAGGATTCACTGGAAGTGCATTACCTATCCATGGGAGGCC
ATGGAGTTCGTGGCCTGGGGGGCCGGATGCGGGCTCCCCACTCCGTTCCCTGATGAA
GCCGGAGCCTCCTGGGCTGGGGGGCGAGAGGACGGAGGAGGCGGGGGGCTGCTGGCC
TCCTACCCCCCCTCAGGCCGTGTCCTGGTGGCGAGACACGGGTACTTGGGG
ACCCCCCAGTGGGTGCCGCCACCCAAATGGAGCCCCCCCACCTGGAGCTGCTG
CAACCCCCCGGGGAGCCCCCCCACCCCTCCTCCGGGCCCTACTGCCACTCAGCAGC
GGGCCCCCACCCCTGCGAGGCCGTGAGTGCCTATGCCAGGAAGAACGGAGCGACG
GCAACGCCGCTGTGGCGCCGGACGGCACCGGGCATTACCTGTGCAACTGGGCTCAGCC
TGCGGGCTTACCAACGCCCTAACGCCAGAACGCCGCTCATCCGCCAAAAAGCGC
CTGCTGGTGAGTAAGCGCGCAGGCACAGTGTGCAGCCACGAGCGTGAAACTGCCAGACA
```

```

TCCACCACCACTCTGGCGTCGCAGCCCCATGGGGACCCCCTGCAACAAACATTAC
GCCTGCGGCCTACTACAAACTGCACCAAGTGAACCGCCCCCTACGATGCGCAAAGAC
GGAATCCAACCGAAACCGCAAAGTTCTCCAAGGGTAAAAGCGGCGCCCCCGGGG
GGGGGAAACCCCTCCGCCACCGCGGGAGGGGGCCTCTATGGGGGAGGGGGGACCCC
TCTATGCCCCCCCGCCGCCGCCGCCCTCAAAGCGACGCTCTGTAC
GCTCTGGCCCCGTGGCTTCCGCCATTCTGCCCTTGAAACTCCGGAGGGTT
TTTGGGGGGGGGGCGGGGGGTTACACGGCCCCCGGGCTGAGCCGCAGATTTAAATA
ATAACTCTGACGTGGCAAGTGGGCCTGCTGAGAAGACAGTGTAAACATAATAATTGCA
CCTCGGCAATTGCAGAGGGTCGATCTCCACTTGGACACAACAGGGCTACTCGTAGGAC
CAGATAAGCACTTGCTCCCTGGACTGAAAAAGAAAGGATTATCTGTTGCTTGTG
GACAAATCCCTGTGAAAGGTAAAAGTCGGACACAGCAATCGATTATTCTGCCGTGTG
AAATTACTGTGAATATTGTAATATATATATATATATATCTGTATAGAACAGCC
TCGGAGGCAGCATGGACCCAGCGTAGATCATGCTGGATTGTACTGCCGGATT

```

Perform a Blast search using blastn (nucleotide search) and the default non-redundant (nr) nucleotide database.

- (a) What are the two main species used to construct the dinosaur DNA sequence?
- (b) Repeat the search with blastx (DNA vs. protein sequence) using the default non-redundant protein sequence database. Look at the top sequence alignment and retrieve the hidden message there (hint: look at the gaps...).

Part 2 – Theory

1. Lesk book question 5.3: The edit distance between the strings agtcc and cgctca is 3, consistent with the following alignment:

```

agtcc
cgctca

```

Find the sequence of three edit operations that convert agtcc to cgctca.

2. **Dynamic programming:**

- (a) Use the Needleman Wunsch global alignment Dynamic programming formula in slide set no. 2 to find the sequence alignment score of the two DNA sequences TACGGGTAT and GGACGTACG. Show the filled dynamic programming matrix using +1 for a match, -1 for a mismatch and -1 for a gap penalty in a way similar to the slide sets.
- (b) Repeat with the Smith-Waterman local alignment algorithm and the same scoring scheme.

3. **Substitution matrices:**

- (a) Given the BLOSUM-62 matrix (see sequence class notes or http://www.ncbi.nlm.nih.gov/IEB/ToolBox/C_DOC/lxr/source/data/BLOSUM62), find the score of the following alignment (assume this is the optimal alignment):


```

THISSEQ
THATSEQ
      
```
- (b) Repeat with the PAM-250 matrix. It can be found in Lesk, page 257, or here: http://www.ncbi.nlm.nih.gov/IEB/ToolBox/C_DOC/lxr/source/data/PAM250

4. **Multiple Sequence Alignment:** Extend the dynamic programming formula to 3 dimensions. What is the run time in this case? How many cases do we have to compare this time?

Hint: this time the matrix is cubic since instead of a 2-dimensional matrix we need to run on a cube of $m \times n \times k$ where m,n, and k are the lengths of the three sequences. Every path goes from one vertex of the cube and traveling inside the cube to another vertex. Try to count how many such paths there can be.