

Biomedical Data Science 2023: Homework Assignment 1

Due: March 6th (Monday) 11:59pm EST

Choose to do either MCDB & MBB (non-programming) or CBB & CS & S&DS (programming) assignment, depending on your academic affiliation. No late submissions will be accepted.

Submission should be done in Canvas

MCDB & MBB (non-programming)

1. (25pt) Multiple sequence alignments (MSA) cannot be efficiently handled using purely dynamic programming. Choose one existing MSA software and describe how it implements MSA. (for example Muscle, clustalW, Kalign, MView, T-coffee etc)
2. (25pt) Align the following two sequences using the Smith-Waterman algorithm (local alignment), with the following scores: Match: 2; Mismatch: -2; Gap: -1. In addition to writing out the alignment matrix, indicate the traceback and write out the final alignment.
Sequence 1: CGACTGAAGCT
Sequence 2: GACATTGACT

3. (25pt) ChIP-seq is a common method to determine protein-DNA interaction on a genome-wide scale. The exact sites of binding must be inferred from sequence reads of the DNA that is purified along with the protein of interest. Describe an algorithm for determining protein-DNA binding sites from ChIP-seq data. See the following citation for a list of example algorithms:

Wilbanks, EG, Facciotti, MT (2010). Evaluation of algorithm performance in ChIP-seq peak detection. PLoS ONE, 5, 7:e11471.

4. (25pt) Machine learning approaches have become extremely useful in the analysis of biological data. Read the paper referenced below, and answer the following questions:
Ghandi, Mahmoud, et al. "Enhanced regulatory sequence prediction using gapped k-mer features." PLoS Comput Biol 10.7 (2014): e1003711
(<https://journals.plos.org/ploscompbiol/article?id=10.1371/journal.pcbi.1003711>)

- What are the researchers trying to predict/infer?
- What information is being used for the prediction? What is the logic behind using these data?
- What preprocessing steps are used to prepare the data for machine learning?
- What is the model the researchers use?
- How do the researchers evaluate their predictions? Were they effective? What biological insight was gained?

5. (Optional extra credit 10 pts) Answer the above questions for question 4 for the following paper:

"Building a Hybrid Physical-Statistical Classifier for Predicting the Effect of Variants Related to Protein-Drug Interactions" Structure. 2019 Sep 3;27(9):1469-1481.e3. doi: 10.1016/j.str.2019.06.001.
(<https://www.sciencedirect.com/science/article/pii/S096921261930200X?via%3Dihub>)

CBB & CPSC & S&DS (programming)

The programming task is to implement the Smith-Waterman local alignment algorithm for protein sequences.

Submit a single script file including all the codes and comments. All supplementary files could be found [here](#) or in Canvas files.

Requirements

- Implement the algorithm with Python or R. You may use the templates below. All scripts must be done from scratch. Standard libraries (I/O, math-related) are allowed. Except Numpy and pandas for python, or similar auxiliary functions in R, other pre-existed libraries are **NOT** allowed.
- The program should automatically read in the similarity matrix file called "blosum62.txt" and input sequences in a file called "input.txt", where each line is a sequence.
- For a given gap penalties, the program should output the best alignment of two sequences. The default gap penalties are as follows: opening gap -2, extension gap = -1
- The output should contain (1) Sequences, where the input is shown (2) Score matrix, where the completed scoring matrix are shown in tab-delimited format (akin to the hand-drawn scoring matrix) (3) Best alignment output as well as the alignment score. Please see file sample-output.txt for details of output format. These will constitute 90% of your grade, with the remaining 10% coming from your programming style (e.g., clear comments).
- All submitted programs should run successfully and be well-commented. Please comment with a Usage line and example line (see template below for examples). We will test and execute the script exactly as you write out in the example line.
- For example, if you have the following in your code:

```
## Usage: python hw1.py -i <input file> -s <score file>  
## Example: python hw1.py -i input.txt -s blosum62.txt
```

Then running `python hw1.py -i sample-input.txt -s blosum62.txt` should give you the exact file output as sample-output.txt (by "exact" we mean using bash `diff -E -b file1 file2` result in no output). And we will run your code with `python hw1.py -i input.txt -s blosum62.txt` and the output of this will be compared to our answer (again with `diff`)

- (Optional extra credit) You can obtain extra 10% points for publishing your program into GitHub as a package. Users should be able to install your package with a simple command such as "pip install git+git://github.com/author/package.git" or "install_github("author/package")". For full extra credit, it is required to have a clear

documentation in your GitHub repository (code example, and example input and output files) as a README markdown file.

Templates

If you use Python, you could use the following code template. Recommend Python 3.x but it would be good as long as the output is the same as expected.

```
#!/usr/bin/python
__author__ = "FirstName LastName"
__email__ = "first.last@yale.edu"
__copyright__ = "Copyright 2021"
__license__ = "GPL"
__version__ = "1.0.0"

### Usage: python hw1.py -i <input file> -s <score file>
### Example: python hw1.py -i input.txt -s blosum62.txt
### Note: Smith-Waterman Algorithm

import argparse

### This is one way to read in arguments in Python.
parser = argparse.ArgumentParser(description='Smith-Waterman Algorithm')
parser.add_argument('-i', '--input', help='input file', required=True)
parser.add_argument('-s', '--score', help='score file', required=True)
parser.add_argument('-o', '--opengap', help='open gap', required=False,
                    default=-2)
parser.add_argument('-e', '--extgap', help='extension gap',
                    required=False, default=-1)
args = parser.parse_args()

### Implement your Smith-Waterman Algorithm
def runSW(inputFile, scoreFile, openGap, extGap):
    ### calculation
    ### write output

### Run your Smith-Waterman Algorithm
runSW(args.input, args.score, args.opengap, args.extgap)
```

If you use R, you could use the following code template.

```
#!/usr/bin/env Rscript

### Usage: Rscript --vanilla hw1.R <input file> <score file>
### Example: Rscript --vanilla hw1.R input.txt blosum62.txt
### Note: Smith-Waterman Algorithm

### This is one way to read in arguments in R
```

```
args = commandArgs(trailingOnly=TRUE)

if (length(args)<2) {
  stop("At least two arguments must be supplied (inputFile, scoreFile).n",
  call.=FALSE) } else if (length(args)>=2) {
  # default gap penalties
  args[3] = -2
  args[4] = -1 }

## Specifying author and email
p <- c(person("FirstName", "LastName", role = "aut", email =
"first.last@yale.edu"))

## Implement your Smith-Waterman Algorithm
runSW <- function(inputFile, scoreFile, openGap = -2, extGap = -1) {
  ### calculation
  ### write output
}

## Run the main function and generate results
runSW(inputFile=args[1], scoreFile=args[2], openGap=args[3],
extGap=args[4])
```