## BMEG3105: Data analytics for personalized genomics and precision medicine

# Lecture 3 Scribing [Sequence data and Dynamic programming], Sep 10 (Wed)

#### **Outline**

- Questions and comments
- Recap
- Sequencing methods
- Sequence Processing
- Sequence alignment and similarity
- Introduction to DP
- Summary and resources

#### A. Review of Comments

#### **Student Feedback:**

- Mostly positive: "good," "interesting," "clear."
- Suggestions for more detailed explanations, especially for Python.

#### B. Recap of Types of Data and Sequence Types

• Biomedical data that we will encounter: genes expression, proteome, molecular & cellular network hospital test result, travel history, lifestyle, social data, etc.

• Types of data:

Sequential data

Data matrix

o Spatial data

Temporal data

Graph/Network data

Text data

Multimodality data

• The purpose of learning programming language: program is like a communication software connects between humans and computers. We need to learn the language to let the translator (Python) translate the programming codes to the codes that machines know.

#### **C. Sequence Data: Foundation of Genomics**

#### Why Sequence Data?

• Follows the central dogma (DNA → RNA → Protein).

• We have to decode the genetic sequence to know the genetic information (some genes may represent a higher possibility of having a disease).

 Phenotype arises from both genotype and environment. And genotypes are determined by the DNA sequence.

#### What Constitutes Sequence Data?

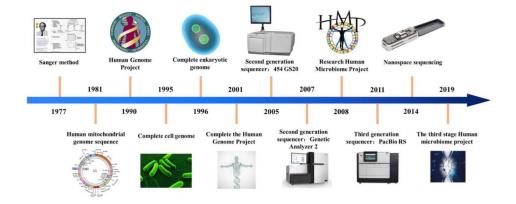
• **DNA sequence:** A, T, C, G (double-stranded, ~3 billion bp)

RNA sequence: A, U, C, G

• **Protein sequence:** 20 amino acids; analyzed with multiple sequence alignment

#### **GETTING Sequences**

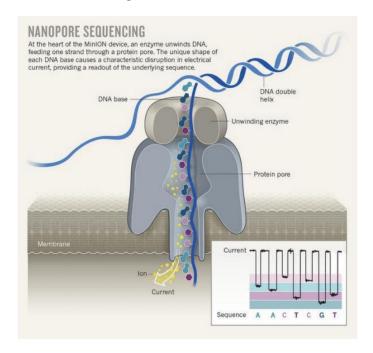
• **DNA/RNA sequencing:** Rapidly developing, moving toward long reads.



\*\*\* the first generation can only sequence ~200bp and it uses radioactive markers (low efficiency). \*\*\* the Human Genome Project sequences the whole genome of a person.

#### • Nanopore sequencing:

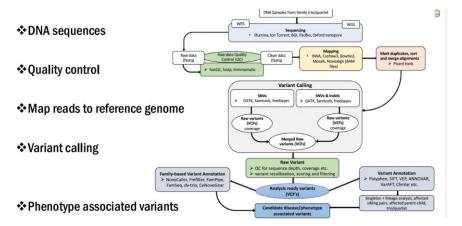
 While each DNA base in the genome passes through the chemical pore, it gives an unique current change.



- o It supports long reads (up to 4Mb) and the error rate is about 1% now, while PacBio is about 0.001%.
- **Protein sequencing:** Uses mass spectrometry—fragments are weighed and assembled like a jigsaw puzzle.

#### D. Processing and Using Sequence Data

#### General Flow to Deal with DNA Raw Data:



#### For Protein:

- Sequence comparison
- Multiple sequence alignment
- Sequence similarity helps predict structure and function, or infer evolutionary homology

#### E. Sequence Comparison & Alignment

**Definition: Sequence Similarity= The alignment score of the best alignment.** 

#### \*\*\*Why Compare Sequences?

• Determine similarity, infer functions, or discover evolutionary mechanisms by identifying the conservative/non-conservative regions, etc.

#### **Pairwise Alignment**

- Align two sequences to maximize similarity.
- Score alignment via:
  - Match
  - o Mismatch (substitution) [may be caused by mutations]

#### **Sample Scoring Matrix & Example**

	Α	С	G	Т
Α	2	-7	-5	-7
С	-7	2	-7	-5
G	-5	-7	2	-7
Т	-7	-5	-7	2

### Gap penalty = -10

- Matches score +2, mismatches -5 or -7, gap penalty -10
- Example:

AGGCCG and ATGC\_G

Alignment score: 2 + (-7) + 2 + 2 + (-10) + 2 = -9

#### F. Introduction to Dynamic Programming (DP)

#### **Enumeration Problem**

• The number of possible alignments grows exponentially with sequence length (e.g. for a 300bp sequence, the combination is  $7*10^{88}$ ; more than atoms in the universe).

#### **DP Solution Overview**

#### Analogy of DP: finding the cheanest way from KAUST to CUHK\_\_\_

independently

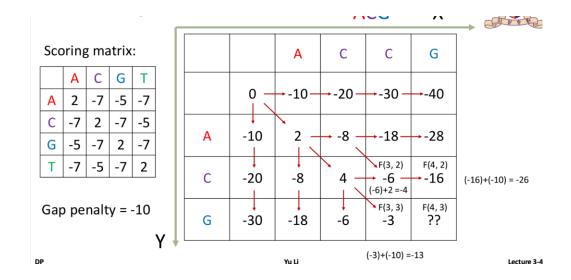
 Adopt the recursion method, reduce the combination of the above example to 300\*300=90000 (efficiency improvement).

#### • Steps in DP:

- o Break the alignment problem into smaller subproblems.
- o Solve these subproblems recursively and optimally.
- o Build up to the global solution.

#### Sequence Alignment with DP: Example in class

• Input: Aligning "ACCG" and "ACG" using the given scoring matrix and gap penalty

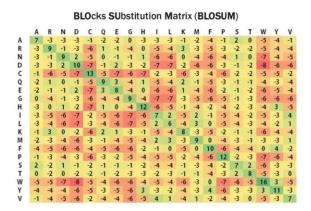


\*\*\*Notes

- Keep the arrows so that we can trace back the original path, also we can identify if there is more than one way to achieve the result.
- To simply put, both the horizontal and vertical moves just add the gap penalty. e.g. if penalty gap is -10, moving rightward a frame just needs to add (-10).

#### H. Summary & Links

- Sequence alignment with scoring matrices quantifies similarity.
- Dynamic programming offers a computationally feasible approach which the calculation to n^2.
- Some other score matrix can be used e.g. BLOSUM



- Website for sequence alignment: <a href="https://www.ebi.ac.uk/Tools/psa/emboss-needle/">https://www.ebi.ac.uk/Tools/psa/emboss-needle/</a>
- We can also use Biopython to perform alignment in Python https://www.ebi.ac.uk/Tools/psa/emboss\_needle/