

BMEG3105 (24/25 fall) - Data analytics for personalized genomics and precision medicine | Week 2 Short Lecture: Sequence data & dynamic programming

Lecturer: Yu LI (李煜) from CSE | Liyu95.com, liyu@cse.cuhk.edu.hk | Friday, Sept 13th 2024

Scriber: Yunwen ZHANG (1155173743)

From last lecture

[Comments on survey inputs](#)

[Last lecture recap](#)

Today's topic: [dynamic programming](#)

[Analogy: Cheapest flight problem](#)

[Resolving a simple example of sequence alignment using DP](#)

[Part 1: finding the optimal score](#)

[Part 2: finding the optimal alignment \(the path\)](#)

From last lecture

Comments on survey inputs

- On discrepancy in perceived difficulty, pace and detail
 - advised that students pay effort to understand the **concepts**, as staying scatterbrained throughout the time leads to bad performance at exams;
 - advised that students show up at the **revision lectures** where he will talk about what the exams will cover.

Last lecture recap

- Data types - the 6 basic types + **multi-modality data** + **unknown data** (i.e. data type unknown before we see the data)
- Python programming - way of communication with computer, essentially = a messaging app + a translating app
- Sequence data (DNA, RNA, protein), what we can do to them (sequencing & sequence alignment), & what we can obtain from them (similarity → possible relations between species)
 - Pairwise alignment: when facing long sequences, possible alignments too much → enumeration fails. Therefore here we are in front of dynamic programming.

Today's topic: dynamic programming

Analogy: Cheapest flight problem

- Problem nature: optimization (comparing)
- Constraint: no direct flight, or breaking is cheaper
- problem reduced to: finding cheapest total transportation cost between routes corresponding to each transfer place - i.e.,

$$\text{Cheapest (KAUST, CUHK)} = \text{Cheapest} \left\{ \begin{array}{l} \text{Cheapest (KAUST, Qatar)} + \text{Cheapest (Qatar, CUHK)} \\ \text{Cheapest (KAUST, Dubai)} + \text{Cheapest (Dubai, CUHK)} \\ \text{Cheapest (KAUST, GZ)} + \text{Cheapest (GZ, CUHK)} \end{array} \right.$$

- such **breaking down of a problem into sub-problems** is called **dynamic programming** (DP).

Resolving a simple example of sequence alignment using DP

- Problem:

Scoring matrix:

	A	C	G	T
A	2	-7	-5	-7
C	-7	2	-7	-5
G	-5	-7	2	-7
T	-7	-5	-7	2

Gap penalty = -10

Input sequences: ACCG
ACG

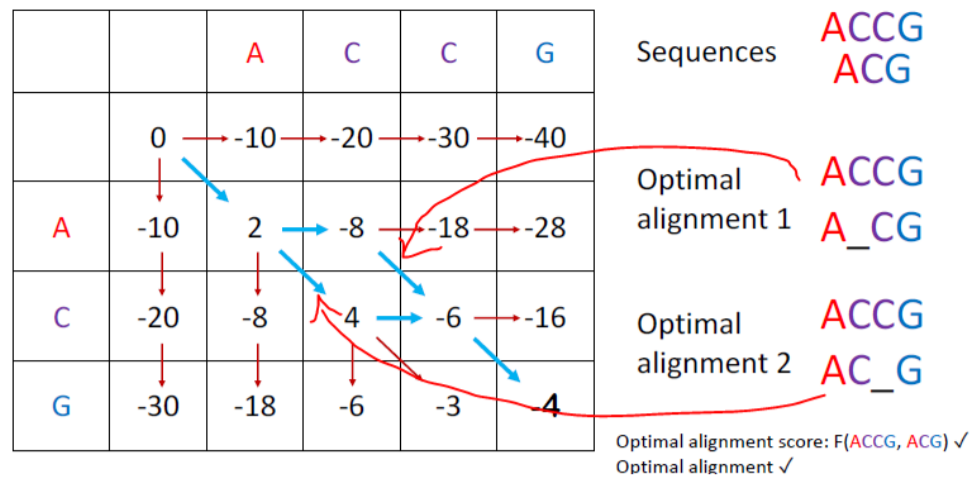
Questions:

Optimal alignment score: $F(\text{ACCG}, \text{ACG})??$

Optimal alignment??

gut feeling: should be ACCG - AC_G or ACCG - A_CG.

▼ Part 1: finding the optimal score



- Sanity check: consistent with our gut feeling.
- Reflection on computation time, for two sequences with length n
 - enumerations requires $\binom{2n}{n} = \frac{(2n)!}{(n!)^2}$ times of computation
 - DP (n by n table) requires $(n + 1)^2 \cdot 3 = 3(n + 1)^2$ times of computation