**Data Analytics for Personalized Genomics and Precision Medicine**

**Lecture 4: Dynamic Programming**

Lecturer: Yu LI (李煜) from CSE

Liyu95.com, liyu@cse.cuhk.edu.hk


15 September 2023

**Lecture outline:**

· Recap from last lecture

· Dynamic Programming


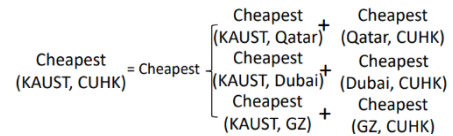**Part 1. Recap from last lecture**

Sequence Data:

· DNA sequence

   · Composed of A, T, C, G bases

   · Consists of complementary double strands


· RNA sequence:

   · Composed of A, U, C, G bases


· Protein sequence:

   · Composed of 20 amino acids

   · Multiple sequence alignment


· To find the best pairwise alignment

   · With two sequences and scoring matrix known

   · Enumeration as a solution

      · Enumerate all the possible alignments between two sequences

      · Calculate scores for all alignments

      · Choose the alignment with the highest score

      · Problem with this solution is that there are too many possible alignments

**Part 2. Dynamic Programming (DP):**

- Break down a problem into smaller sub-problems
- Solve these sub-problems optimally and recursively
- Use these optimal solutions to construct the optimal solution for the original problem

Flight problem example:

- Direct flight not always cheap
- Finite pre-destination
- Break the flight into several flight trips
- Compare the sum of flight trips for each option and choose the cheapest option

$$\text{Cheapest}_{(KAUST, CUHK)} = \text{Cheapest} \begin{bmatrix} \text{Cheapest}_{(KAUST, Qatar)} + \text{Cheapest}_{(Qatar, CUHK)} \\ \text{Cheapest}_{(KAUST, Dubai)} + \text{Cheapest}_{(Dubai, CUHK)} \\ \text{Cheapest}_{(KAUST, GZ)} + \text{Cheapest}_{(GZ, CUHK)} \end{bmatrix}$$

Similarly with Optimal alignment score:

- Each base either aligns to a gap or another base
- Alignment score equals to the sum of score for each alignment pair

Sequence alignment with DP:

- Consider the possibilities of the last pair of the alignment
- Calculate alignment score and choose the best option
- Break down the best option and repeat the last step

ACCG and ACG alignment example:

1. Consider the last base pair. ( _ represent blank)

$$F(ACCG, ACG) = \text{Best} \begin{bmatrix} F(ACC, ACG) + F(G, \_) \\ F(ACCG, AC) + F(\_, G) \\ \colorbox{yellow}{F(ACC, AC) + S(G, G)} \end{bmatrix}$$

$$F(ACCG, ACG) = \text{Best} \begin{bmatrix} F(ACC, ACG) - 10 \\ F(ACCG, AC) - 10 \\ F(ACC, AC) + 2 \end{bmatrix}$$

**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

➔ Choose the option with highest last base pair score, i.e. F(ACC, AC) + 2

2. Break down unknown part the option chosen

$$F(ACC, AC) = \text{Best} \begin{bmatrix} F(AC, AC) + F(C, \_) \\ F(ACC, A) + F(\_, C) \\ \colorbox{yellow}{F(AC, A) + S(C, C)} \end{bmatrix} = \text{Best} \begin{bmatrix} F(AC, AC) - 10 \\ F(ACC, A) - 10 \\ F(AC, A) + 2 \end{bmatrix}$$

The Formula is reduced to boundary case

$$F(AC, A) = \text{Best} \begin{cases} F(AC, \_) + F(\_, A) \\ F(A, \_) + S(C, A) \\ \colorbox{yellow}{$F(C, \_) + S(A, A)$} \end{cases}$$

$$= \text{Best} \begin{cases} -20 - 10 = -30 \\ -10 - 7 = -17 \\ 2 - 10 = -8 \end{cases}$$

The optimal alignment is ACCG AC_G or ACCG A_CG

Table representation of the method:

| | Gap | A | C | C | G |
|---|---|---|---|---|---|
| Gap | 0 | -10 | -20 | -30 | -40 |
| A | -10 | 2<br><br>F(A, A) | -8 | -18 | -28 |
| C | -20 | -8 | 4<br><br>F(A, A)+F(C, C) | -6<br><br>F(AC, AC)+F(C, _) /<br><br>F(AC, A)+F(C, C) | -16 |
| G | -30 | -18 | -6 | -3 | -4<br><br>F(ACC, AC)+F(G,G) |

Optimal alignment 1:

ACCG

A_CG

Optimal alignment 2:

ACCG

AC_G

- The arrows preserve the path information
- Arrows in blue highlight the direction that gives the best alignment score
- Trace back the arrows to get the optimal alignment