

**Data analytics for personalized genomics and precision medicine****Lecture3 scribing**

Lecturer: YuLI

Wednesday, 14 September 2022

Scriber: Kam Hei Man 1155144587

**Sequences in biology:**

- DNA sequence, RNA sequence, protein sequence

**Nanopore sequencing:**

- DNA goes through a chemical pore and read the sequence by detecting different signal given by different base.
- +ve: no need to break sequence but able to read very long sequence.

**Protein sequencing:**

- Step 1. Break long sequence into short pieces.
- Step 2. Each short pieces can be determined by their weight (mass spectrometry).
- Step 3. Assemble the short pieces back to the raw sequence.

**Raw data (DNA sequence) handling:**

- Quality control
- Map reads to reference genome
- Variant calling
- Phenotype associated variants

**Sequence alignment:**

- To compare more than 1 sequence by finding the similarity between sequences.
- Example: Pairwise sequence alignment.

**Pairwise sequence alignment:**

- Step 1. Arrange 2 sequences.
- Step 2. Maximize the similarity by inserting gap.
- Step 3. Calculate the maximum sequence alignment score (i.e. score that represent the alignment with the highest similarity) by the defined scoring matrix.

## Scoring method:

- Method 1: Enumeration
  - Problem: Too many possibility !!!
- Method 2: Dynamic programming

## Dynamic programming:

- Reduce big problem to smaller sub-problems
- Aim: reduce all problem to boundary case so that values are known and scoring matrix can be used
- Example:

### Method 1

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

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

$F(\text{G}, \_) = F(\_, \text{G}) = \text{gap penalty}$   
 $S(\text{G}, \text{G}) = \text{match}$



According to scoring matrix:

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



Solve sub-problem:

$$\begin{aligned}
 & \text{F(ACC, AC)} = \text{Best} \begin{cases} \text{F(AC, AC)} + \text{F(C, \_)} \\ \text{F(ACC, A)} + \text{F(\_, C)} \\ \text{F(AC, A)} + \text{S(C, C)} \end{cases} \\
 & \quad \quad \quad \downarrow \\
 & \text{F(ACC, AC)} = \text{Best} \begin{cases} \text{F(AC, AC)} + (-10) \\ \text{F(ACC, A)} + (-10) \\ \text{F(AC, A)} + 2 \end{cases} \\
 & \quad \quad \quad \downarrow
 \end{aligned}$$

All are boundary cases !!!

$$\begin{aligned}
 & \text{F(AC, A)} = \text{Best} \begin{cases} \text{F(AC, \_)} + \text{F(\_, A)} \\ \text{F(A, A)} + \text{F(C, \_)} \\ \text{F(A, \_)} + \text{S(C, A)} \end{cases} \\
 & \quad \quad \quad \downarrow \\
 & -8 = \text{F(AC, A)} = \text{Best} \begin{cases} (-20) + (-10) = -30 \\ 2 + (-10) = -8 \\ (-10) + (-7) = -17 \end{cases} \quad \text{Best= Highest score= -8 (-30 < -17 < -8)}
 \end{aligned}$$

.....

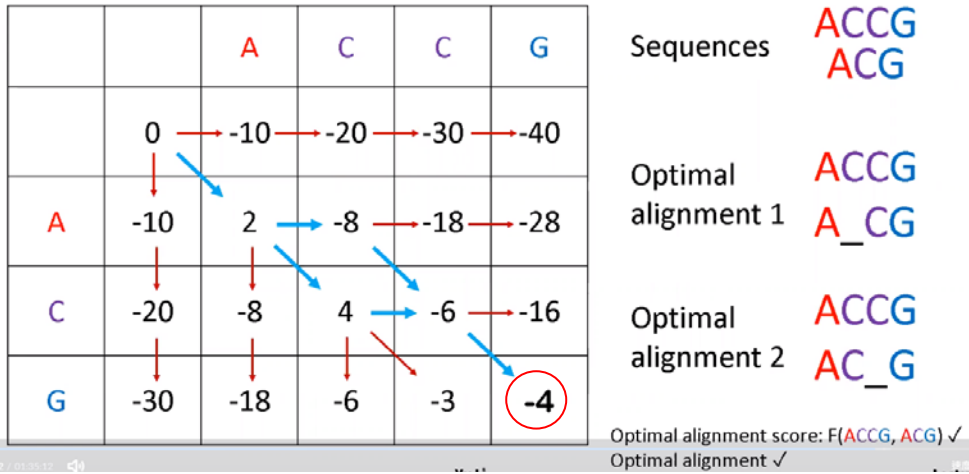
Method 2

Table representation:

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



- Step 1. Fill in the table according to the scoring matrix.
- Step 2. Draw arrows that have the smallest value difference between the previous one and the one on  $\downarrow \rightarrow \searrow$  direction.
- We can find:
  - Arrow = the alignment arrangement of 2 sequences
  - -4 = sequence alignment score