

Lecture 4: Introduction Sep 13

Lecturer: Yu Li

What are the sequence data?

-DNA: Composed of A,T,C,G, complementary double strand, approximately 3 billion of these base pairs

-RNA of AUCG

-Protein sequence: Usually composed of 20 amino acids Ø Multiple sequence alignment

How can we find the best alignment?

-Straightforward way: enumeration >> list down all the possibilities >> count the score >> find the highest

-Problem: too many alignments > need dynamic programming

How do we do and why do we need Dynamic Programming?

-Definition of sequence similarity:

Match: A <-> A

Mismatch (Substitution): G <-> T

Gap (Insertion or deletion): C <-> _

-Sequence alignment score:

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

-Break main problems into sub-problems

-Solve the sub-problems optimally and recursively

- Utilize these optimal solutions to build the best overall solution for the initial problem

-Take a flight problem example:

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

-Direct flight is expensive

-Divide it into several connecting flights

-Compare each trip/way and choose the cheapest one

ACCG, ACG for example, there are 7 bases, so a maximum 6 bases

1. Consider the last pair

According to the scoring matrix:

$$\begin{array}{l}
 F(\text{ACCG}, \text{ACG}) = \text{Best} \left\{ \begin{array}{l} F(\text{ACC}, \text{ACG}) + F(\text{G}, _) \\ F(\text{ACCG}, \text{AC}) + F(_, \text{G}) \\ F(\text{ACC}, \text{AC}) + S(\text{G}, \text{G}) \end{array} \right. \Rightarrow F(\text{ACCG}, \text{ACG}) = \text{Best} \left\{ \begin{array}{l} F(\text{ACC}, \text{ACG}) + (-10) \\ F(\text{ACCG}, \text{AC}) + (-10) \\ F(\text{ACC}, \text{AC}) + 2 \end{array} \right.
 \end{array}$$

from 7 to 6 or 5

Highest score: $F(\text{ACC}, \text{AC}) + 2$, Devide $F(\text{ACC}, \text{AC})$ to

$$\begin{array}{l}
 F(\text{ACC}, \text{AC}) = \text{Best} \left\{ \begin{array}{l} F(\text{AC}, \text{AC}) + F(\text{C}, _) \\ F(\text{ACC}, \text{A}) + F(_, \text{C}) \\ F(\text{AC}, \text{A}) + S(\text{C}, \text{C}) \end{array} \right. \Rightarrow F(\text{ACC}, \text{AC}) = \text{Best} \left\{ \begin{array}{l} F(\text{AC}, \text{AC}) + (-10) \\ F(\text{ACC}, \text{A}) + (-10) \\ F(\text{AC}, \text{A}) + 2 \end{array} \right.
 \end{array}$$

- c c

$$\begin{array}{l}
 F(\text{ACC}, \text{AC}) = \text{Best} \left\{ \begin{array}{l} F(\text{AC}, \text{AC}) + (-10) \\ F(\text{ACC}, \text{A}) + (-10) \\ F(\text{AC}, \text{A}) + 2 \end{array} \right. = \text{Best} \left\{ \begin{array}{l} (-20) + (-10) = -30 \\ 2 + (-10) = -8 \\ (-10) + (-7) = -17 \end{array} \right.
 \end{array}$$

-8+2=-6

The optimal alignment is ACCG AC_G or ACCG A_CG

The table represents the method:

		A	C	C	G
	0	-10	-20	-30	-40
A	-10	2	-8	-18	-28
C	-20	-8	4	-6	-16
G	-30	-18	-6	-3	-4

$F(A, A)$
 $F(A, A) + F(C, C)$
 $F(AC, AC) + F(C, -)$ or $F(AC, A) + F(C, C)$
 $F(ACC, AC) + F(G, G)$

Follow the two paths,

		A	C	C	G
	0	-10	-20	-30	-40
A	-10	2	-8	-18	-28
C	-20	-8	4	-6	-16
G	-30	-18	-6	-3	-4

A/A
 C/C
 C/C
 G/G

Two optimal alignments are ACCG A_CG and ACCG AC_G