BMEG3105: Data Ana for Personalized Geno

2025-2026 Semester 1

Lecture -- 4 Assembly & Mapping

Lecturer: Professor Yu Li

1. Review of DP (Dynamic Programming)

A. main goal: find the DNA base pair with the minimum value of alignment

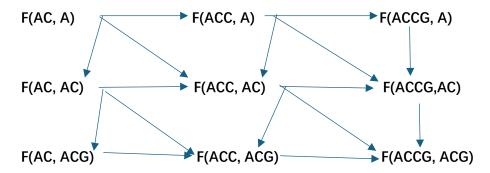
score. E g: F (ACCG, ACG)

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

Scoring matrix (gap -10)



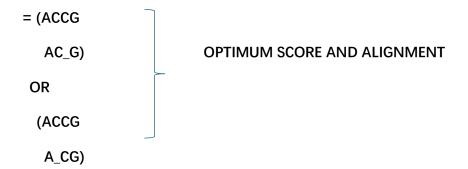
- B. Construct a DP Table.
- 1. Break it into smaller sub-base.



2. Construct the table according to the sub-base and scoring matrix.

		А	С	С	G
	0	-10 (A _) =-10(GAP)	-20	-30	-40
Α	-10 [*]	1 2	-8	-18	-28
С	-20	-8 (A_ AC) =2-10 =-8	4	-6 (ACC AC_) =2+2-10 =-6	-16
G	-30	-18	-6	-3	-4

The one we care about



C. Summary.

Good: -- Can identify sequence similarity

- --can find the optimum alignment
- --easy way to solve the sequence recursively
- --DP table shows a clear sub-sequence and construction path for understanding

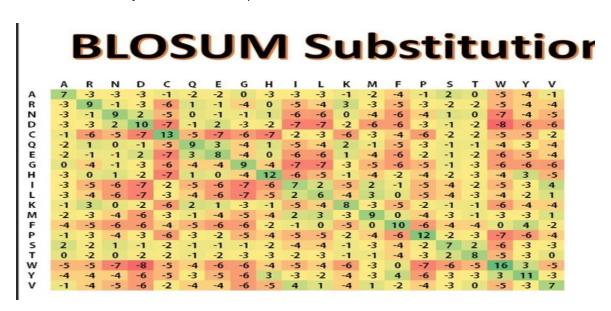
Limitation: --The calculation becomes complicated when the length of sequence is large



No. of possible alignment=(2n)!/(n!)*2

Difficult to construct a DP table

- --For global alignment, we only care about the bottom right corner one.
- --Sometimes, there is mismatch which is because of mutation (sequence change)
- -- Gap occur from insertion, deletion, gene duplication. (Gene mutation)
- --The scoring matrix may be different in different sources. (depends on how we define the similarity between two sequence)



D. Appendix

--Online Website for sequence alignment

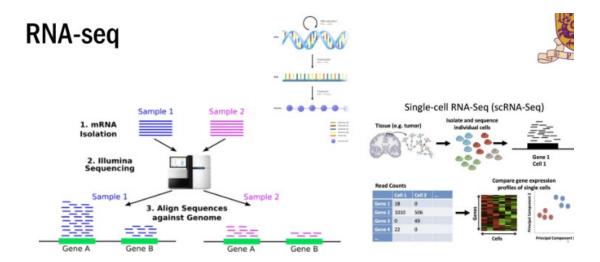
https://www.ebi.ac.uk/Tools/psa/emboss_needle/

https://biopython.org

3. Assembly and mapping.

A. Reasons for sequence data

- 1. Central dogma (DNA—RNA—Protein)
- 2. DNA sequence contains genetic information
- 3. Phenotype is related to Genotype and environment, and genotype is determined by the sequence
- --genetic variation between all people only accounts for 0.01%
- --1% of the gene is responsible for protein coding
- --gene expression different make the phenotype different



Sequence data → data matrix

1. Map the short read to the genome

2. Count the no. of reads to construct a gene expression matrix

B. Genome Assembly

Why need it?

→ Human genome is very long (~3*10^9 bp), and we need to get the genome from the long reading.

Principle: The 200bp short reads have some overlap regions, we can assembly them by assembly the overlap regions and get the interested genome.

Eg. TAA,AAT,ATG,TGC,GCC,CCA,CAT,ATG,TGG,GGA,GAT,ATG

→TAA

ATG

TGC······

C. Mapping

.

Main goal: find out the number of differences between different base in the genome during sliding each other. And find out the least differences match.

TAATGCCATGGATG

TAA,CCA,GAT,GCC,CCA,ATG

TAATGCCATGGATG

first time

CCA

2

TAATGCCATGGATG

CCA

starts from the second neuclotide and see the difference

→ TAATGCCATGGATG

CCA

233320233233

--Finally, we find that the least difference is in the 6^{th} place which has 0 different.

3. The Upcoming Topic

-- Data exploration and data cleaning