

# Sciz - Assignment - V

## MSA & phylogeny

NAGA MANIOLAR

2021101128

### MSA

①

MSA is a tool used to align  $\geq 3$  sequences simultaneously.

#### Applications:

- ~~It~~ used for phylogenetic analysis, functional annotation, structural analysis and identification of regulatory elements
- It provides more information than pairwise alignment as it identifies conserved regions & variations across multiple sequences
- This information can help identify functional domains, structural features and regular elements that may not be evident from pairwise comparisons.
- MSA provides a more accurate representation of the evolutionary relationships between sequences
- MSA can help identify convergent & divergent evolution across different lineages.

(2)

The sum-of-pairs (sp) score is a measure of the quality of a MSA, calculated by summing the scores of all pairs of aligned residues in the MSA

However it has limitations, as it does not consider gaps and the variability of residues in columns.

Additionally it doesn't consider the importance of certain residues for protein functions.

It also tends to overweight the contributions of differences.

Other Techniques A tree-form technique can be used to score. In this technique, the MSA software finds a phylogenetic tree showing the links among the sequences, the total lengths of the tree branches can be computed using the substitutions in the MSA column.

A simpler tree, with one of the <sup>n</sup> sequences serving as the ancestor of all the others, can be used instead.



③

## Progressive Alignment Approach

### Steps

- ① All sequences are pairwise aligned to generate a matrix of pairwise similarity scores.
- ② A phylogenetic tree is constructed based on the pairwise similarity scores.
- ③ A guide tree is constructed based on the phylogenetic tree.
- ④ Sequences are aligned one by one, starting from the most closely related sequences and gradually adding more distant sequences using the guide tree as a guide.
- ⑤ The alignment is refined to improve its quality.

### Drawbacks

- NOT globally optimal
- Errors are propagated throughout the final result
- $O(n^3)$  time complexity makes it unsuitable for datasets with a large number of sequences.

The Related Neighbour Joining Technique, which lowers the constraints for linking tree nodes can lessen the shortcomings of the progressive

alignment approach as a result, the time complexity is reduced to  $O(n^2 \log n)$ .

(ii)  $L=50$

The alignment of  $N$  sequences takes  $= (2L)^{N-2}$   
 $= 10^{2N-4}$

$$\begin{aligned}\text{Time} &= 5 \text{ billion years} \\ &= 5 \times 10^9 \text{ years} \\ &= 5 \times 10^9 \times 365 \times 86400 \text{ sec} \\ &= 15768 \times 10^{13} \text{ s}\end{aligned}$$

$$10^{2N-4} = 15768 \times 10^{13}$$

$$\Rightarrow 2N-4 = 17.19$$

$$\boxed{N = 10.595}$$

$$\lfloor N \rfloor = 10$$

$\therefore$  The computer can align  
sequences in 5 billion years



⑥

~~S<sub>1</sub> GATTC A~~

match — +1

mismatch — -1

indel — -1

(S<sub>1</sub>, S<sub>2</sub>, S<sub>3</sub>, S<sub>4</sub>)

$$\sqrt{4 \times 2 = \frac{4 \times 2}{2} = 2}$$

if possible pairs, we need to do pairwise alignment.

(S<sub>1</sub>, S<sub>4</sub>), (S<sub>1</sub>, S<sub>2</sub>), (S<sub>1</sub>, S<sub>3</sub>), (S<sub>2</sub>, S<sub>3</sub>), (S<sub>2</sub>, S<sub>4</sub>), (S<sub>3</sub>, S<sub>4</sub>)

① S<sub>1</sub> and S<sub>2</sub>

	G	A	T	T	C	A
G	0	-1	-2	-3	-4	-5
A	-1	0	-1	-2	-3	-4
T	-2	-1	0	0	-1	-2
C	-3	-1	-1	0	0	-1
T	-2	-1	0	0	-1	-2
G	-1	-2	-3	-4	-5	-6
A	-2	-3	-4	-5	-6	-7

GATTC A  
A-TCTGA

score = 1

mismatches = 3

② S<sub>1</sub> and S<sub>3</sub>

	G	A	T	T	C	A
G	0	-1	-2	-3	-4	-5
A	-1	0	-1	-2	-3	-4
T	-2	-1	0	0	-1	-2
A	-2	0	-1	-2	-3	-4
T	-3	-1	0	0	-1	-2
A	-4	-2	-1	-2	-3	-4
T	-5	-3	-2	-3	-4	-5
T	-6	-4	-3	-4	-5	-6

GATTC A  
GATAT A

GATTC A  
GATATT

score = 1

mismatches = 3

8

$S_1$  and  $S_4$

	G	A	T	T	C	A
G	0	-1	-2	-3	-4	-5
T	-1	0	-1	0	-2	-3
C	-2	-1	0	0	0	-1
A	-3	-2	-1	-1	0	0
G	-4	-3	-2	-2	-1	-1
C	-5	-4	-3	-3	-1	-2

Score = -1  
 GATTC - A  
 G - TCAAC

Score = -2  
 mismatch = 5

9

$S_2$  and  $S_3$

	G	T	C	T	G	A
G	0	-1	-2	-3	-4	-5
A	-1	0	-1	-2	-3	-4
T	-2	0	-1	-2	-3	-4
C	-3	-1	0	0	-1	-2
G	-4	-2	0	0	-1	-2
T	-5	-3	-1	-1	0	-1
C	-6	-4	-2	-2	0	-1

G-TCTGA  
 GATATT

Score = -1  
 mismatch = 4

10  $S_2$  and  $S_4$

	G	T	C	T	G	A
G	0	-1	-2	-3	-4	-5
T	-1	0	-1	-2	-3	-4
C	-2	0	1	2	-1	-2
A	-3	-1	1	3	2	1
G	-4	-2	0	2	2	1
T	-5	-3	1	1	1	3
C	-6	-4	-2	0	0	2

GTCTGA  
 GCTAGC

Score = 2  
 mismatch = 2



⑥  $S_3$  and  $S_4$

	G	A	T	A	T	T
G	-1	-2	-3	-4	-5	-6
T	-2	0	1	0	-1	-2
C	-3	-1	-1	0	-1	-2
A	-4	-2	0	-1	1	0
G	-5	-3	-1	-1	0	-1
C	-6	-4	-2	-2	-1	-1

GAT-ATT

G-TCAGC

Score = -1

mismatch = 4

Score (pairwise) Matrix

	$S_1$	$S_2$	$S_3$	$S_4$
$S_1$				
$S_2$	1			
$S_3$	1	-1		
$S_4$	-2	2	-1	

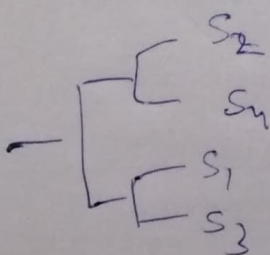
distance (pairwise)

	$S_1$	$S_2$	$S_3$	$S_4$
$S_1$			0	2
$S_2$	3		1	1
$S_3$	③	4		1
$S_4$	5	②	4	

cluster division

$S_2$  and  $S_4$

$S_1$  and  $S_3$



# Final MSA Alignment

S<sub>2</sub> G - T C T G A

S<sub>4</sub> G - T C A G C

S<sub>1</sub> G A T - T C A

S<sub>3</sub> G A T A T - T

## SpB score

col 1 = 6

col 2 = -1 -1 -1 -1 +1 = -4

col 3 = 6

col 4 = -4

col 5 = 0

col 6 = -4

col 7 = -4

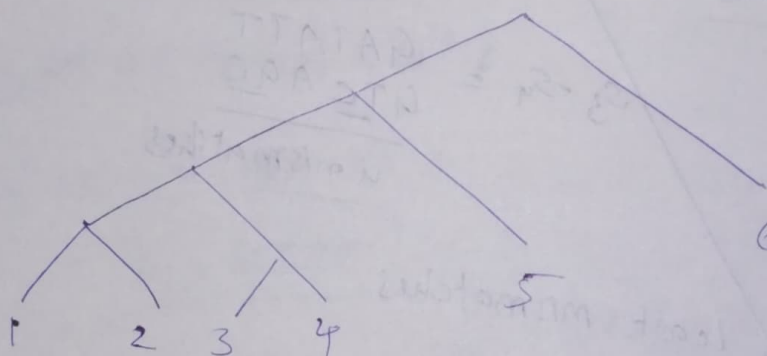
∴ Total score = -4



7) Given

	Site			
Species	1	2	3	4
1	T	C	A	A
2	G	C	A	T
3	T	T	T	T
4	G	A	T	A
5	G	A	A	C
6	A	T	A	G

The tree follows =  $((((1,2),3),4),5),6)$



$$(1,2) = \begin{Bmatrix} G \\ T \end{Bmatrix} \text{CA} \begin{Bmatrix} A \\ T \end{Bmatrix}$$

$$(3,4) = \begin{Bmatrix} T \\ G \end{Bmatrix} \begin{Bmatrix} A \\ T \end{Bmatrix} \text{T} \begin{Bmatrix} A \\ T \end{Bmatrix}$$

$$((1,2),3,4) = \begin{Bmatrix} T \\ G \end{Bmatrix} \begin{Bmatrix} A \\ C \\ T \end{Bmatrix} \begin{Bmatrix} A \\ T \end{Bmatrix} \begin{Bmatrix} A \\ T \end{Bmatrix}$$

$$(((1,2),3,4),5) = \text{GAA} \begin{Bmatrix} A \\ C \\ T \end{Bmatrix}$$

$$((((1,2),3,4),5),6) = \begin{Bmatrix} A \\ G \end{Bmatrix} \begin{Bmatrix} A \\ T \end{Bmatrix} \text{A} \begin{Bmatrix} A \\ G \\ T \end{Bmatrix}$$

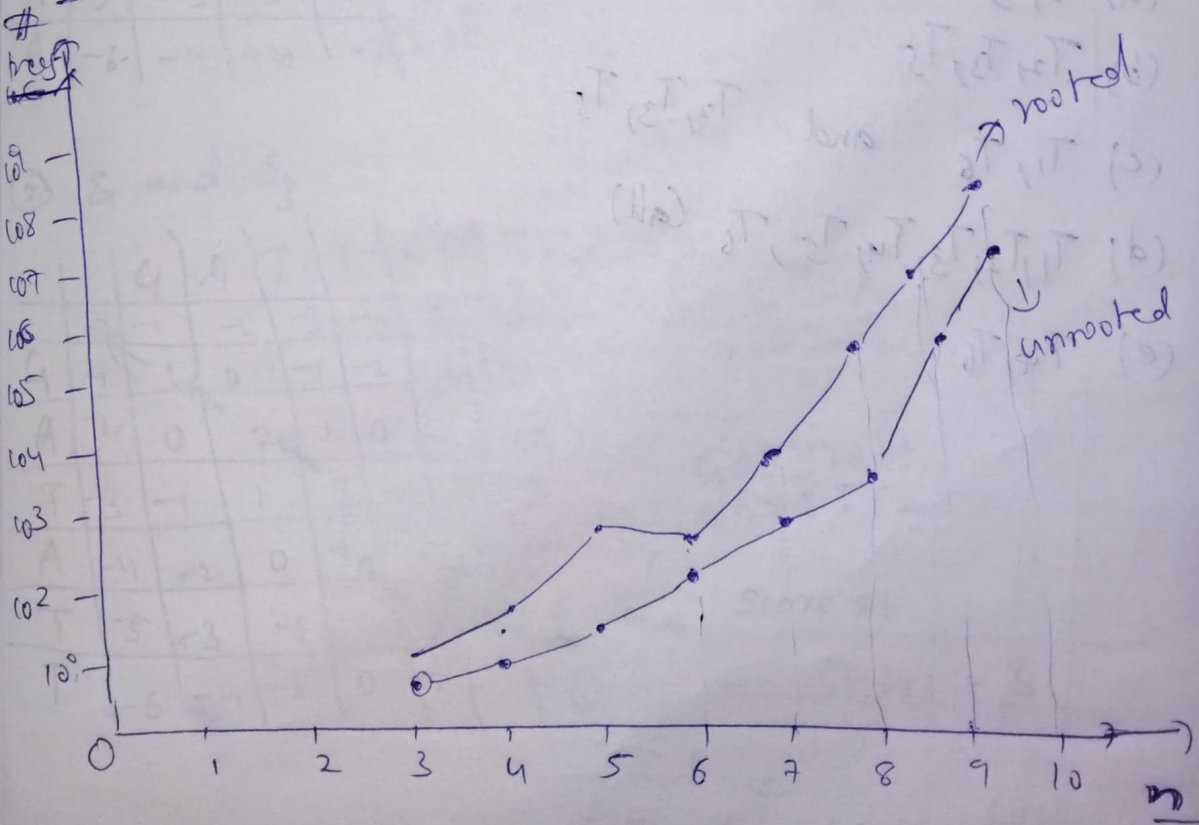
$$\therefore \text{Score} = 2+3+2+1+3 = \boxed{11}$$

8

Table 1

n	unrooted trees	rooted trees
	$\frac{(2n-5)!}{2^{n-3} \cdot (n-3)!}$	$\frac{(2n-3)!}{2^{n-2} \cdot (n-2)!}$
3	1	3
4	3	15
5	15	105
6	105	945
7	945	10395
8	10395	135135
9	135135	2027025
10	2027025	34459425

Graph





9

(a)  $T_2, T_3$

(b)  $T_2, T_3, T_5$

(c)  $T_1, T_6$  and  $T_2, T_3, T_5$

(d)  $T_1, T_2, T_3, T_4, T_5, T_6$  (all)

(e)  $T_4, T_6$