

3) I GTGCTG CACG
 II AGCTGCAAGC
 III GTGCTG CACT
 IV GTATCACA CG
 V GTATCA CAT A

a) I II III IV V

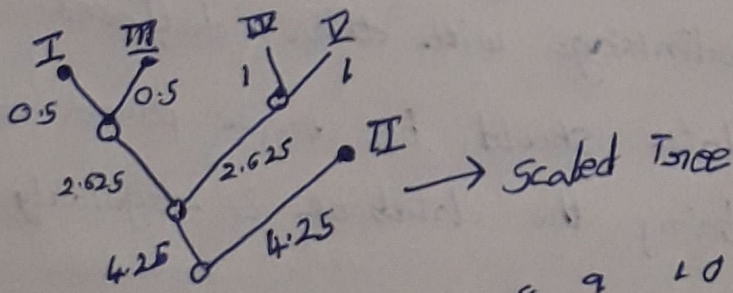
I	-				
II	9	-			
III	9	-			
IV	4	8	5	-	
V	6	8	6	2	-

I	III	II	IV	V
I	-			
II	9	-		
IV	4.5	8	-	
V	6	8	2	-

(I, II), (IV, V)

(I, II), (IV, V)
 II 8.5

I	III	II	IV	V
I	-			
II	9	-		
IV	5.25	8	-	



b)

	1	2	3	4	5	6	7	8	9	10
A	1	0	2	0	0	2	1	5	0	1
G	4	1	2	0	1	2	0	0	1	2
C	0	0	1	2	2	1	4	0	3	1
T	0	4	0	3	2	0	0	0	1	1

→ Freq. matrix

Use the formula $\ln \left[\frac{(n_{ij} + P_i)}{P_i} \right] / (N + 1)$ $N \rightarrow \text{no. of sites}$

For Ex: If $n_{ij} = 1$ then $\ln \left(\frac{1.25}{(5+1) \times 0.25} \right) = \ln \left(\frac{1.25}{1.5} \right) = -0.18$

Similarly calculate the elements of weight matrix using other values

	1	2	3	4	5	6	7	8	9	10
A	-0.182	-1.79	0.405	-1.79	-1.79	0.405	-0.182	1.252	-1.79	-0.182
G	1.04	-0.18	0.405	-1.79	-0.182	0.405	-1.79	-1.79	-0.182	0.405
C	-1.79	-1.79	-0.182	0.405	0.405	-0.182	1.04	-1.79	0.773	-0.182
T	-1.79	1.04	-1.79	0.773	0.405	-1.79	-1.79	-1.79	-0.182	-0.182

3c) Position 3: G freq: 0.4
A freq: 0.4
C freq: 0.2

Using entropy measure

$$C.S = 0.4 \ln 0.4 + 0.4 \ln 0.4 + 0.2 \ln 0.2$$

$$= -1.05492$$

Position 8: A freq: 1
 $C.S = 1 \ln 1 = 0$

3d) score: 2

	A	G	C	T	G	C	A	A	G	C
A	0	0	0	0	0	0	0	0	0	0
G	0	0	2	0	2	0	0	0	2	0
T	0	0	0	1	2	0	1	0	0	1
C	0	0	0	0	4	0	0	0	2	0
T	0	0	0	2	0	3	6	3	0	4
G	0	0	0	0	1	3	3	5	2	1
C	0	0	0	0	1	6	3	4	4	1
T	0	0	0	2	0	3	8	5	2	6
G	0	0	0	0	4	1	9	10	7	4
C	0	0	0	0	1	3	2	7	12	9
T	0	2	0	0	1	3	2	7	9	11
A	0	0	4	2	0	2	4	9	11	8
C	0	0	2	3	1	2	1	6	8	13
G	0	0	2	3	1	2	1	6	8	10

local alignment

G C T G C A C G
 G C T G C A A G

$$2+2+2+2+2+2-1+2=13$$

1) a) Life sciences, is used to get ~~info~~ information of ~~entities~~ biological entities

Computing is used for storing it in a computer so that the ~~reqd~~ ~~info~~ information is easily accessible

Maths, statistics is used to propose different methods for calculating different parameter.

Physics is used to define the forces that are involved in biological molecules.

IT (Information Tech) is used to collect large ~~sets~~ sets of data

b) i) To ~~understand~~ ^{predict} the protein structures & the amino acids involved in it, so that we can ~~find~~ find the function of that respective protein.

ii) To align ~~large~~ large no. of sequences & find the conserved subset of sequences

iii) Phylogeny Analysis

iv) To understand the freq. of mutations, insertions & deletions

v) To construct diff. databases.

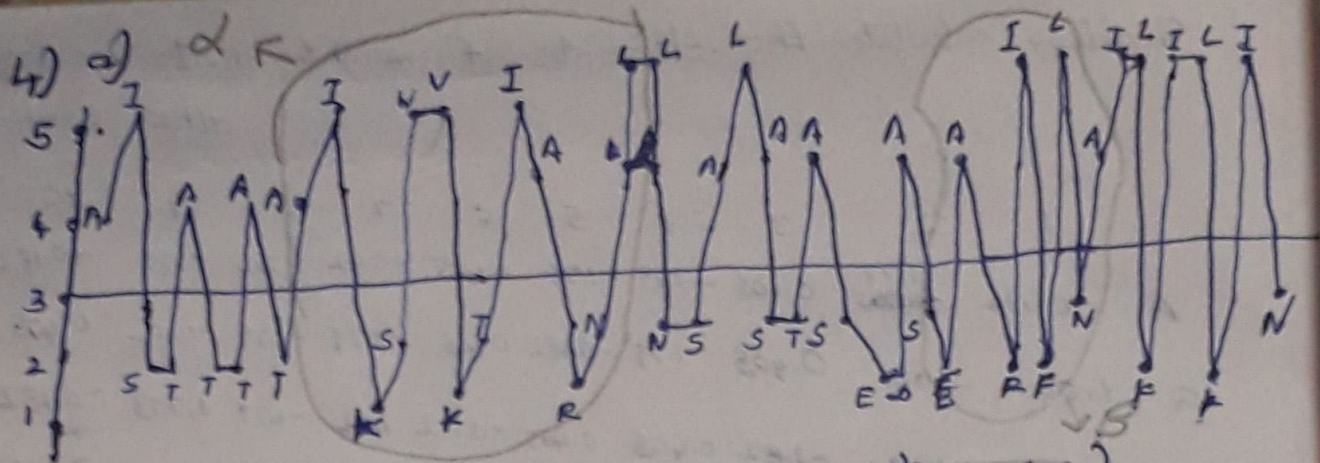
vi) Also for structure based drug design

vii) To understand protein folding & its stability

2) a) Database ^{contains} ~~has~~ collection of information in a computer readable form

d) DDBJ, EMBL

- c) i) The data ~~should be~~ ^{which is} in a proper order & format
- ii) Also organize the database ~~continuously~~ with proper definitions (to maintain)
- iii) Different routes for ~~retrieval~~ ^{recovery} of data & also to use this in other programs
- iv) Proper presentation of results.
- v) Also the links for original publications
- vi) Interlinkage with other ^y databases.
- b) i) The data should be in a proper order
- ii) Maintaining the database & regularly updating it.



alpha segment (AIKSWVKTIARNL) ~~LNLSA~~

$$a_1 = \frac{4+5+5+5+0}{4} = 4.75 \rightarrow \frac{H_{A1} + H_{N5} + H_{S9} + H_{L13} + H_{L17}}{54}$$

Similarly $a_2 = \frac{5+5+4+0}{3} = 4.67$

$$a_3 = \frac{1+1+1+2}{3} = 1.251$$

Similarly $a_4 = \frac{2+2+2+2}{3} = 2.2$

$$\frac{(H_{K3} + H_{K7} + H_{R11} + H_{L15})}{43}$$

Polar of amphipathicity = $(a_1 + a_2) - (a_3 + a_4)$
(P.O.A)

$$= \frac{11.75 + 11.75}{2} - \frac{1.25 + 2.2}{2} = 11.75 + 4.67 - (1.25 + 2.2)$$

beta segment

EARIKLNA

$$b_1 = \frac{1+1+1+2}{4} = 1.25$$

$$b_2 = \frac{4+5+5+4}{4} = 4.5$$

P.O.A = $|b_1 - b_2|$
= 3.25

$$= 11.75 + 4.67 = 6.42$$

b) alpha AIKSWVKTIARNL ~~LNLSA~~

First 6: AIKSWV = ~~1.45 + 1.00 + 1.07 + 0.79 + 1.14~~ + 1.14

Next 4: KTIA = ~~0.5 + 0.5 + 1.00 + 1.45~~ = 0.5 + 0 + 0.5 + 1

First 6: AIKSWV = $1 + 0.5 + 0.5 + 0 + 1 + 1 = 4 = 4$

Next 4: KTIA = $1.07 + 0.82 + 1.00 + 1.45 > 4$

RNLL = $0.79 + 0.73$

+ 1.34 + 1.34 < 4

Terminate

TIAR = $0.82 + 1.00 + 1.45 + 0.79 > 4$

IARN = $1.00 + 1.145 + 0.79 + 0.73 > 4$

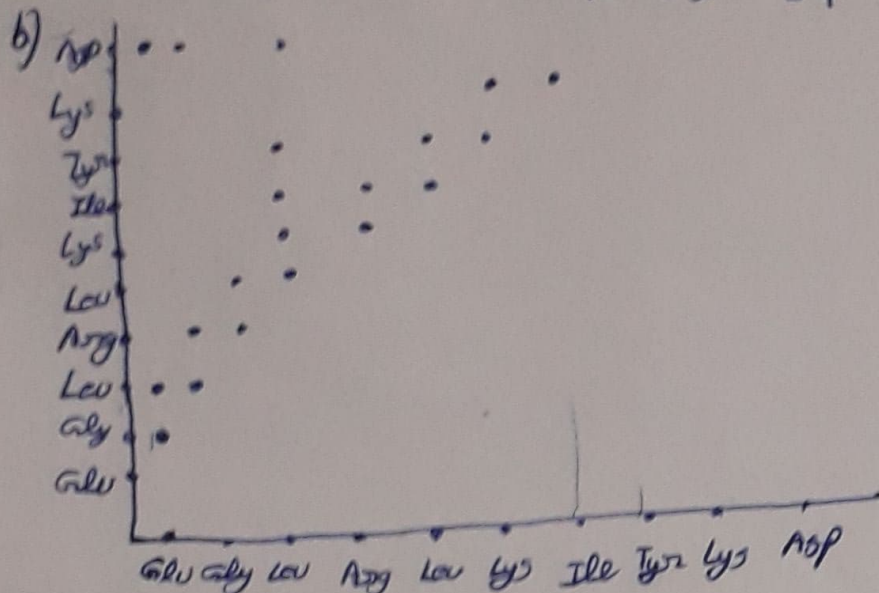
ARNL = $1.45 + 0.79 + 0.73 + 1.34 > 4$

5) a) Contact maps represents the distance b/w all possible residue pair.

Short contacts have sequence distance $< \pm 3$

medium contacts have S.D. ± 3 or ± 4

long-range contact have S.D. $> \pm 4$



c) i) S.R contacts (Glu¹¹, Gly¹²) & (Arg¹⁴, Lys¹⁶)

ii) M.R contacts (Arg¹⁴, Tyr¹⁸) & (Arg¹⁴, Ile¹⁷)

iii) L.R contacts (Gly¹², Asp²⁰) & (Leu¹³, Asp²⁰)

d) Surrounding hydrophobicity of the residue. Lys¹⁶

$$h_i = \sum n_{ij} h_j$$

$$= (2 \times 5) + (1 \times 1) + (1 \times 1) + (1 \times 4) + (1 \times 2) + (1 \times 1) + (1 \times 5) + (1 \times 1)$$

$$= 27$$

e) b)

	Predicted ^t	
	+	-
Exptl ^t		
+	7	3
-	6	44

$$TP=7 \quad FN=3$$

$$FP=6 \quad TN=44$$

$$\text{Sensitivity} = \frac{TP}{TP+FN} = \frac{7}{7+3} = 0.7$$

$$\text{Accuracy} = \frac{TP+TN}{TP+FN+TN+FP}$$

$$= \frac{13}{60} = 0.21667$$

$$\text{Specificity} = \frac{TN}{TN+FP} = \frac{44}{44+6} = \frac{44}{50} = 0.88$$