

offset 5

$$S_{CDERD} = 1 + 1 - 3 + 0 - 4$$

= very tiny value (indeterminate)

15

①

2

Offset position 3 of the given sequence when aligned with probe shows highest log odds score = 18.544. Indeed, this offset position shows the presence of conserved motif 'NKEDE'.

(6) Needleman Wunsch algorithm is employed to perform global sequence alignment.

The algorithm involves three stages

05

- Matrix Initialization
- Matrix fill
- Trace back.

The matrix will be initialised with the given scoring scheme, Match = 5, Mismatch = -3 and INDEL = -6.

The matrix will be filled/owned as per the following condition:

0.5

$$S_{i,j} = \max \begin{cases} S_{i,j} + 1 \\ S_{i,j-1} + 1 \\ S_{i-1,j} + S(x_i, y_j) \end{cases}$$

②

Alignment 1

INDIANA
INDIGO

Alignment 2

INDIANA
INDIGO

Alignment 3

INDIANA
INDIGO

①
=5

_____ X _____

7) (a) In CLUSTALW, global sequence alignment of given sequences are performed for which the E-value statistic is not available.

(2) (b) In order to calculate the weighted ~~to~~ sure for the given MSA, the weight for each sequence needs to be calculated from the guide tree as follows:

Let W_i be the weight for the i^{th} sequence.

$$W_1 = 0.3 + \left(\frac{0.04}{3}\right) + \left(\frac{0.06}{4}\right) = 0.3283$$

$$W_2 = 0.02 + \left(\frac{0.28}{2}\right) + \left(\frac{0.04}{3}\right) + \left(\frac{0.06}{4}\right) = 0.1883$$

$$W_3 = 0.02 + \left(\frac{0.28}{2}\right) + \left(\frac{0.04}{3}\right) + \left(\frac{0.06}{4}\right) = 0.1883$$

$$W_4 = 0.34 + \left(\frac{0.06}{4}\right) = 0.355$$

$$W_5 = 0.4$$

(15) For eg, the weight of sequence 1 (W_1) is calculated by summing up the distance from sequence (1) to the root in such a way that it is a particular distance

(14)

→ BLOSUM matrices are based on sequence identity
 ∴ larger the sequence identity greater is the divergence.

① In BLOSUM, the smallest number will be most suitable for distantly related sequence \Rightarrow BLOSUM 30

→ PAM matrices are based on evolutionary distance
 ∴ larger the evolutionary distance greater is the divergence

① In PAM, the largest number will be most suitable for aligning distantly related sequences \Rightarrow PAM 250

4) (a) To calculate the log odds score (S), we can simply add individual alignment score of the residues from the PAM 250 matrix

$$\begin{aligned} S &= S(R,R) + S(K,K) + S(R,D) + S(K,E) + S(R,R) + \\ &\quad S(K,K) + S(R,D) + S(K,E) \\ &= 6 + 5 + (-1) + 0 + 6 + 5 + (-1) + 0 \\ &= 20 \end{aligned}$$

① (b) To determine the statistical significance of an alignment of two random sequences of length equal to query & db sequence that could achieve the given score, we have to calculate the E-value (E).

From Karlin-Altschul equation,

$$E = K m n e^{-\lambda S} \rightarrow (1)$$

① K & λ are constants that depend on the scoring matrix.

$$3) (a) \quad M_{a,b} = \frac{\lambda m_b A_{a,b}}{\sum_{a \rightarrow b} A_{a,b}} \rightarrow (1)$$

$$M_{b,a} = \frac{\lambda m_a A_{a,b}}{\sum_{a \rightarrow b} A_{a,b}} \rightarrow (2)$$

where $M_{a,b} \neq M_{b,a} \rightarrow$ Mutation probability matrix

$\lambda \rightarrow$ proportionality constant that depends on the slowing matrix & accounts for the unsaturated evolutionary time scale

$m_a \neq m_b \rightarrow$ mutability of amino acids A & B, respectively

$A_{a,b} \rightarrow$ Accepted point mutation matrix

From (1) & (2), it is clear that

$$m_a \neq m_b$$

Therefore, $M_{a,b} \neq M_{b,a}$ are not symmetric

(b)

Alignment

E D G E
A N G E

Amino acids A D E G N

No of changes 1 1 1 0 1

Total composition 1 1 3 2 1

Mutability = $\frac{\text{No of changes}}{\text{Total composition}}$ 1 1 $\frac{1}{3}$ 0 1

$$0.5 \quad M_{AA} = 1 - \sum_{B \neq A} M_{BA} = 1 - \frac{\lambda m_A \sum_{B \neq A} A_{BA}}{\sum_{B \neq A} A_{BA}} = 1 - \lambda m_A$$

$$= 1 - \lambda \quad [\because m_A = 1 \text{ from the above table}]$$

(8)

$$\begin{aligned}
 f_{T,T}^1 &= 1 \times \frac{1}{2} = \frac{1}{2} ; f_{S,T}^2 = \frac{1}{2} \times 1 = \frac{1}{2} \\
 f_{A,T}^1 &= \frac{1}{3} \times 1 + \frac{1}{3} \times \frac{1}{2} = \frac{1}{3} + \frac{1}{6} = \frac{3}{6} = \frac{1}{2} \\
 f_{S,S}^2 &= \frac{1}{3} \times \frac{1}{2} = \frac{1}{6} ; f_{D,S}^2 = \frac{1}{3} \times \frac{1}{2} = \frac{1}{6} \\
 f_{D,T}^1 &= \frac{1}{3} \times 1 + \frac{1}{3} \times \frac{1}{2} = \frac{1}{3} + \frac{1}{6} = \frac{3}{6} = \frac{1}{2} \\
 f_{E,T}^2 &= \frac{1}{3} \times 1 + \frac{1}{3} \times \frac{1}{2} = \frac{1}{3} + \frac{1}{6} = \frac{3}{6} = \frac{1}{2} \\
 f_{E,S}^2 &= \frac{1}{3} \times \frac{1}{2} = \frac{1}{6}
 \end{aligned}$$

No A & T in Column 3, 4 & 5. Hence these columns are ignored.

Column 6: Cluster 1 $\Rightarrow \{T\}$; Cluster 2 $\Rightarrow \{\frac{1}{2} A, \frac{1}{2} A\}$
 Cluster 3 $\Rightarrow \{\frac{1}{3} U, \frac{1}{3} A, \frac{1}{3} L\}$

$$\begin{aligned}
 f_{A,T}^6 &= \frac{1}{2} \times 1 + \frac{1}{3} \times 1 = \frac{1}{2} + \frac{1}{3} = \frac{5}{6} \\
 f_{M,T}^6 &= \frac{1}{2} \times 1 = \frac{1}{2} ; f_{U,T}^6 = \frac{1}{3} \times 1 = \frac{1}{3} ; f_{L,T}^6 = \frac{1}{3} \times 1 = \frac{1}{3} \\
 f_{A,A}^6 &= \frac{1}{2} \times \frac{1}{3} = \frac{1}{6} ; f_{M,A}^6 = \frac{1}{2} \times \frac{1}{3} = \frac{1}{6} \\
 f_{A,U}^6 &= \frac{1}{2} \times \frac{1}{3} = \frac{1}{6} ; f_{M,U}^6 = \frac{1}{2} \times \frac{1}{3} = \frac{1}{6} \\
 f_{A,L}^6 &= \frac{1}{2} \times \frac{1}{3} = \frac{1}{6} ; f_{M,L}^6 = \frac{1}{2} \times \frac{1}{3} = \frac{1}{6}
 \end{aligned}$$

1.5 Since there are three clusters with six columns in the given alignment, total number of possible combination is given by

$$\sum_{a,b} f_{a,b} = 3 \times 6 = 18$$

0.5 Score (A,T) is given by

$$S_{A,T} = \log_2 \left(\frac{q_{A,T}}{2 P_A P_T} \right) \rightarrow (1)$$

(4)

From (1),
 In order to calculate $S_{A,T}$,
 $q_{A,T}$ (non-random model),

$$q_{A,T} = \frac{\sum_{\text{all columns}} f_{A,T}}{\sum_{\text{all columns}} f_A}$$

0.5

PA

⇒

$m \neq n \rightarrow$ length of the query & db, respectively
 $S \rightarrow$ log odds score

$$E = 0.09 \times 250 \times 250 e^{-0.229 \times 20}$$

①

$$= 57.6837$$

Lower the E value (closer to 0), greater is the significance of the alignment.

①

Based on this, the given alignment is not statistically significant.

3

5) (a) Log odds score in bits (S) = $\log_2 \left(\frac{\text{Observed frequency of}}{\text{Background frequency of}} \right)$

① Background frequency of amino acid = $\frac{1}{20} = 0.05$.

Amino acid	C1	C2	C3	C4	C5
N	3.585	1	1	1	1
K	0	3.8074	1	0	1
C	1	1	3.321	2	1
D	1	1	1	3.8074	-∞
E	-∞	1	-∞	1	+

①

2

(b) We have to calculate
 aligning the first seq
 offset.
 For example, offset 1
 D E N K C D E K D
 D E N K C

$$\Rightarrow q_{T,A} = 0.0741 \quad [\text{from previous calculation}]$$

$$\Rightarrow q_{T,D} = \frac{\sum_{\text{all columns}} f_{T,D}}{\sum_{\text{all columns}} f_{a,b}} = \frac{f_{T,D}^2}{18} = \frac{\frac{1}{2}}{18} = \frac{1}{36} = 0.0278$$

$$\Rightarrow q_{T,E} = \frac{\sum_{\text{all columns}} f_{T,E}}{\sum_{\text{all columns}} f_{a,b}} = \frac{f_{T,E}^2}{18} = \frac{\frac{1}{2}}{18} = \frac{1}{36} = 0.0278$$

$$\Rightarrow q_{T,M} = \frac{\sum_{\text{all columns}} f_{T,M}}{\sum_{\text{all columns}} f_{a,b}} = \frac{f_{T,M}^6}{18} = \frac{\frac{1}{2}}{18} = \frac{1}{36} = 0.0278$$

$$\Rightarrow q_{T,L} = \frac{\sum_{\text{all columns}} f_{T,L}}{\sum_{\text{all columns}} f_{a,b}} = \frac{f_{T,L}^6}{18} = \frac{\frac{1}{3}}{18} = \frac{1}{54} = 0.0185$$

Substitute these values in (4)

$$\downarrow P_T = 0.0278 + \frac{1}{2} \left[0.0185 + 0.056 + 0.0741 + 0.0278 + 0.0278 + 0.0278 + 0.085 \right]$$

$$= 0.0278 + 0.12525$$

$$= 0.15305$$

$$S_{A,T} = \log_2 \left(\frac{q_{A,T}}{2 P_A P_T} \right)$$

$$= \log_2 \left(\frac{0.0741}{2 \times 0.09265 \times 0.15305} \right)$$

$$= 1.3856$$

1) The Compositional Complexity (K) is given by

$$K = \frac{1}{L} \log_N \left(\frac{L!}{\prod_i n_i!} \right) \rightarrow (1)$$

where $L \rightarrow$ Length of the given sequence

$N \rightarrow$ No of different kinds of residues

$N = 4$ for nucleotide sequences

$n_i \rightarrow$ No of nucleotides of i^{th} kind

$$(i) K(\text{GCGACT}) = \frac{1}{6} \log_4 \left(\frac{6!}{n_A! n_C! n_G! n_T!} \right)$$

$$n_A = 1 \quad n_G = 2$$

$$n_T = 1 \quad n_C = 2$$

$$= \frac{1}{6} \log_4 \left(\frac{6!}{1! 2! 2! 1!} \right)$$

$$= \frac{1}{6} \log_4 \left(\frac{6 \times 5 \times 4 \times 3 \times 2 \times 1}{2! \times 2!} \right)$$

$$= 0.6243$$

$$(ii) K(\text{TATATA}) = \frac{1}{6} \log_4 \left(\frac{6!}{n_A! n_C! n_G! n_T!} \right)$$

$$n_A = 3 \quad n_G = 0$$

$$n_T = 3 \quad n_C = 0$$

$$= \frac{1}{6} \log_4 \left(\frac{6!}{3! \times 3!} \right)$$

$$= \frac{1}{6} \log_4 \left(\frac{6 \times 5 \times 4 \times 3 \times 2 \times 1}{3! \times 3!} \right)$$

$$= 0.3601$$

The sequence TATATA ($K = 0.3601$) is less complex than

the sequence GCGACT (0.6243)

(1)

From (1),

In order to calculate $S_{A,T}$, we have to calculate $q_{A,T}$ (non-random model), q_{PA} & $p_{A,T}$ (random model).

$$q_{A,T} = \frac{\sum_{\text{all columns}} f_{A,T}}{\sum_{\text{all columns}} f_{A,B}} = \frac{f_{A,T}^2 + f_{A,T}^6}{18} = \frac{\frac{3}{6} + \frac{5}{6}}{18}$$

0.5

$$= \frac{8}{108} = 0.0741$$

$$p_A = q_{AA} + \frac{1}{2} \sum_i q_{A,i} \rightarrow (2)$$

$$= q_{AA} + \frac{1}{2} [q_{A,h} + q_{A,S} + q_{A,L} + q_{A,T}] \rightarrow (3)$$

$$\Rightarrow q_{AA} = \frac{\sum_{\text{all columns}} f_{A,A}}{\sum_{\text{all columns}} f_{A,B}} = \frac{f_{A,A}^6}{18} = \frac{\frac{1}{6}}{18} = \frac{1}{108} = 0.0093$$

$$\Rightarrow q_{A,h} = \frac{\sum_{\text{all columns}} f_{A,h}}{\sum_{\text{all columns}} f_{A,B}} = \frac{f_{A,h}^1 + f_{A,h}^6}{18} = \frac{\frac{1}{2} + \frac{1}{6}}{18} = \frac{\frac{4}{6}}{18} = \frac{4}{108}$$

$$= 0.037$$

$$\Rightarrow q_{A,S} = \frac{\sum_{\text{all columns}} f_{A,S}}{\sum_{\text{all columns}} f_{A,B}} = \frac{f_{A,S}^1 + f_{A,S}^2}{18} = \frac{\frac{1}{2} + \frac{1}{6}}{18} = \frac{\frac{4}{6}}{18} = \frac{4}{108}$$

$$= 0.037$$

(5)

α - INDEL penalty
 $S(x_i, y_j)$ - Match or mismatch score for aligning x_i vs y_j .

Matrix Initialization & filling

$j \rightarrow$
 $i \downarrow$

	-	I	N	D	I	G	O
-	-	0	-6	-12	-18	-24	-30
I	-6	5	-1	4	-2	3	-3
N	-12	-1	10	-7	9	14	8
D	-18	-7	4	-2	14	18	12
I	-24	-13	-8	3	8	12	16
A	-30	-19	-4	-3	2	6	10
N	-36	-25	-9	-9	2	6	10
A	-42	-31	-20	-9	2	6	10

Trace back

	-	I	N	D	I	G	O
-	-	0	-6	-12	-18	-24	-30
I	-6	5	-1	4	15	9	3
N	-12	-1	10	-7	20	14	8
D	-18	-7	4	-2	14	18	12
I	-24	-13	-8	3	8	12	16
A	-30	-19	-4	-3	2	6	10
N	-36	-25	-9	-9	2	6	10
A	-42	-31	-20	-9	2	6	10

shared by more than one sequences, then that distance should be divided by the number of shared sequences.

Weighted average score (S) is calculated as follows:

$$= \sum_{i,j} \frac{W_i W_j S(x_i, y_j)}{N}$$

where W_i & W_j are the weights of sequences i & j , respectively.

$S(x_i, y_j)$ → Score for aligning character x_i & y_j

N → Number of partial alignments

$$= \frac{W_1 \times W_4 \times S(A, I) + W_1 \times W_5 \times S(A, V) + W_2 \times W_4 \times S(E, I) + W_2 \times W_5 \times S(E, V) + W_3 \times W_4 \times S(D, I) + W_3 \times W_5 \times S(D, V)}{6}$$

$$= \frac{0.3283 \times 0.355 \times 5 + 0.3283 \times 0.4 \times 8 + 0.1883 \times 0.355 \times (-7) + 0.1883 \times 0.4 \times (-5) + 0.1883 \times 0.355 \times (-3) + 0.1883 \times 0.4 \times (-5)}{6}$$

$$= \frac{0.5837 + 1.0506 - 0.4679 - 0.3766 - 0.2005 - 0.3766}{6}$$

$$= 0.2117 / 6$$

$$= 0.0353$$

①

⑮

-3

$$\Rightarrow q_{A,M} = \frac{\sum_{\text{all columns}} f_{A,M}}{\sum_{\text{all columns}} f_{A,B}} = \frac{f_{A,M}^6}{18} = \frac{1/6}{18} = \frac{1}{108} = 0.0093$$

$$\Rightarrow q_{A,L} = \frac{\sum_{\text{all columns}} f_{A,L}}{\sum_{\text{all columns}} f_{A,B}} = \frac{f_{A,L}^6}{18} = \frac{1/6}{18} = \frac{1}{108} = 0.0093$$

$$\Rightarrow q_{A,T} = \frac{\sum_{\text{all columns}} f_{A,T}}{\sum_{\text{all columns}} f_{A,B}} = \frac{f_{A,T}^2 + f_{A,T}^6}{18} = \frac{3/6 + 5/6}{18} = \frac{8}{108} = 0.0741$$

Substitute these values in (3)

$$\begin{aligned} \textcircled{1} P_A &= q_{AA} + \frac{1}{2} [q_{A,G} + q_{A,S} + q_{A,M} + q_{A,L} + q_{A,T}] \\ &= 0.0093 + \frac{1}{2} [0.037 + 0.037 + 0.0093 + 0.0093 + 0.0741] \\ &= 0.0093 + 0.08335 = 0.09265 \end{aligned}$$

$$P_T = q_{T,T} + \frac{1}{2} [q_{T,G} + q_{T,S} + q_{T,A} + q_{T,D} + q_{T,E} + q_{T,M} + q_{T,L}] \rightarrow (4)$$

$$\Rightarrow q_{T,T} = \frac{\sum_{\text{all columns}} f_{T,T}}{\sum_{\text{all columns}} f_{A,B}} = \frac{f_{T,T}^2}{18} = \frac{1}{18} = \frac{1}{36} = 0.0278$$

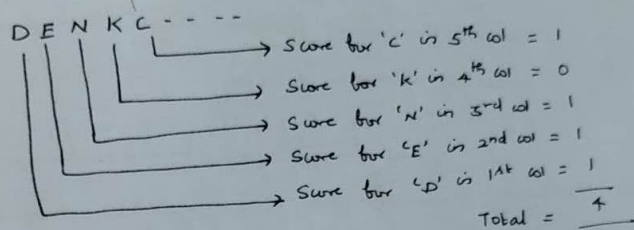
$$\Rightarrow q_{T,S} = \frac{\sum_{\text{all columns}} f_{T,S}}{\sum_{\text{all columns}} f_{A,B}} = \frac{f_{T,S}^1 + f_{T,S}^2}{18} = \frac{1/2 + 1/2}{18} = \frac{1}{18} = 0.056$$

$$\Rightarrow q_{T,G} = \frac{\sum_{\text{all columns}} f_{T,G}}{\sum_{\text{all columns}} f_{A,B}} = \frac{f_{T,G}^6}{18} = \frac{1/3}{18} = \frac{1}{54} = 0.0185 \quad \textcircled{6}$$

(b) We have to calculate the sum of log odds score after aligning the given sequence with the profile for different offset.

For example, offset 1

D E N K C D E R D (given sequence)



Sum of all scores will give the log odds score for aligning 'DENKC' with the profile at offset 1.

Similarly log odds score will be calculated for different offsets by sliding the window by one nucleotide.

offset 2

$$S_{ENKCD} = -\infty + 1 + 1 + 2 - \infty = \text{Some very long value (indeterminate)}$$

offset 3

$$S_{NKCDE} = 3.585 + 3.807 + 3.321 + 3.807 + 4 = 18.520$$

offset 4

$$S_{KCDEK} = 0 + 1 + 1 + 1 + 1 = 4$$