

SPARSH AGARWAL

9075905142

A27 (a) $P(X=x, Y=y, O=y)$

$$= P(X=x, Y=y | Y) P(O=Y)$$

$$= \left(\prod_{i=1}^{\min(3,4)} q/x_i \right) \left(\prod_{i=4}^4 p x_i y_{i-3} \right) \left(\prod_{j=\max(2,1)}^4 q/y_j \right) P(O=4)$$

$$= \left(\sum_{i=1}^3 q_{x_i} \right) \left(p_{x,y} \right) \left(\sum_{j=2}^4 q_{y_j} \right) P(O=Y)$$

$$= (q^3) (s) (q^3) (1-\theta)^3 \theta$$

$$= q_1^6 s (1-\theta)^3 \theta$$

$$= \left(\frac{1}{4}\right)^6 \cdot \left(\frac{2}{3}\right)^3 \cdot \frac{2}{3}$$

$$= 5 \left(\frac{1}{4}\right)^4 \left(\frac{1}{3}\right)^4 = 0.000010334$$

$$\textcircled{b} P(X=x, Y=y) = \sum_{o=1}^{\infty} P(X=x, Y=y, O=o)$$

$$= \sum_{o=1}^4 P(X=x, Y=y, O=o) + \sum_{o=5}^{\infty} P(X=x, Y=y, O=o)$$

⇒ Dividing the answer in 2 parts

\Rightarrow Dividing the answer in 2 parts

$$\sum_{O=5}^{\infty} P(X=x, Y=y, O=o) = \left(\prod_{i=1}^4 q_i n_i \right) \begin{pmatrix} \text{not valid} \\ i=5 \end{pmatrix} \left(\prod_{j=1}^4 q_j y_j \right) (1-o)^0$$

$$+ \quad " \quad " \quad " \quad (1-o)^1 o$$

$$+ \quad " \quad " \quad " \quad (1-o)^2 o$$

$$+ \quad " \quad " \quad " \quad (1-o)^3 o$$

$$+ \quad " \quad " \quad " \quad (1-o)^4 o + \dots \dots (1-o)^6 o]$$

$$= q^8 (1-\theta)^4 \theta [1 + (1-\theta) + (1-\theta)^2 + (1-\theta)^3 + \dots]$$

$$= q^8 (1-\theta)^4 \theta \left[\frac{1 - (1-\theta)^\infty}{1 - (1-\theta)} \right]$$

$$= q^8 (1-\theta)^4 \theta$$

⇒ First part, $\sum_{i=1}^4 P(X=x, Y=y, O=i)$

$$= P(X=x, Y=y | O=1) P(O=1) + P(X=x, Y=y | O=2) P(O=2) + P(X=x, Y=y | O=3) P(O=3) + P(X=x, Y=y | O=4) P(O=4)$$

$$= (P_{x_1 y_1} \times P_{x_2 y_2} \times P_{x_3 y_3} \times P_{x_4 y_4}) \theta + q^2 (P_{x_2 y_1} \times P_{x_3 y_2} + P_{x_4 y_3}) (1-\theta) \theta + q^4 (P_{x_3 y_1} \times P_{x_4 y_2}) (1-\theta)^2 \theta + q^6 (1-\theta)^3 \theta (P_{x_4 y_1})$$

$$= (m \times m \times s \times s) \theta + q^2 (m \times s \times m) (1-\theta) \theta + q^4 (m \times m) (1-\theta)^2 \theta + q^6 s (1-\theta)^3 \theta$$

$$= m^2 s^2 \theta + m^2 q^2 s (1-\theta) \theta + q^4 m^2 (1-\theta)^2 \theta + q^6 s (1-\theta)^3 \theta$$

~~0.000000723~~

∴ Total answer is

$$= m^2 s^2 \theta + m^2 q^2 s (1-\theta) \theta + q^4 m^2 (1-\theta)^2 \theta + q^6 s (1-\theta)^3 \theta + q^8 (1-\theta)^4 \theta$$

$$= \frac{1}{1382976} + \frac{1}{2370816} + \frac{1}{12192768} + \frac{1}{143536} + \frac{1}{331776}$$

$$= 0.000000723 + 0.000000422 + 0.000000082 + 0.000005167 + 0.000003014$$

$$= 0.000009408$$

$$P(O=4 | X=x, Y=y)$$

$$= P(O=4, X=x, Y=y) / P(X=x, Y=y)$$

\therefore taking answers from first two parts

$$= \frac{0.000010334}{0.000009408}$$

$$(d) P(O \leq 4 | X=x, Y=y)$$

using previous parts

$$P(O \leq 4 | X=x, Y=y) = \frac{\sum_{o=1}^4 P(O=o, X=x, Y=y)}{P(X=x, Y=y)}$$

$$= \frac{0.000009408 - 0.000003014}{0.000009408}$$

(e) considering that q_i^n is valid for ~~for~~ $n=1$, there are 16 possibilities of $(X=x, Y=y)$

$$\sum_{o=1}^{\infty} P(X=x, Y=y, O=o) = \sum_{o=1}^{\infty} P(X=x, Y=y, O=o) + \sum_{o=2}^{\infty} P(X=x, Y=y, O=o)$$

$$= \theta \left(\prod_{i=1}^n q_i x_i \right) \left(\prod_{i=1}^n p_i x_i, y_i \right) \left(\prod_{j=\max(2,1)}^{\pi} q_{y_j} \right)$$

$$+ (1-\theta)^0 + (1-\theta)^2 \theta + (1-\theta)^3 \theta + \dots$$

$$= \theta m + (1-\theta) \theta \left[\frac{1 - (1-\theta)^{\infty}}{1 - (1-\theta)} \right]$$

$$= \theta m + 1 - \theta$$

∴ the total probability = 1 = 16 (0_m + 1 - 0)

$$\therefore 16 \left(0_m + \frac{2}{3} \right) = 1$$

$$0_m + \frac{2}{3} = \frac{1}{16}$$

$$\cancel{0_m} \frac{m+2}{3} = \frac{1}{16}$$

$$m = \frac{3}{16} - 2$$

$$m = \frac{-29}{16}$$

A3) (b) Seed methods like Blast are very fast & reliable in a statistical sense. (very good for large sequences) whereas, DP methods like Needleman-Wunsch alignment are ~~are~~ relatively slow & computational steps increase as the square or cube of the sequence lengths.

(a) In protein alignment, we look for ~~transpos~~ transpositions, inversions, deletions, insertions & substitutions. For short proteins we consider only substitutions & insertion/deletions which are represented as ~~mismatches~~ match/mismatch & gaps respectively. But for long proteins we also consider transposition & inversion.

Some insertion & deletions may not significantly affect the structure of protein \therefore we need to look for similarities instead of just match/mismatch.

We look for partial matches (i.e. some amino ~~acid~~ acid pairs are more substitutable than others).

We represent this similarity in form of a score matrix, where we produce all possible scores for different alignments feasible.

Then we select the alignment with more score.