

Solutions - Series 4

Genomics and bioinformatics - Week 4

1 Global alignment

1.1 The table

Using the given scoring matrix M , one deduces the following rules to compute the 3 intermediate scores in each cell:

- Upper neighbour cell score - 2
- Left neighbour cell score - 2
- Upper-left neighbour cell score $\begin{cases} +1 & \text{if nucleotides are identical} \\ -1 & \text{if nucleotides are different} \end{cases}$

The highest score is then retained and the corresponding arrow is labelled. Here is the resulting table:

	-	G	A	A	T	T	C	A	G
-	0	-2	-4	-6	-8	-10	-12	-14	-16
G	-2	1	-1	-3	-5	-7	-9	-11	-13
G	-4	-1	0	-2	-4	-6	-8	-10	-10
A	-6	-3	0	1	-1	-3	-5	-7	-9
T	-8	-5	-2	-1	2	0	-2	-4	-6
C	-10	-7	-4	-3	0	1	1	-1	-3
G	-12	-9	-6	-5	-2	-1	0	0	0

1.2 Backtracking

One then applies the traceback process to the obtained table to deduce the best alignment between the two reads. The traceback begins with the bottom-right cell and is completed when the top-left cell of the table is reached. Note that several traceback paths are possible and thus in general there may be more than one optimum alignment. The best alignments are given by the paths that have the maximum score (here $F_{68} = 0$). A traceback path is highlighted below; a second one would go through the -1 at coordinate (3,4).

	-	G	A	A	T	T	C	A	G
-	0	-2	-4	-6	-8	-10	-12	-14	-16
G	-2	1	-1	-3	-5	-7	-9	-11	-13
G	-4	-1	0	-2	-4	-6	-8	-10	-10
A	-6	-3	0	1	-1	-3	-5	-7	-9
T	-8	-5	-2	-1	2	0	-2	-4	-6
C	-10	-7	-4	-3	0	1	1	-1	-3
G	-12	-9	-6	-5	-2	-1	0	0	0

1.3 Alignment

An optimal alignment is easily obtained by using the following arrow rules:

Left/Right = Deletion, Up/Down = Insertion, and Diagonal = Match.

An optimal alignment is

```

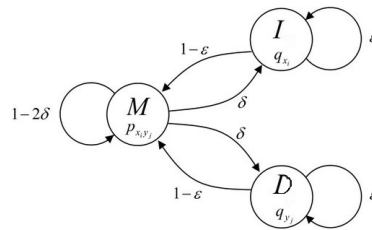
G A A T T C A G
M M M M D M D M
G A A T T C A G
|   |   |   |   |
G G A T - C - G

```

2 HMM

2.1 Finding the Emission and Transition Probabilities

The corresponding HMM is described in the following figure:



1. Transitions: from $d = \log_2(\delta)$ and $d = -2$ we deduce $\delta = 2^{-2} = \frac{1}{4} = \varepsilon$.
2. Emissions: $p(x)$ is the probability of choosing one nucleotide at random: $p(x) = p(y) = 1/4$.
Inverting the formula for $S(x, y)$, one gets $p(x, y) = \frac{1}{16} \cdot 2^{\frac{S(x, y)}{\lambda}}$:

$$p(x, x) = \frac{1}{16} \cdot 2^{\log_2(3)} = \frac{1}{16} \cdot 3 = \frac{3}{16},$$




$$p(x, y) = \frac{1}{16} \cdot 2^{-\log_2(3)} = \frac{1}{16} \cdot 2^{\log_2(1/3)} = \frac{1}{16} \cdot \frac{1}{3} = \frac{1}{48} \text{ if } x \neq y.$$

2.2 Constructing the Three Matrices (VM,VD and VI)

Then we can construct the three matrices for VM, VD and VI by the Viterbi algorithm.
Matrix V_M :

	-----	G	A	A	T	T	C	A	G
-----	0	$-\infty$	$-\infty$	$-\infty$	$-\infty$	$-\infty$	$-\infty$	$-\infty$	$-\infty$
G	$-\infty$	1	-3	-5	-7	-9	-11	-13	-13
G	$-\infty$	-1	0	-2	-4	-6	-8	-10	-10
A	$-\infty$	-5	0	1	-3	-5	-7	-7	-11
T	$-\infty$	-7	-4	-1	2	0	-4	-6	-8
C	$-\infty$	-9	-6	-3	-2	1	1	-3	-5
G	$-\infty$	-9	-8	-5	-4	-1	0	0	0

Matrix V_D :

	-----	G	A	A	T	T	C	A	G
-----	$-\infty$	-2	-4	-6	-8	-10	-12	-14	-16
G	$-\infty$	$-\infty$	-1	-3	-5	-7	-9	-11	-13
G	$-\infty$	$-\infty$	-3	-2	-4	-6	-8	-10	-12
A	$-\infty$	$-\infty$	-7	-2 	-1	-3	-5	-7	-9
T	$-\infty$	$-\infty$	-9	-6	-3 	0	-2	-4	-6
C	$-\infty$	$-\infty$	-11	-8	-5	-4	-1 	1	-3
G	$-\infty$	$-\infty$	-11	-10	-7	-6	-3	-2	-2

Matrix V_I :

	-----	G	A	A	T	T	C	A	G
-----	$-\infty$	$-\infty$	$-\infty$	$-\infty$	$-\infty$	$-\infty$	$-\infty$	$-\infty$	$-\infty$
G	-2	$-\infty$	$-\infty$	$-\infty$	$-\infty$	$-\infty$	$-\infty$	$-\infty$	$-\infty$
G	-4	-1	-5	-7	-9	-11	-13	-15	-15
A	-6	-3	-2	-4	-6	-8	-10	-12	-12
T	-8	-5	-2	-1	-5	-7	-9	-9	-13
C	-10	-7	-4	-3	0	-2	-6	-8	-10
G	-12	-9	-6	-5	-2	-1	-1	-5	-7

2.3 Deducing the Alignments

The final matrix is formed by taking the maximum value from the three elements of each matrix. Note that it is exactly the matrix of Exercise 1. Remember (annotate) which matrix the maximum came from. Then start at the bottom-right cell. If the max came from V_M , go one step back diagonally (up-left direction); if it came from V_D , go one step to the left; if from V_I , go up one step; if there are multiple choices for the maximum, this is what opens the possibility to have multiple alignments with the same probability.

	----	G	A	A	T	T	C	A	G
----	0	-2	-4	-6	-8	-10	-12	-14	-16
G	-2	1	-1	-3	-5	-7	-9	-11	-13
G	-4	-1	0	-2	-4	-6	-8	-10	-10
A	-6	-3	0	1	-1	-3	-5	-7	-9
T	-8	-5	-2	-1	2	0	-2	-4	-6
C	-10	-7	-4	-3	0	1	1	-1	-3
G	-12	-9	-6	-5	-2	-1	0	0	0

There are two possible optimal alignments (same as in Exercise 1).

G	A	A	T	T	C	A	G	OR	G	A	A	T	T	C	A	G
G	G	A	T	-	C	-	G		G	G	A	-	T	C	-	G
M	M	M	M	D	M	D	M		M	M	M	D	M	M	D	M

This procedure may seem more complicated than the Needleman-Wunsch algorithm to find the same result, but the point is that HMMs can be used in a wide range of other problems, providing efficient solutions thanks to the Viterbi algorithm. In this particular case, one can show that there is an equivalence between the recursion formulas for both algorithms.