

## Exercise sheet 7: BLAT

---

In the first step of the BLAT (Blast-like Alignment Tool) algorithm, regions that are likely to be homologous are detected. In this exercise sheet, we will investigate the search stage of the BLAT algorithm on the example of the mouse genome.

We want to determine whether a region from the human genome aligns to a part of the mouse genome. Therefore, a homologous region to our query sequence will be detected during the search stage.

We assume that:

- the human genome  $G$  is approximately 2.9 billion bases long.
- the mouse genome  $G$  is approximately 2.5 billion bases long.
- the match ratio  $M$  between homologous areas of both species is 98% for DNA and 89% for protein alignments.
- in our example we assume that homologous areas  $H$  are typically 50 bases long.
- our query sequence  $Q$  is GTCCTCGGAACCAGGACCTCGGCGTGGCCTAGCG.

### Exercise 1

For the  $K$ -mer sizes  $K = 7$  and  $K = 14$ , respectively.

**Question 1A** What is the probability of having a perfect match between a specific  $K$ -mer in a homologous region and a  $K$ -mer in the query sequence?

**Formula**

$$p_1 = M^K$$

**Solution** For  $K = 7$ :

$$p_1 = 0.98^7 = 0.8681$$

For  $K = 14$ :

$$p_1 = 0.98^{14} = 0.7536$$

**Question 1B** What is the probability that at least one non-overlapping  $K$ -mer in the homologous region matches perfectly with the corresponding  $K$ -mer in the query sequence?

**Formula**

$$P = 1 - (1 - p_1)^T = 1 - (1 - M^K)^T, \text{ with } T = \left\lfloor \frac{H}{K} \right\rfloor$$

**Solution** For  $K = 7$ :

$$P = 1 - (1 - 0.98^7)^{\lfloor \frac{50}{7} \rfloor} = 0.9999...$$

For  $K = 14$ :

$$P = 1 - (1 - 0.98^{14})^{\lfloor \frac{50}{14} \rfloor} = 0.9850...$$

**Question 1C** Calculate the number of False Positives (FPs), i.e. the number of  $K$ -mers that are expected to match by chance.

**Formula**

$$F = (Q - K + 1) \times \left( \frac{G}{K} \right) \times \left( \frac{1}{A} \right)^K,$$

$A$  : the alphabet size

$Q$  : the query length in bases

$G$  : the genome size in bases

**Solution** For  $K = 7$ :

$$F = (34 - 7 + 1) \times \frac{2500000000}{7} \times \frac{1}{4}^7 = 610351.5625$$

For  $K = 14$ :

$$F = (34 - 14 + 1) \times \frac{2500000000}{14} \times \frac{1}{4}^{14} = 13.9698$$

**Question 1D** Observe the True positive rate (TPR) and the number of False Positives for the 7-mers and 14-mers that you computed in part 1B and 1C. What observation do you make?

**Solution** When increasing the  $K$ -mer size both the TPR and the number of False Positives are reduced. But the number of False Positives reduces drastically, compared to the TPR.

## Exercise 2

In order to increase the True Positive Rate, we want to allow single mismatches when checking for exact matches.

For the  $K$ -mer sizes  $K = 7$  and  $K = 14$ , respectively.

**Question 2A** What is the probability that at least one non-overlapping  $K$ -mer in the homologous region matches perfectly with the corresponding  $K$ -mer in the query sequence? Given that we allow one mismatch.

**Formula**

$$P = 1 - (1 - (M^K + K \times M^{K-1} \times (1 - M)))^T, \text{ with } T = \left\lfloor \frac{H}{K} \right\rfloor$$

**Solution** For  $K = 7$ :

$$P = 1 - (1 - (0.98^7 + 7 * 0.98^6 * 0.02))^{\lfloor \frac{50}{7} \rfloor} \approx 1$$

For  $K = 14$ :

$$P = 1 - (1 - (0.98^{14} + 14 * 0.98^{13} * 0.02))^{\lfloor \frac{50}{14} \rfloor} \approx 1$$

**Question 2B** Calculate the number of False Positives (FPs), i.e. the number of  $K$ -mers that are expected to match by chance. Given that we allow one mismatch.

**Formula**

$$F = (Q - K + 1) \times \left( \frac{G}{K} \right) \times \left[ \left( \frac{1}{A} \right)^K + K \times \left( \frac{1}{A} \right)^{K-1} \times \left( 1 - \frac{1}{A} \right) \right],$$

 $A$  : the alphabet size $Q$  : the query length in bases $G$  : the genome size in bases**Solution** For  $K = 7$ :

$$F = (34 - 7 + 1) \times \frac{2500000000}{7} \times \left( \frac{1}{4}^7 + 7 \times \frac{1}{4}^6 \times \left( 1 - \frac{1}{4} \right) \right) = 1.3427... * 10^7$$

For  $K = 14$ :

$$F = (34 - 14 + 1) \times \frac{2500000000}{14} \times \left( \frac{1}{4}^{14} + 14 \times \frac{1}{4}^{13} \times \left( 1 - \frac{1}{4} \right) \right) = 600.7031$$

**Question 2C** What development do we see when observing the TPR and FPs for the updated formulae?

**Solution** We observe an increase in both the TPR and the number of False positives.

**Exercise 3**

Finally, we want to reduce the number of False Positive results. To that end, instead of requiring single perfect matches, we now want at least  $n$  perfect matches.

For the  $K$ -mer sizes  $K = 7$  and  $K = 14$ , respectively.

**Question 3A** What is the probability that at least 2 non-overlapping  $K$ -mer in the homologous region match perfectly with corresponding  $K$ -mers in the query sequence?

**Formula**

$$p_n = (M^K)^n \times (1 - M^K)^{T-n} \times \frac{T!}{n! \times (T-n)!} \text{ with } T = \left\lfloor \frac{H}{K} \right\rfloor P = P_n + P_{n+1} + \dots + P_T$$

**Solution** For  $K = 7$ :

$$P_2 = 0.98^{7 \times 2} \times (1 - 0.98^7)^{\lfloor \frac{50}{7} \rfloor - 2} \times \frac{\frac{50}{7}!}{2! \times (\frac{50}{7} - 2)!} \dots$$

$$\begin{aligned} P &= P_2 + P_3 + P_4 + P_5 + P_6 + P_7 \\ &= 0.000631 + 0.006925 + 0.045591 + 0.180074 + 0.395142 + 0.371601 \\ &= 0.999967 \dots \end{aligned}$$

For  $K = 14$ :

$$P_2 = 0.98^{14 \times 2} \times (1 - 0.98^{14})^{\lfloor \frac{50}{14} \rfloor - 2} \times \frac{\frac{50}{14}!}{2! \times (\frac{50}{14} - 2)!} \dots$$

$$\begin{aligned} P &= P_2 + P_3 \\ &= 0.419776 + 0.428050 \\ &= 0.847827 \dots \end{aligned}$$

**Question 3B** We can observe a decrease in the TPR. Does it still make sense to use this method? Why?

**Solution** While it is true that this also decreases the TPR. It drastically decreases the amount of False Positives, which greatly improves the overall results. In real examples, the size of the homologous regions is typically higher than 50 nucleotides that we chose for this example. When considering our example with  $H = 100$ , we would end up with a TPR close to 1.0, while having nearly no False Positives.