

Assignment Report

Paarth Gupta
15477

In this assignment we try to find out the most probable configuration (of out of the 4 discussed in class) that leads to color blindness.

I have tried to follow the similar approach discussed in the class. We have 4 characters A, C, G and T so for each one of them we create a binary column from the BWT column.

Now I have use the concept of Δ -milestones (with $\Delta = 100$).
So to keep a Δ -rank milestone array we need $\frac{\Delta}{n}$ bytes and for 4 chars we need $\frac{4\Delta}{n}$ bytes.

This is a one time process done so that we can obtain the rank of a character in $O(\Delta)$ time.

Now our aim is to find for each read either a perfect match or a match with at most 2 mismatched characters, using the rank array. We try to recursively search for the range of indices in which the suffix has an exact match.

Now we note down the count of these matches in 2 different sets i.e Red and green and further in them an which number they lie.

Finding probable configuration :-

eg for Red : $[R_1, R_2, R_3, R_4, R_5, R_6]$
 Green : $[G_1, G_2, G_3, G_4, G_5, G_6]$

$$P(\{Red, Green\} | Config_i) = \prod_{i=2}^6 \binom{R_i + G_i}{R_i} p_{config_i}^{R_i} (1 - p_{config_i})^{G_i} \quad - (1)$$

Running 3M reads was not possible system was crashing + Google Colab was also crashing. Results are obtained by partitioning and that too only an around 1.7 M reads (high density), 1M from the back and 0.7 M from the front. After analysing that backside contains much more matches.

It is observed that all the exams are present after 2.8M only.

R_i	[97	237	107.5	167.5	303.5	235]
G_i	[97	287	154.5	140.5	358.5	235]

Now using equation 1 we get :-

$P(R, G Config 1)$	$= 2.1294 \times 10^{-44}$	$\log P_1$	$= -100.55$
$P(R, G Config 2)$	$= 0$	$\log P_2$	$= -\infty$
$P(R, G Config 3)$	$= 1.266 \times 10^{-36}$	$\log P_3$	$= -82.657$
$P(R, G Config 4)$	$= 3.0427 \times 10^{-62}$	$\log P_4$	$= -141.647$

Maximum likelihood is obtained on Config 3. i.e. [33% 33% 100% 100%

Solution for select query

select query :- to find the i^{th} 1 from a binary array.

- Considering that we have an array of size $\frac{n}{\Delta}$ which stores Δ -rank milestones and will contain the indices of the corresponding ranked 1 in the original binary array.

Now as we have this Δ -rank milestone array, if we are given i using this Δ -rank milestone array we can directly jump to the index in the original array till we reach i^{th} (desired) 1.

With the i in hand, we can use the binary search as the array is sorted to find the correct location in $O(\log n)$ time.

If we analyse the technique we can observe that traversal to the required 1 may take more time than $O(\Delta)$ but in amortized setting we say that it's $O(\Delta)$. Also to improve this time complexity we use another array of size $\frac{n}{\Delta}$ which will contain a reference to the next one.

This setting has a time complexity of $O(\Delta)$ and a total space complexity of $O(2 \frac{n}{\Delta})$.