

Workshop Systems Analysis

Juan Diego Lozada Gonzalez

COD: 2022020014

Universidad Distrital Francisco Jose de Caldas

Systems Analysis

Carlos Andres Sierra

March 14, 2024

Workshop 1

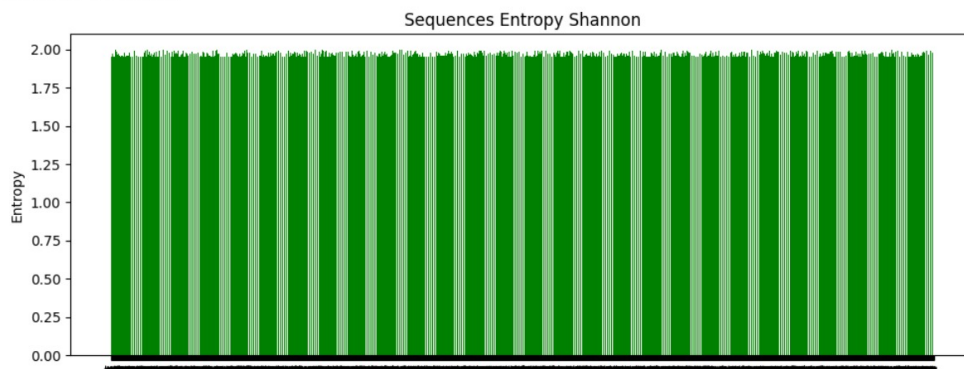
In this workshop we calculate the entropy from a sequence of nitrogen bases (A,G,T,C) with the goal to analyze why is important the entropy in the area of data, the Shannon Entropy have many uses but we are using to measure the information content in a set of symbols

$$H(X) = - \sum_{i=1}^n p(x_i) \log_2 p(x_i) \quad (1)$$

The equation (1) is the formula. After some results we conclude that is we have 50,000 of database size around the 20% continue after the filter, in my case i use 1.95 of filter. And we conclude that between more near the number of 2, the size will decreases.

```

2      CCATCAGGCAACGTGGT  1.971335
3      TCTGTGTGCAATCCGAT  1.954247
4      CACGGCATTTAAGA    1.959190
...
11178  CGGTGTACCAACTC    1.959190
11179  TGCTGGTCGATCGACATCA 1.993774
11180  CAACTACGTGGT     2.000000
11181  GACTCAGATC       1.970951
11182  TTAGTCAGACCGTGG   1.965596
[11183 rows x 2 columns]
```



And at the moment of get the motifs 6 and 8, i realize that if we choose a motif more big the combinations will decrease as you can see



- And after get the motifs we conclude choosing the size of the motif



```

13 s ✓ [176] print(get_motif(6, filtered_sequences))

      ('AGCTCA', 64)
```



```

2 m ✓ print(get_motif(8, filtered_sequences))

      ('GTATACGC', 11)
```



In conclusion, the Shannon Entropy help us to find patrons and in this case those patrons represent genetics diseases and the Shannon Entropy can make the difference depending the treshold you choose