

I529 HW1-Section1

Chia-Hsuan Chou

Jan. 29, 2016

I. Goal

The goal is to generate random sequences and calculate the mean and standard deviation of every nucleotide in generated sequences.

II. Procedure

For RANSEQ1.py and RANSEQ2.py procedure, they are basically the same. The only difference is the way to generate the random sequences. First, they both read the sequence from input FASTA file. However, RANSEQ1.py calculate the frequency of every nucleotide and based on that frequency to generate random sequences. For RANSEQ2.py, it shuffles the nucleotides appearing in the input sequence and generate new sequences based on the permutation. Next, they all calculate the standard deviation and mean for every nucleotide based on the generated sequences. Finally, they all print the information above on the screen and write the information in an output file.

III. Result

```
1 Random sequence generation:
2 TTCCTTTTTT
3 TTCACATTAA
4 GTCCATGATA
5 CTGTATTTTA
6 CTCCTATTTT
7 TTATTTATGT
8 AGGCAGTTTT
9 GCTTCITTTT
10 TCCTTATATT
11 TTAATATCTT
12 The mean of A is 0.1800
13 The mean of C is 0.1800
14 The mean of G is 0.0800
15 The mean of T is 0.5600
16 The standard deviation of A is 0.116619
17 The standard deviation of C is 0.097980
18 The standard deviation of G is 0.097980
19 The standard deviation of T is 0.149666
```

Figure 1. The output file from RANSEQ1.py

```
1 Random sequence permutation:
2 ATTTTCCATG
3 TTGCTACTAT
4 TCCTTGTTAA
5 ATCCTATTG
6 TTCTTAAGT
7 ATGCCTATTT
8 TGCTTCATAT
9 TATGTCTATC
10 CGTCTATTAT
11 TTTCAATGTAC
12 The mean of A is 0.2000
13 The mean of C is 0.2000
14 The mean of G is 0.1000
15 The mean of T is 0.5000
16 The standard deviation of A is 0.000000
17 The standard deviation of C is 0.000000
18 The standard deviation of G is 0.000000
19 The standard deviation of T is 0.000000
```

Figure 2. The output file from RANSEQ2.py

IV. Discussion

For the result of RANSEQ1.py, because the program generate the sequences based on the frequency of every input nucleotide, the mean and standard deviation will not be the same every time when we generate it. However, for the result of RANSEQ2.py, the program generate the sequences based on shuffling the input nucleotide, the mean will same as the mean of input sequence and the standard deviation will be 0 all the time.

I think that for random generate sequences, the way which RANSEQ1.py do will be better. Because for the random sequence permutation, there is no difference between the standard deviation of every nucleotide and the mean is exactly the same as input sequence and there is no any statistics analysis to do. If we think about evolution, it is not possible that every nucleotide numbers will be exactly the same via evolution. However, the frequency of every nucleotide might be the same and it can give us some clues of predicting new gene.