# REPORT

## 1-) Pseudo Code

PROBES( input.txt, probeLen)

1. subSeqDict = **dict**()

2. nuc = {'A', 'T', 'G', 'C'}

3. dnaSequences = **list**()

4. lengthOfValidSequence = 27000

5. **for** every line (**i**) **in** input.txt

6.     i = i[:lengthOfValidSequence]    #every i will have 27000 char

7.     **if** every char in i is equal to one of the nuc's chars

8.         add i to the dnaSequences

9. dnaSequences = **list**(**set**(dnaSequences))   #in order to eliminate same sequences, I use set()

10. **for** every dna sequence  **in** dnaSequences

11.     tempSet = set()

12.     **for** every proper subsequence of current dna sequence

13.         **if** current subsequence is already in subSeqDict and **not in** tempSet

14.             increase value of this subsequence in subSeqDict by 1

15.         **else**

16.             add current subsequence as key into subSeqDict and assign it's value as 1

17.         add current subsequence to the  tempSet

18. maxKeyValue = 1

19. **for** key, result **in** subSeqDict.items()

21.     find maxKeyValue

22. resultList = **list**()

23. **for** key, result **in** subSeqDict.items()

24.     add key to the resultList which has value that equal to maxKeyValue

25. write  number of valid sequences, maximum number sequences that the probes are found,

26. length of the returning list of the probes function and resultList to the output.txt file

**NOTE:** I use set() at line 9. Thus, if you run this algorithm for a limited number(10,500,1000, etc.) of valid DNA sequences instead of whole valid DNA sequences, this limited number of selected valid DNA sequences will be random. As a result, the result of the output.txt file can vary for each run. However, if you run for whole valid DNA sequences, the result will be always the same.

**2-)Results Table**

| The number of valid sequences | The number of sequences probe is found | Probe length | Time in seconds |
|---|---|---|---|
| 10 | 22823 | 90 | 10.49 |
| 50 | 12839 | 90 | 7.18 |
| 100 | 7334 | 90 | 7.77 |
| 250 | 1865 | 90 | 13.05 |
| 500 | 466 | 90 | 22.86 |
| 1000 | 53 | 90 | 42.21 |
| 10 | 22393 | 100 | 8.74 |
| 50 | 11835 | 100 | 6.75 |
| 100 | 6407 | 100 | 7.47 |
| 250 | 1467 | 100 | 13.51 |
| 500 | 323 | 100 | 23.41 |
| 1000 | 42 | 100 | 43.69 |