Algorithms

```python
 1 import random
 2 import math
 3 def create_sequence():
 4     nucleotid_bases=['A','C','G','T']
 5     size_sequence = random.randint(10,20)
 6     new_sequence = [nucleotid_bases[random.randint(0,3)] for i in range (size_sequence)]
 7     return "".join(new_sequence)
 8 print(create_sequence())
 9 def create_database():
10     db_size = 5000
11     data_base = [create_sequence() for i in range(db_size)]
12     return data_base
13 print(create_database())
14 from itertools import combinations
15 def get_combinations (n,sequence,bases):
16     if n == 1:
17         return [sequence+bases[i] for sequence in sequence for i in range(len(bases))]
18     else:
19         sequence_ = [s+bases[i] for s in sequence for i in range(len(bases))]
20         return get_combinations(n-1,sequence_,bases)
21 def count_motif(motif,sequence_db):
22     count = 0
23     for sequence in sequence_db:
24         count += sequence.count(motif)
25         return count
26 def get_motif(motif_size,sequences_db):
27     nucleotid_bases = ['A','C','G','T']
28     combinations = get_combinations(motif_size, [""], nucleotid_bases)
29     max_counter = 0
30     motif_winner = ""
31     for motif_candidate in combinations:
32         temp_counter = count_motif(motif_candidate,sequences_db)
33         if temp_counter > max_counter:
34             max_counter = temp_counter
35             motif_winner = motif_candidate
36     return motif_winner, max_counter
37 print(get_motif(6, create_database()))
38
```

First of all, the code that the professor made in class using recursion and control structures

```
37 motif_winner, _ = get_motif(6, create_database())
38 motif_winnerDos, _ = get_motif(8, create_database())
39 print("Initial motif winner of 6:", motif_winner)
40 print("Initial motif winner of 8:", motif_winnerDos)
41
42 def shannon_entropy(sequence):
43     counts = {base: sequence.count(base) for base in ['A','C','G','T']}
44     total_bases = sum(counts.values())
45     entropy = 0
46     for count in counts.values():
47         if count > 0:
48             probability = count / total_bases
49             entropy -= probability * math.log2(probability)
50     return entropy
51 print("Entropy of initial motif winner 6:", shannon_entropy(motif_winner))
52 print("Entropy of initial motif winner 8:", shannon_entropy(motif_winnerDos))
53
54 while shannon_entropy(motif_winner) <1.7 or shannon_entropy(motif_winner) >2:
55     motif_winner, _ = get_motif(6, create_database())
56 while shannon_entropy(motif_winnerDos) <2 or shannon_entropy(motif_winnerDos) >2.5:
57     motif_winnerDos, _ = get_motif(8, create_database())
58 print("Motif winner 6 with entropy :", motif_winner)
59 print("Final entropy:", shannon_entropy(motif_winner))
60 print("Motif winner 8 with entropy :", motif_winnerDos)
61 print("Final entropy:", shannon_entropy(motif_winnerDos))
62
```

Lately I save the motif winner to use the string in the calculate of entropy and last I create a cycle that depends of the value of the entropy this will be changed meanwhile the strings change and will stop when the algorithm found a string whit a entropy between 1.7 and 2.

Conclusions

- If the value of the entropy is lower than 1 the string will have so many repeated bases or letters
- I define the maximum entropy in 2 because higher is so difficult to found a string with that entropy
- The value of the entropy changes depends of the size of the string