
Biology Meets Programming

Bioinformatics for Begginers

Amin Kasrou Aouam

2020-10-10

Contents

Biology Meets Programming: Bioinformatics for Beginners **3**

Week 1 3

 DNA replication 3

Week 2 5

 DNA replication (II) 5

Week 3 10

 The circadian clock 10

Week 4 17

 The circadian clock (II) 17

Vocabulary 28

Biology Meets Programming: Bioinformatics for Beginners

Week 1

DNA replication

1. Origin of replication (ori)

Locating an ori is key for gene therapy (e.g. viral vectors), to introduce a therapeutic gene.

2. Computational approaches to find ori in *Vibrio Cholerae*

a) Exercise: find Pattern

We'll look for the **DnaA box** sequence, using a sliding window, in that case we will use this function to find out how many times a sequence appears in the genome:

```
1 def PatternCount(Text, Pattern):
2     count = 0
3     for i in range(len(Text)-len(Pattern)+1):
4         if Text[i:i+len(Pattern)] == Pattern:
5             count = count+1
6     return count
```

For the second part, we're going to calculate the frequency map of the sequences of length k :

```
1 def FrequentWords(Text, k):
2     words = []
3     freq = FrequencyMap(Text, k)
4     m = max(freq.values())
5     for key in freq:
6         if freq[key] == m:
7             words.append(key)
8     return words
9
10
11 def FrequencyMap(Text, k):
12     freq = {}
13     n = len(Text)
14     for i in range(n - k + 1):
15         Pattern = Text[i:i + k]
16         freq[Pattern] = 0
17     for i in range(n - k + 1):
18         Pattern = Text[i:i + k]
19         freq[Pattern] += 1
20     return freq
```

b) Exercise: Find the reverse complement of a sequence

We're going to generate the reverse complement of a sequence, which is the complement of a sequence, read in the same direction (5' -> 3').

```
1 def ReverseComplement(Pattern):
2     Pattern = Reverse(Pattern)
3     Pattern = Complement(Pattern)
4     return Pattern
5
6
7 def Reverse(Pattern):
8     reversed = Pattern[::-1]
9     return reversed
10
11
12 def Complement(Pattern):
13     compl = ""
14     complement_letters = {"A": "T", "T": "A", "C": "G", "G": "
15         C"}
16     for char in Pattern:
17         compl += complement_letters[char]
18     return compl
```

After using our function on the *Vibrio Cholerae*'s genome, we realize that some of the frequent *k-mers* are reverse complements of other frequent ones.

c) Exercise: Find a subsequence within a sequence

We're going to find the occurrences of a subsequence inside a sequence, and save the index of the first letter in the sequence.

```
1 def PatternMatching(Pattern, Genome):
2     positions = []
3     for i in range(len(Genome)-len(Pattern)+1):
4         if Genome[i:i+len(Pattern)] == Pattern:
5             positions.append(i)
6     return positions
```

We find out that the *9-mers* with the highest frequency appear in cluster. There is strong statistical evidence that our subsequences are *DnaA boxes*.

3. Computational approaches to find ori in any bacteria

Now that we're pretty confident about the *DnaA boxes* sequences that we found, we are going to check if they are a common pattern in the rest of bacterias. We're going to find the occurrences of the sequences in *Thermotoga petrophila*:

```
1 def PatternCount(Text, Pattern):
2     count = 0
3     for i in range(len(Text)-len(Pattern)+1):
4         if Text[i:i+len(Pattern)] == Pattern:
5             count = count+1
6     return count
```

We observe that there are **no** occurrences of the sequences found in *Vibrio Cholerae*. We can conclude that different bacterias have different *DnaA boxes*.

We have to try another computational approach, find clusters of *k-mers* repeated in a small interval.

Week 2

DNA replication (II)

1. Replication process

The *DNA polymerases* start replicating while the parent strands are unraveling. On the lagging strand, the DNA polymerase waits until the replication fork opens around 2000 nucleotides, and because of that it forms Okazaki fragments. We need 1 primer for the leading strand and 1 primer per Okazaki fragment for the lagging strand. While the Okazaki fragments are being synthesized, a *DNA ligase* starts joining the fragments together.

2. Computational approach to find ori using deamination

As the lagging strand is always waiting for the helicase to go forward, the lagging strand is mostly in single-stranded configuration, which is more prone to mutations. One frequent form of mutation is **deamination**, a process that causes cytosine to convert into thymine. This means that cytosine is more frequent in half of the genome.

a) Exercise: count the occurrences of cytosine

We're going to count the occurrences of the bases in a genome and include them in a symbol array.

```

1 def SymbolArray(Genome, symbol):
2     array = {}
3     n = len(Genome)
4     ExtendedGenome = Genome + Genome[0:n//2]
5     for i in range(n):
6         array[i] = PatternCount(ExtendedGenome[i:i+(n//2)],
7                                 symbol)
8     return array
9
10 def PatternCount(Text, Pattern):
11     count = 0
12     for i in range(len(Text)-len(Pattern)+1):
13         if Text[i:i+len(Pattern)] == Pattern:
14             count = count+1
15     return count

```

After executing the program, we realize that the algorithm is too inefficient.

- b) Exercise: find a better algorithm for the previous exercise

This time, we are going to evaluate an element $i+1$, using the element i .

```

1 def FasterSymbolArray(Genome, symbol):
2     array = {}
3     n = len(Genome)
4     ExtendedGenome = Genome + Genome[0:n//2]
5     array[0] = PatternCount(symbol, Genome[0:n//2])
6     for i in range(1, n):
7         array[i] = array[i-1]
8         if ExtendedGenome[i-1] == symbol:
9             array[i] = array[i]-1
10        if ExtendedGenome[i+(n//2)-1] == symbol:
11            array[i] = array[i]+1
12    return array
13
14
15 def PatternCount(Text, Pattern):
16     count = 0
17     for i in range(len(Text)-len(Pattern)+1):
18         if Text[i:i+len(Pattern)] == Pattern:
19             count = count+1
20    return count

```

It's a viable algorithm, with a complexity of $O(n)$ instead of the previous $O(n^2)$. In *Escherichia Coli* we plotted the result of our program:

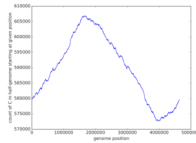


Figure 1: Symbol array for Cytosine in E. Coli Genome]

We can conclude that ori is located around position 4000000, because that's where the Cytosine concentration is the lowest, which indicates that the region stays single-stranded for the longest time.

3. The Skew Diagram

Usually scientists measure the difference between G - C, which is **higher on the lagging strand** and **lower on the leading strand**.

a) Exercise: Synthetize a Skew Array

We're going to make a Skew Diagram, for that we'll first need a Skew Array. For that purpose we wrote:

```

1 def SkewArray(Genome):
2     Skew = []
3     Skew.append(0)
4     for i in range(0, len(Genome)):
5         if Genome[i] == "G":
6             Skew.append(Skew[i] + 1)
7         elif Genome[i] == "C":
8             Skew.append(Skew[i] - 1)
9         else:
10            Skew.append(Skew[i])
11    return Skew

```

We can see the utility of a Skew Diagram looking at the one from *Escherichia Coli*:

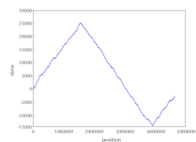


Figure 2: Symbol array for Cytosine in E. Coli Genome]

Ori should be located where the skew is at its minimum value.

b) Exercise: Efficient algorithm for locating ori

Now that we know more about ori's skew value, we're going to construct a better algorithm to find it:

```
1 def MinimumSkew(Genome):
2     positions = []
3     skew = SkewArray(Genome)
4     minimum = min(skew)
5     return [i for i in range(0, len(Genome)) if skew[i] ==
6             minimum]
7
8 def SkewArray(Genome):
9     Skew = []
10    Skew.append(0)
11    for i in range(0, len(Genome)):
12        if Genome[i] == "G":
13            Skew.append(Skew[i] + 1)
14        elif Genome[i] == "C":
15            Skew.append(Skew[i] - 1)
16        else:
17            Skew.append(Skew[i])
18    return Skew
```

4. Finding *DnaA* boxes

When we look for *DnaA* boxes in the minimal skew region, we can't find highly repeated 9-mers in *Escherichia Coli*. But we found approximate sequences that are similar to our 9-mers and only differ in 1 nucleotide.

a) Exercise: Calculate Hamming distance

The Hamming distance is the number of mismatches between 2 strings.

```
1 def HammingDistance(p, q):
2     count = 0
3     for i in range(0, len(p)):
4         if p[i] != q[i]:
5             count += 1
6     return count
```

b) Exercise: Find approximate patterns

Now that we have our Hamming distance, we use it to find the approximate sequences:


```
1 def ApproximatePatternMatching(Text, Pattern, d):
2     positions = []
3     for i in range(len(Text)-len(Pattern)+1):
4         if Text[i:i+len(Pattern)] == Pattern:
5             positions.append(i)
6         elif HammingDistance(Text[i:i+len(Pattern)], Pattern)
7             <= d:
8             positions.append(i)
9     return positions
10
11 def HammingDistance(p, q):
12     count = 0
13     for i in range(0, len(p)):
14         if p[i] != q[i]:
15             count += 1
16     return count
```

c) Exercise: Count the approximate patterns

The final part is counting the approximate sequences:

```
1 def ApproximatePatternCount(Pattern, Text, d):
2     count = 0
3     for i in range(len(Text) - len(Pattern) + 1):
4         if (
5             Text[i : i + len(Pattern)] == Pattern
6             or HammingDistance(Text[i : i + len(Pattern)],
7                 Pattern) <= d
8         ):
9             count += 1
10    return count
11
12 def HammingDistance(p, q):
13     count = 0
14     for i in range(0, len(p)):
15         if p[i] != q[i]:
16             count += 1
17     return count
```

After trying out our ApproximatePatternCount in the hypothesized ori region, we find a frequent *k-mer* with its reverse complement in *Escherichia Coli*. We've finally found a computational method to find ori that seems correct.

Week 3

The circadian clock

Variation in gene expression permits the cell to keep track of time.

1. Computational approaches to find regulatory motifs

a) Exercise: Find the most common nucleotides in each position

We are going to create a $t \times k$ Motif Matrix, where t is the k -mer string. In each position, we'll insert the most frequent nucleotide, in upper case, and the nucleotide in lower case (if there's no popular one). Our goal is to select the **most** conserved Matrix, i.e. the Matrix with the most upper case letters. We'll use a $4 \times k$ Count Matrix, one row for each base.

```
1 def Count(Motifs):
2     k = len(Motifs[0])
3     count = {'A': [0]*k, 'C': [0]*k, 'G': [0]*k, 'T': [0] * k}
4     t = len(Motifs)
5     for i in range(t):
6         for j in range(k):
7             symbol = Motifs[i][j]
8             count[symbol][j] += 1
9     return count
```

Now that we have a Count Matrix, we will generate a Profile Matrix, which has the frequency of the nucleotide instead of the count:

```
1 def Profile(Motifs):
2     t = len(Motifs)
3     profile = Count(Motifs)
4     for key, v in profile.items():
5         v[:] = [x / t for x in v]
6     return profile
7
8
9 def Count(Motifs):
10    k = len(Motifs[0])
11    count = {'A': [0]*k, 'C': [0]*k, 'G': [0]*k, 'T': [0] * k}
12    t = len(Motifs)
13    for i in range(t):
14        for j in range(k):
15            symbol = Motifs[i][j]
16            count[symbol][j] += 1
17    return count
```

b) Exercise: Form the most frequent sequence of nucleotides

Finally, we can form a Consensus string, to get a candidate regulatory motif:

```
1 def Consensus(Motifs):
2     consensus = ""
3     count = Count(Motifs)
4     k = len(Motifs[0])
5     for j in range(k):
6         m = 0
7         frequentSymbol = ""
8         for symbol in "ACGT":
9             if count[symbol][j] > m:
10                m = count[symbol][j]
11                frequentSymbol = symbol
12            consensus += frequentSymbol
13    return consensus
14
15
16 def Count(Motifs):
17     k = len(Motifs[0])
18     count = {'A': [0]*k, 'C': [0]*k, 'G': [0]*k, 'T': [0] * k}
19     t = len(Motifs)
20     for i in range(t):
21         for j in range(k):
22             symbol = Motifs[i][j]
23             count[symbol][j] += 1
24    return count
```

After obtaining the Consensus string, all we need to do is obtain the total score of our selected *k*-mers:

```
1 def Score(Motifs):
2     score = 0
3     for i in range(len(Motifs)):
4         score += HammingDistance(Motifs[i], Consensus(Motifs))
5     return score
6
7
8 def Consensus(Motifs):
9     consensus = ""
10    count = Count(Motifs)
11    k = len(Motifs[0])
12    for j in range(k):
13        m = 0
14        frequentSymbol = ""
15        for symbol in "ACGT":
16            if count[symbol][j] > m:
17                m = count[symbol][j]
18                frequentSymbol = symbol
19        consensus += frequentSymbol
20    return consensus
21
22
23 def Count(Motifs):
24    k = len(Motifs[0])
25    count = {'A': [0]*k, 'C': [0]*k, 'G': [0]*k, 'T': [0] * k}
26    t = len(Motifs)
27    for i in range(t):
28        for j in range(k):
29            symbol = Motifs[i][j]
30            count[symbol][j] += 1
31    return count
32
33
34 def HammingDistance(p, q):
35    count = 0
36    for i in range(0, len(p)):
37        if p[i] != q[i]:
38            count += 1
39    return count
```

c) Exercise: Find a set of k -mers that minimize the score

Applying a brute force approach for this task is not viable, we'll use a Greedy Algorithm. We first have to determine the probability of a sequence:

```
1 def Pr(Text, Profile):
2     probability = 1
3     k = len(Text)
4     for i in range(k):
5         probability *= Profile[Text[i]][i]
6     return probability
```

We'll use this function to find the most probable *k*-mer in a sequence:

```
1 def ProfileMostProbableKmer(text, k, profile):
2     kmer = ""
3     keys = ["A", "C", "G", "T"]
4     d = dict(zip(keys, profile))
5     prob = -1
6     for i in range(len(text)-k+1):
7         if (Pr(text[i:i+k], profile) > prob):
8             prob = Pr(text[i:i+k], profile)
9             kmer = text[i:i+k]
10    return kmer
11
12
13 def Pr(Text, Profile):
14     probability = 1
15     k = len(Text)
16     for i in range(k):
17         probability *= Profile[Text[i]][i]
18    return probability
```

Now we're finally ready to assemble all the pieces and implement a Greedy Motif Search Algorithm:

```

1  def GreedyMotifSearch(Dna, k, t):
2      BestMotifs = []
3      for i in range(0, t):
4          BestMotifs.append(Dna[i][0:k])
5      n = len(Dna[0])
6      for i in range(n-k+1):
7          Motifs = []
8          Motifs.append(Dna[0][i:i+k])
9          for j in range(1, t):
10             P = Profile(Motifs[0:j])
11             Motifs.append(ProfileMostProbableKmer(Dna[j], k, P
12             ))
13             if Score(Motifs) < Score(BestMotifs):
14                 BestMotifs = Motifs
15         return BestMotifs
16
17 def Score(Motifs):
18     score = 0
19     for i in range(len(Motifs)):
20         score += HammingDistance(Motifs[i], Consensus(Motifs))
21     return score
22
23 def Consensus(Motifs):
24     consensus = ""
25     count = Count(Motifs)
26     k = len(Motifs[0])
27     for j in range(k):
28         m = 0
29         frequentSymbol = ""
30         for symbol in "ACGT":
31             if count[symbol][j] > m:
32                 m = count[symbol][j]
33                 frequentSymbol = symbol
34         consensus += frequentSymbol
35     return consensus
36
37
38 def Count(Motifs):
39     k = len(Motifs[0])
40     count = {'A': [0]*k, 'C': [0]*k, 'G': [0]*k, 'T': [0] * k}
41     t = len(Motifs)
42     for i in range(t):
43         for j in range(k):
44             symbol = Motifs[i][j]
45             count[symbol][j] += 1
46     return count
47
48
49 def HammingDistance(p, q):
50     count = 0
51     for i in range(0, len(p)):
52         if p[i] != q[i]:
53             count += 1
54     return count
55
56
57 def Profile(Motifs):
58     t = len(Motifs)
59     profile = Count(Motifs)
60     for key, v in profile.items():

```

d) Motifs in tuberculosis

Tuberculosis is an infectious disease, caused by a bacteria called *Mycobacterium tuberculosis*. The bacteria can stay latent in the host for decades, in hypoxic environments. Our Greedy Algorithm can help us identify a motif that might be involved in the process.

The transcription factor behind this behaviour is **DosR**, we'll identify the binding sites:

```

1 def GreedyMotifSearch(Dna, k, t):
2     BestMotifs = []
3     for i in range(0, t):
4         BestMotifs.append(Dna[i][0:k])
5     n = len(Dna[0])
6     for i in range(n - k + 1):
7         Motifs = []
8         Motifs.append(Dna[0][i : i + k])
9         for j in range(1, t):
10            P = Profile(Motifs[0:j])
11            Motifs.append(ProfileMostProbableKmer(Dna[j], k, P
12                ))
13            if Score(Motifs) < Score(BestMotifs):
14                BestMotifs = Motifs
15        return BestMotifs
16
17 def Score(Motifs):
18     score = 0
19     for i in range(len(Motifs)):
20         score += HammingDistance(Motifs[i], Consensus(Motifs))
21     return score
22
23
24 def Consensus(Motifs):
25     consensus = ""
26     count = Count(Motifs)
27     k = len(Motifs[0])
28     for j in range(k):
29         m = 0
30         frequentSymbol = ""
31         for symbol in "ACGT":
32             if count[symbol][j] > m:
33                 m = count[symbol][j]
34                 frequentSymbol = symbol
35         consensus += frequentSymbol
36     return consensus
37
38
39 def Count(Motifs):
40     k = len(Motifs[0])
41     count = {"A": [0] * k, "C": [0] * k, "G": [0] * k, "T":
42         [0] * k}
43     t = len(Motifs)
44     for i in range(t):
45         for j in range(k):
46             symbol = Motifs[i][j]
47             count[symbol][j] += 1
48     return count
49
50 def HammingDistance(p, q):
51     count = 0
52     for i in range(0, len(p)):
53         if p[i] != q[i]:
54             count += 1
55     return count
56
57
58 def Profile(Motifs):
59     t = len(Motifs)

```


Our algorithm is pretty fast, but it's not optimal, and that's just a characteristic of Greedy Algorithms, they trade optimality for speed.

Week 4

The circadian clock (II)

The French mathematician Laplace estimated the probability that the sun will not raise tomorrow, this approach plays an important role in statistics, we cannot simplify the probability of a low-probability event to zero.

1. Motif finding with pseudocounts

In order to account for this problem, bioinformaticians often substitute zeroes with small numbers called pseudocounts.

a) Exercise: Create a count matrix with pseudocounts

We are going to generate the count matrix, while adding 1 to each value as a pseudocount. As a starting point we'll use our `Count(Motifs)` function, and we'll tweak it to achieve our objective.

```
1 def CountWithPseudocounts(Motifs):
2     k = len(Motifs[0])
3     count = {'A': [1]*k, 'C': [1]*k, 'G': [1]*k, 'T': [1] * k}
4     t = len(Motifs)
5     for i in range(t):
6         for j in range(k):
7             symbol = Motifs[i][j]
8             count[symbol][j] += 1
9     return count
```

We only modified the third line of the function, by starting the count at 1 instead of 0. With this simple adjustment, we have solved the problem.

b) Exercise: Create a profile matrix with pseudocounts

Now that we have our count matrix, we can generate a profile matrix. We'll adjust our `Profile(Motifs)` function for this purpose.

```
1 def ProfileWithPseudocounts(Motifs):
2     t = len(Motifs) + 4
3     profile = CountWithPseudocounts(Motifs)
4     for key, v in profile.items():
5         v[:] = [x / t for x in v]
6     return profile
7
8
9 def CountWithPseudocounts(Motifs):
10    k = len(Motifs[0])
11    count = {'A': [1]*k, 'C': [1]*k, 'G': [1]*k, 'T': [1] * k}
12    t = len(Motifs)
13    for i in range(t):
14        for j in range(k):
15            symbol = Motifs[i][j]
16            count[symbol][j] += 1
17    return count
```

We have only modified the value of t , adding 4 because the total sum of each column is different. Now that we have pseudocounts, we initialize each cell to 1, and because we have 4 rows, the total sum is now $t+4$;

c) Exercise: An improved Greedy Motif Search algorithm

We have all the required functions to construct a Greedy Motif Search algorithm with pseudocounts. As with all the previous exercises, we'll start with our `GreedyMotifSearch(Motifs)` function.

```

1  def GreedyMotifSearchWithPseudocounts(Dna, k, t):
2      BestMotifs = []
3      for i in range(0, t):
4          BestMotifs.append(Dna[i][0:k])
5      n = len(Dna[0])
6      for i in range(n - k + 1):
7          Motifs = []
8          Motifs.append(Dna[0][i : i + k])
9          for j in range(1, t):
10             P = ProfileWithPseudocounts(Motifs[0:j])
11             Motifs.append(ProfileMostProbableKmer(Dna[j], k, P
12             ))
13             if Score(Motifs) < Score(BestMotifs):
14                 BestMotifs = Motifs
15         return BestMotifs
16
17 def ProfileWithPseudocounts(Motifs):
18     t = len(Motifs) + 4
19     profile = CountWithPseudocounts(Motifs)
20     for key, v in profile.items():
21         v[:] = [x / t for x in v]
22     return profile
23
24
25 def CountWithPseudocounts(Motifs):
26     k = len(Motifs[0])
27     count = {'A': [1]*k, 'C': [1]*k, 'G': [1]*k, 'T': [1] * k}
28     t = len(Motifs)
29     for i in range(t):
30         for j in range(k):
31             symbol = Motifs[i][j]
32             count[symbol][j] += 1
33     return count
34
35
36 def Score(Motifs):
37     score = 0
38     for i in range(len(Motifs)):
39         score += HammingDistance(Motifs[i], Consensus(Motifs))
40     return score
41
42
43 def Consensus(Motifs):
44     consensus = ""
45     count = Count(Motifs)
46     k = len(Motifs[0])
47     for j in range(k):
48         m = 0
49         frequentSymbol = ""
50         for symbol in "ACGT":
51             if count[symbol][j] > m:
52                 m = count[symbol][j]
53                 frequentSymbol = symbol
54         consensus += frequentSymbol
55     return consensus
56
57
58 def Count(Motifs):
59     k = len(Motifs[0])
60     count = {"A": [0] * k, "C": [0] * k, "G": [0] * k, "T":

```

All we had to do was replace the function call `Profile(Motifs)` with `ProfileWithPseudo-counts(Motif)`.

The performance of the new algorithm is better, but we are still not satisfied. We are going to try a different motif finding algorithm.

2. Randomized Motif Search

A randomized algorithm sounds like a very bad idea, luck doesn't seem to be a good scientific asset. Well, as nonintuitive as it might sound they are used in many cases.

We will consider **Monte Carlo algorithms** for our use case.

a) Exercise: Generate the most probable *k*-mers

We'll start off our pipeline by crafting a function that outputs the most probable *k*-mers using a profile matrix, and a DNA sequence.

We'll reuse our `ProfileMostProbableKmer` and `Pr` functions:

```
1 def Motifs(Profile, Dna):
2     kmer_list = []
3     for sequence in Dna:
4         kmer = ProfileMostProbableKmer(text=sequence, k=len(
5             Profile), profile=Profile)
6         kmer_list.append(kmer)
7     return kmer_list
8
9
10 def ProfileMostProbableKmer(text, k, profile):
11     kmer = ""
12     keys = ["A", "C", "G", "T"]
13     d = dict(zip(keys, profile))
14     prob = -1
15     for i in range(len(text)-k+1):
16         if (Pr((text[i:i+k]), profile) > prob):
17             prob = Pr(text[i:i+k], profile)
18             kmer = text[i:i+k]
19     return kmer
20
21
22 def Pr(Text, Profile):
23     probability = 1
24     k = len(Text)
25     for i in range(k):
26         probability *= Profile[Text[i]][i]
27     return probability
```

We obtain a list of the most probable *k*-mer for each DNA sequence.

A common approach to is starting from a collection of randomly chosen *k-mers* Motifs, construct a profile matrix and use this matrix to generate a new collection of *k-mers*.

We'll iterate many times, hoping to get a better result and seeing if the score keeps improving. This is basically what a Randomized Motif Search does, now that we understand the underlying mechanism, we'll implement a random number generator.

b) Exercise: Implement a random *k-mer* generator

In order to implement our random *k-mer* generator, we need a way to generate random integers. Each integer will be used as the index for a character in our sequences.

Python's **random** module is our choice of tool for this exercise, specifically the **randint** function.

```
1 import random
2
3
4 def RandomMotifs(Dna, k, t):
5     random_motifs = []
6     t = len(Dna)
7     l = len(Dna[0])
8     for i in range(t):
9         random_index = random.randint(1, l-k)
10        random_motifs.append(Dna[i][random_index:random_index+
11                               k])
11    return random_motifs
```

We are ready to develop our Randomized Search algorithm, we just need to iterate over the generation of random motifs until we stop getting good results.

c) Exercise: Random Motif Search algorithm

Finally, we only have to assemble our functions to iterate through the random motifs while the score keeps improving.

```

1  import random
2
3
4  def RandomizedMotifSearch(Dna, k, t):
5      M = RandomMotifs(Dna, k, t)
6      BestMotifs = M
7      while True:
8          Profile = ProfileWithPseudocounts(M)
9          M = Motifs(Profile, Dna)
10         if Score(M) < Score(BestMotifs):
11             BestMotifs = M
12         else:
13             return BestMotifs
14
15
16  def RandomMotifs(Dna, k, t):
17      random_motifs = []
18      t = len(Dna)
19      l = len(Dna[0])
20      for i in range(t):
21          random_index = random.randint(1, l-k)
22          random_motifs.append(Dna[i][random_index:random_index+
23                                 k])
24      return random_motifs
25
26  def Motifs(Profile, Dna):
27      kmer_list = []
28      for sequence in Dna:
29          kmer = ProfileMostProbableKmer(text=sequence, k=len(
30              Profile["A"]), profile=Profile)
31          kmer_list.append(kmer)
32      return kmer_list
33
34
35  def ProfileMostProbableKmer(text, k, profile):
36      kmer = ""
37      keys = ["A", "C", "G", "T"]
38      d = dict(zip(keys, profile))
39      prob = -1
40      for i in range(len(text)-k+1):
41          if (Pr((text[i:i+k]), profile) > prob):
42              prob = Pr(text[i:i+k], profile)
43              kmer = text[i:i+k]
44      return kmer
45
46
47  def Pr(Text, Profile):
48      probability = 1
49      k = len(Text)
50      for i in range(k):
51          probability *= Profile[Text[i]][i]
52      return probability

```

```

53
54
55  def ProfileWithPseudocounts(Motifs):
56      t = len(Motifs) + 4
57      profile = CountWithPseudocounts(Motifs)
58      for key, v in profile.items():
59          v[:] = [x / t for x in v]

```

Each time that we execute our algorithm, we'll get a different solution. We're going to do an experiment by running our algorithm multiple times and checking the final score.

d) Exercise: Evaluate performance of the algorithm

In order to evaluate the performance of this algorithm, we are going to run 100 iterations while keeping the motif with the greatest motif.

```

1  import random
2  Dna = ["
      GCGCCCCGCCCCGACAGCCATGCGCTAACCCCTGGCTTCGATGGCGCCGGCTCAGTTAGGGCCGGAAGTCCC
      ",
3  "
      CCGATCGGCATCACTATCGGTCCTGCGGCCGCCATAGCGCTATATCCGGCTGGTAAAATCAATTGACAAC
      ",
4  "
      ACCGTCGATGTGCCCGGTGCGGCCGCGTCCACCTCGGTCATCGACCCACGATGAGGACGCCATCGGCCGC
      ",
5  "
      GGGTCAGGTATATTTATCGCACACTTGGGCACATGACACACAAGCGCCAGAATCCCGACCGAACCGAGCA
      ",
6  "
      GTAGGTCAAACCGGGTGTACATACCCGCTCAATCGCCCAGCACTTCGGGCAGATCACCGGTTTCCCGGT
      ",
7  "
      CCGTGGCGACGCTGTTCCGCCGAGCGTGCGTGACGACTTCGAGCTGCCCGACTACACCTGGTGACCACC
      ",
8  "
      GGCCAACTGCACCGCGCTCTTGATGACATCGGTGGTCACCATGGTGTCCGGCATGATCAACCTCCGCTGT
      ",
9  "
      GTACATGTCCAGAGCGAGCCTCAGCTTCTGCGCAGCGACGGAACTGCCCACTCAAAGCCTACTGGGCGC
      ",
10 "
      GGCAGCTGTCGGCAACTGTAAGCCATTTCTGGGACTTTGCTGTGAAAAGCTGGGCGATGTTGTGGACCTG
      ",
11 "
      TCAGCACCATGACCGCCTGGCCACCAATCGCCCGTAACAAGCGGGACGTCCGCGACGACGCGTGCGCTAGC
      "]
12 t = 10
13 k = 15
14 N = 100
15
16
17 def RandomizedMotifSearch(Dna, k, t):
18     M = RandomMotifs(Dna, k, t)
19     BestMotifs = M
20     while True:
21         Profile = ProfileWithPseudocounts(M)
22         M = Motifs(Profile, Dna)
23         if Score(M) < Score(BestMotifs):
24             BestMotifs = M
25         else:
26             return BestMotifs
27
28
29 def RandomMotifs(Dna, k, t):
30     random_motifs = []
31     t = len(Dna)
32     l = len(Dna[0])
33     for i in range(t):
34         random_index = random.randint(1, l-k)
35         random_motifs.append(Dna[i][random_index:random_index+k])
36     return random_motifs
37
38
39 def Motifs(Profile, Dna):
40     kmer_list = []

```


As we can see, this algorithm finds a pretty good solution.

3. Gibbs Sampling

Thanks to our previous experiments, we now consider random search algorithms as adequate tools for our problem. The caveat of our algorithm is that it discards all previous motifs in each iteration.

A more cautious alternative is Gibbs Sampling, which discards a single *k-mer* from the current set of motifs at each iteration.

a) Exercise: Normalize the probabilities

The algorithm chooses randomly at each iteration which *k-mer* will be dropped, and it chooses the replacement of that *k-mer*.

We use pseudocounts to generate the next profile matrix, which gives us some extravagant probabilities. We need to normalize them, so their sum equals 1.

```
1 def Normalize(Probabilities):
2     d = {}
3     for k,v in Probabilities.items():
4         d[k] = Probabilities[k]/sum(Probabilities.values())
5     return d
```

b) Exercise: Simulate rolling a die

Now that our probabilities are normalized, we can focus on simulating a weighted die.

```
1 import random
2
3 def WeightedDie(Probabilities):
4     count = 0
5     p = random.uniform(0,1)
6     for k,v in Probabilities.items():
7         count = count+v
8         if p < count:
9             return k
```

Our function returns a single element, in the next step we'll make it a subroutine of a larger function.

c) Exercise: Generate a *k-mer*

We have all the necessary tools to generate a *k-mer* based on the profile matrix.

```
1 import random
2
3
4 def Pr(Text, Profile):
5     probability = 1
6     k = len(Text)
7     for i in range(k):
8         probability *= Profile[Text[i]][i]
9     return probability
10
11
12 def Normalize(Probabilities):
13     d = {}
14     for k,v in Probabilities.items():
15         d[k] = Probabilities[k]/sum(Probabilities.values())
16     return d
17
18
19 def WeightedDie(Probabilities):
20     count = 0
21     p = random.uniform(0,1)
22     for k,v in Probabilities.items():
23         count = count+v
24         if p < count:
25             return k
26
27
28 def ProfileGeneratedString(Text, profile, k):
29     n = len(Text)
30     probabilities = {}
31     for i in range(0,n-k+1):
32         probabilities[Text[i:i+k]] = Pr(Text[i:i+k], profile)
33     probabilities = Normalize(probabilities)
34     return WeightedDie(probabilities)
```

d) Exercise: Implement the Gibbs Sampling algorithm

Finally, we can implement the Gibbs Sampling algorithm:

```

1  import random
2
3  def GibbsSampler(Dna, k, t, N):
4      Motifs = RandomMotifs(Dna, k ,t)
5      BestMotifs = Motifs[:]
6      for i in range(N):
7          p = random.randint(0,t-1)
8          del Motifs[p]
9          profile = ProfileWithPseudocounts(Motifs)
10         Motifs.insert(p, ProfileGeneratedString(Dna[p],
11             profile, k))
12         if Score(Motifs) < Score(BestMotifs):
13             BestMotifs = Motifs
14     return BestMotifs
15
16 def RandomMotifs(Dna, k, t):
17     random_motifs = []
18     t = len(Dna)
19     l = len(Dna[0])
20     for i in range(t):
21         random_index = random.randint(1, l-k)
22         random_motifs.append(Dna[i][random_index:random_index+
23             k])
24     return random_motifs
25
26 def ProfileWithPseudocounts(Motifs):
27     t = len(Motifs) + 4
28     profile = CountWithPseudocounts(Motifs)
29     for key, v in profile.items():
30         v[:] = [x / t for x in v]
31     return profile
32
33
34 def CountWithPseudocounts(Motifs):
35     k = len(Motifs[0])
36     count = {'A': [1]*k, 'C': [1]*k, 'G': [1]*k, 'T': [1] * k}
37     t = len(Motifs)
38     for i in range(t):
39         for j in range(k):
40             symbol = Motifs[i][j]
41             count[symbol][j] += 1
42     return count
43
44
45 def Pr(Text, Profile):
46     probability = 1
47     k = len(Text)
48     for i in range(k):
49         probability *= Profile[Text[i]][i]
50     return probability
51
52
53 def Normalize(Probabilities):
54     d = {}
55     for k,v in Probabilities.items():
56         d[k] = Probabilities[k]/sum(Probabilities.values())
57     return d
58
59

```

Vocabulary

- k-mer: subsequences of length k in a biological sequence
- Frequency map: sequence \rightarrow frequency of the sequence