



T.C.
MARMARA UNIVERSITY
FACULTY OF ENGINEERING
COMPUTER ENGINEERING DEPARTMENT

CSE4065 – Introduction to Computational Genomics

Project 1

GROUP MEMBERS

150117068 - Murathan Akman

150120048 – Yasin Tuğra Taş

```
def randomized_motif_search(dna, k):
    best_motifs = create_random_motifs(dna, k)
    best_score = score_motifs(best_motifs, k)
    score_history = [best_score]

    while True:
        i = random.randint(0, len(dna[0]) - k)
        motifs = [d[i:i+k] for d in dna]
        current_score = score_motifs(motifs, k)

        if current_score < best_score:
            best_motifs, best_score = motifs, current_score
            score_history = [best_score] # Reset since improvement occurred
        else:
            score_history.append(current_score)

        if len(score_history) > 50 and all(s >= best_score for s in score_history[-50:]): # Stuck
            break

    return best_motifs, best_score
```

RUN:

#1

```
k= 9
-----
Best Score: {43}
Best Motifs: ['CGTCTTAGT', 'TCTGCACCT', 'AACACAAGT', 'AGTATGAAT', 'AATACCGCG', 'GAACGAAGT', 'CATACAAGC', 'ACCAATTGC', 'GTGACATTA', 'ATTCAGATT']
Consensus: AATACAAGT
Execution Time: 2.6941299438476562e-05

k= 10
-----
Best Score: {49}
Best Motifs: ['GCACGTCTAGT', 'TCTGCACCTG', 'AACACAAGTA', 'AGTATGAATA', 'AATACCGCGT', 'GAACGAAGCT', 'CATACAAGCT', 'ACCAATTGCA', 'GTGACATTAT', 'ATTCAGATTA']
Consensus: AATACAAGTA
Execution Time: 2.9087066650390625e-05

k= 11
-----
Best Score: {56}
Best Motifs: ['GCACGTCTTAG', 'TAATCTGCACC', 'TCGAACACAAG', 'TCGAGTATGAA', 'TGAATACCGC', 'GACGAACGAAC', 'CGCCATACAAG', 'GATACCAATTG', 'TCTGTGACATT', 'TAAATTCAGAT']
Consensus: TAAATACAAAG
Execution Time: 3.0994415283203125e-05
```

#2

```
k= 9
-----
Best Score: {46}
Best Motifs: ['GCACGTCTT', 'TAATCTGCA', 'TCGAACACA', 'TCGAGTATG', 'TGGAATACC', 'GACGAACGA', 'CGCCATACA', 'GATACCAAT', 'TCTGTGACA', 'TAAATTTCAG']
Consensus: TAAATACACA
Execution Time: 2.503395080564062e-05

k= 10
-----
Best Score: {52}
Best Motifs: ['GTCTTAGTGTG', 'AGGTAGCAAC', 'GGCACCAACC', 'CAATAAACGC', 'GTAGCTGCGT', 'GGCCGACAGC', 'GTGGAGTGGC', 'GTCAAGGAGC', 'TGGGAAAACG', 'CCAAGCTGAGT']
Consensus: GGGGAAGAGC
Execution Time: 2.5987625122070312e-05

k= 11
-----
Best Score: {55}
Best Motifs: ['GTCTTAGTGTG', 'CTGCACCTGAG', 'ACACAAGTATG', 'GTATGAATATC', 'ATACCGCGTCG', 'AACGAAGTGA', 'ATACAAGCTGC', 'CCAATTGCAGG', 'TGACATTATTA', 'TTCAGATTACC']
Consensus: ATACAAGTATG
Execution Time: 2.7179718017578125e-05
```

#3

```
k= 9
-----
Best Score: {46}
Best Motifs: ['CAATAAGC', 'TGAGGAATA', 'ATGGTGTGG', 'GGCTAAGCA', 'CTATCGACA', 'ATGTATTTT', 'GAACAGAAC', 'AAGTAAATC', 'ATAGAAGCC', 'AAGCAAGTC']
Consensus: AAATAAATC
Execution Time: 3.2901763916015625e-05

k= 10
-----
Best Score: {50}
Best Motifs: ['GTCTTAGTGT', 'CTGCACCTGA', 'ACACAAGTAT', 'GTATGAATAT', 'ATACCGCGTC', 'AACGAAGTGA', 'ATACAAGCTG', 'CCAATTGCAG', 'TGACATTATT', 'TTCAGATTAC']
Consensus: ATACAAGTAT
Execution Time: 3.0994415283203125e-05

k= 11
-----
Best Score: {50}
Best Motifs: ['CCGCACGTCTT', 'ATTAATCTGCA', 'TGTCGAACACA', 'TATCGAGTATG', 'AATGGAATACC', 'TAGACGAACGA', 'CACGCCATACA', 'CCGATACCAAT', 'GTTCTGTGACA', 'GGTAAATTTCAG']
Consensus: CATAAAATACA
Execution Time: 2.9802322387695312e-05
```

#4

```
k= 9
-----
Best Score: {44}
Best Motifs: ['GTCTTAGTGT', 'CTGCACCTG', 'ACACAAGTA', 'GTATGAATA', 'ATACCGCGT', 'AACGAAGT', 'ATACAAGCT', 'CCAATTGCA', 'TGACATTAT', 'TTCAGATTA']
Consensus: ATACAAGTA
Execution Time: 2.384185791015625e-05

k= 10
-----
Best Score: {53}
Best Motifs: ['AGCCGGGAGA', 'CAACACCAAGT', 'GCTCCGAGGC', 'CAAAGCAGGA', 'TGAGACGTAA', 'TCATCAGCT', 'CGTCTCATCA', 'TAGCGATGGA', 'GCATGGTGGC', 'TCAGAAGGTA']
Consensus: TCACACGGGA
Execution Time: 2.574920654296875e-05

k= 11
-----
Best Score: {56}
Best Motifs: ['CACGTCTTAGT', 'AATCTGCACCT', 'CGAACACAAGT', 'CGAGTATGAAT', 'GGAATACCGG', 'ACGAACGAAGT', 'GCCATACAAGC', 'ATACCAATTGC', 'CTGTGACATTA', 'AAATTTCAGAT']
Consensus: AAAATACAAGT
Execution Time: 2.6702880859375e-05
```

#5

```
k= 9
-----
Best Score: {46}
Best Motifs: ['GATAGCTGA', 'GTAGCAACA', 'CACCAACTT', 'ATAAACGCC', 'AGCTGCGTC', 'CCGACAGCA', 'GGAGTGGCA', 'CAAGGAGCA', 'GGAAAAGT', 'AGCTGAGTT']
Consensus: GGAAGAGCA
Execution Time: 2.9087066650390625e-05

k= 10
-----
Best Score: {52}
Best Motifs: ['CACGTCTTAG', 'AATCTGCACC', 'CGAACACAAG', 'CGAGTATGAA', 'GGAATACCGC', 'ACGAACGAAC', 'GCCATACAAG', 'ATACCAATTG', 'CTGTGACATT', 'AAATTTCAGAT']
Consensus: AAAATACAAG
Execution Time: 3.0994415283203125e-05

k= 11
-----
Best Score: {56}
Best Motifs: ['GCACGTCTTAG', 'TAATCTGCACC', 'TCGAACACAAG', 'TCGAGTATGAA', 'TGAATACCGC', 'GACGAACGAAC', 'CGCCATACAAG', 'GATACCAATTG', 'TCTGTGACATT', 'TAAATTTCAGAT']
Consensus: TAAATACAAAG
Execution Time: 3.528594970703125e-05
```

2. Gibbs Sampler

In this algorithm, we again randomly selected motifs from each row. These motifs are the best motif and the best score at the beginning. In this algorithm we have a counter to count the number of iterations. In a while loop, we choose a random motif at each iteration. We then calculated the probability of each motif based on the profile matrix. Again, we calculate the score of new motifs and perform the same operations as in the random motif search. The end condition of this algorithm is that the algorithm checks the best score every 50 iterations. If the best score is not increased, the algorithm terminates.

```
# Gibbs Sampler
def gibbs_sampler(dna, k):
    motifs = random_motifs(dna, k)
    best_motifs = motifs[:]
    best_score = score(motifs)
    scores = [best_score] # Keep track of the scores

    while True:
        start_time = time.time()

        i = random.randint(0, len(dna) - 1)
        motifs_except_i = motifs[:i] + motifs[i+1:]
        profile = profile_with_pseudocounts(motifs_except_i)
        motifs[i] = weighted_random_kmer(dna[i], k, profile)
        current_score = score(motifs)

        if current_score < best_score:
            best_motifs = motifs[:]
            best_score = current_score
            scores.append(best_score) # Update
        else:
            scores.append(current_score) # Keep

        if len(scores) > 50 and scores[-1] == scores[-51]:
            break # Stuck

    end_time = time.time()
    execution_time = end_time - start_time

    return best_motifs, best_score, execution_time
```

RUN:

#1

```
k= 9
-----
Best Score: {32}
Best Motifs: ['ACGTCTTAG', 'AGATCTAGG', 'CTGGCTTAT', 'AGATTTTT', 'CGCTTTAT', 'TGGGCTTGT', 'AGAGCTTGT', 'TCACCTGG', 'AAAGCATGG', 'CAATCTGGT']
Consensus: AGATCTTGT
Execution Time: 0.0004248619 seconds

k= 10
-----
Best Score: {41}
Best Motifs: ['TCCCAAGCCC', 'TTGCAGTCTT', 'AACCTAAGCT', 'ATCCAAAGGC', 'ACCTAGAACC', 'TCAGGAAGTG', 'CACCAAAGTA', 'CACCTAGGTC', 'AATCGAAGAT', 'CCACAAAGCA']
Consensus: AACCAAGGCC
Execution Time: 0.0003848076 seconds

k= 11
-----
Best Score: {5}
Best Motifs: ['GCCACAAAGCA', 'CCCACAAAGCA', 'CCCACAAAGCA', 'TCCACAAAGCA', 'TCCACAAAGCA', 'ACCACAAAGCA', 'CCCACAAAGCA', 'CCCACAAAGCA', 'TCCACAAAGCA', 'CCCACAAAGCA']
Consensus: CCCACAAAGCA
Execution Time: 0.0003869534 seconds
```

#2

```
k= 9
-----
Best Score: {0}
Best Motifs: ['CACAAAGCA', 'CACAAAGCA', 'CACAAAGCA', 'CACAAAGCA', 'CACAAAGCA', 'CACAAAGCA', 'CACAAAGCA', 'CACAAAGCA', 'CACAAAGCA', 'CACAAAGCA']
Consensus: CACAAAGCA
Execution Time: 0.0003590584 seconds

k= 10
-----
Best Score: {39}
Best Motifs: ['TGTTTCTCT', 'AGACTCCTTA', 'ACTAACCTAA', 'CCATGCCACA', 'CGTTTGCTCA', 'GGTGCCGGA', 'AGTGCCGAA', 'AGCAGTCACT', 'GGATACGTAA', 'ACGTGCCTAA']
Consensus: AGATGCCTAA
Execution Time: 0.0004088879 seconds

k= 11
-----
Best Score: {41}
Best Motifs: ['GTTTCATCAAGT', 'AGCCATATTTG', 'CGGCTCTGTC', 'ATTCACTTAC', 'CTCCAATTGTA', 'CTTAAGTTTTA', 'ATTCACTGATA', 'CTCCATCCTTG', 'GTTAATCGAAG', 'CTCCATGTTGG']
Consensus: CTTCATGTTG
Execution Time: 0.0003840923 seconds
```

#3

```
k= 9
-----
Best Score: {27}
Best Motifs: ['ATAAGCCTA', 'CGGCAGAAA', 'CCATACTAA', 'ATAAACTAA', 'ACACTGCAA', 'ATAATGCTA', 'ACAAAGCAA', 'CTATTGCTA', 'ACAAACCGA', 'ACAAAGCAA']
Consensus: ACAAGCAA
Execution Time: 0.0003631115 seconds

k= 10
-----
Best Score: {32}
Best Motifs: ['CGATGATCGG', 'TAACGATAGC', 'TTCCGATCGC', 'TTACGACAGA', 'TAACATGCA', 'TTACGATAAC', 'CAACAACGGT', 'TAGCGATGGA', 'GAACATTCGC', 'TAATCATGGT']
Consensus: TAACGATGCG
Execution Time: 0.0003862381 seconds

k= 11
-----
Best Score: {40}
Best Motifs: ['CTGATAGCTGA', 'ATATCCGAGA', 'CTTATCGTTGA', 'CTGGCGGTTC', 'ACTTTCGTACA', 'GTGTTGATGA', 'CAGAACGTTGA', 'GTTATTGTTGA', 'GGGTATTAGGA', 'GCGGTCTATGA']
Consensus: CTGATCGTTGA
Execution Time: 0.0004138947 seconds
```

#4

```
k= 9
-----
Best Score: {35}
Best Motifs: ['TGACACATC', 'CTACACATC', 'CGTCAATGC', 'CTAAGCATC', 'CGTCACGCT', 'CTTCTTGCA', 'CGACTCCGC', 'TTGCTAGTG', 'CCACATGTG', 'GGACATGTC']
Consensus: CGACACGTC
Execution Time: 0.0003700256 seconds

k= 10
-----
Best Score: {39}
Best Motifs: ['AGGCCGGGCT', 'CTTAGTGGG', 'TTGCTCGGT', 'GTGCGGACAC', 'CTGCGTCTT', 'CGGCTGCCGG', 'CTGCTGTGCT', 'GGTTGCGCT', 'CCTCTGGCTT', 'CTGAGTTCGT']
Consensus: CTGCGGGCGT
Execution Time: 0.0003497601 seconds

k= 11
-----
Best Score: {41}
Best Motifs: ['CTGATCCGCGA', 'ATGATGAGTTT', 'AAGAACAGCCA', 'ACCATGCAATA', 'AACAAGCGACT', 'ATGATGCAATA', 'CTCATCAGAGA', 'ATGTAGACACA', 'ATGGTGCGCTA', 'CACATCCAACT']
Consensus: ATGATGCGACA
Execution Time: 0.0003719330 seconds
```

#5

```
k= 9
-----
Best Score: {33}
Best Motifs: ['GATCCGGCA', 'ACTCCGAAT', 'ACCCCGCTG', 'AAAGGCCAT', 'AACTCCAAT', 'AATCTGACA', 'AATTTGCCA', 'AATATGCCT', 'ACTCTCCAC', 'AATCCCCAC']
Consensus: AATCCGCAT
Execution Time: 0.0003957748 seconds

k= 10
-----
Best Score: {11}
Best Motifs: ['TGCCACAAAG', 'CCCCACAAAG', 'CCCCACAAAG', 'GTCCACAAAG', 'ATCCACAAAG', 'AACCACAAAG', 'GCCACAAAG', 'CCCCACAAAG', 'TTCCACAAAG', 'CCCCACAAAG']
Consensus: CCCCCAAAG
Execution Time: 0.0003471375 seconds

k= 11
-----
Best Score: {40}
Best Motifs: ['TAGCTGACACA', 'CACCACTACA', 'TTCCAGACTCA', 'CAGGAATACT', 'TAGTCGCTCC', 'ATGCAGAGATG', 'CACCGAACACA', 'TACGTGCCAAG', 'TACCATACTTA', 'TACCGAATTCG']
Consensus: TACCAGACACA
Execution Time: 0.0003669262 seconds
```

3. Median String

With the Median String problem, we aimed to find the k-mer motif that is closest to all sequences within the DNA sequence. In this problem, we measured the “nearest” term by the total Hamming distance between the selected k-mer and the closest k-mer in each DNA sequence. Median String found the k-mer that minimizes this total distance.

```
# Median String
def median_string(dna, k):
    min_distance = float('inf')
    execution_time = time.time()
    for kmer in all_kmers(k):
        distance = sum(min(hamming_distance(kmer, seq[i:i+k]) for i in range(len(seq) - k + 1)) for seq in dna)
        if distance < min_distance:
            min_distance = distance
            median_kmer = kmer

    execution_time = time.time() - execution_time
    return median_kmer, execution_time
```

RUN:

```
Median String for k=9: CACAAAGCA
Median String for k=10: CCACAAAGCA
Median String for k=11: CCCACAAAGCA
```

3. Outputs

# of RUN	Randomized Motif Search Score	k value		
1	43	9		
	49	10		
	56	11		
2	46	9		
	52	10		
	55	11		
3	46	9		
	50	10		
	59	11		
4	44	9		
	53	10		
	56	11		
5	46	9		
	52	10		
	56	11		
			Average Score	k value
			45	9
			51	10
			56	11

# of RUN	Gibbs Sampler Score	k value		
1	32	9		
	41	10		
	5	11		
2	0	9		
	38	10		
	41	11		
3	27	9		
	32	10		
	40	11		
4	35	9		
	39	10		
	41	11	Average Score	k value
5	33	9	25	9
	11	10	32	10
	40	11	33	11

When we compare these 3 algorithms, we can make the following conclusions:

- No significant difference was observed between the Gibbs Sampler and Randomized Motif Search algorithms in terms of execution time.
- It is understood from the scores in the tables above that Gibbs Sampler has better flexibility in avoiding local optima than the Randomized Motif Search algorithm for this project.
- It has been observed that Gibbs Sampler and Randomized Motif Search algorithms follow a parallel performance as the k value for k-mer increases.
- It is necessary to open separate parentheses for the Median String algorithm. Since this algorithm works with brute force logic, it can directly reach the perfect result and global optimum. However, as k value increases for k-mer, it reaches an unacceptable execution time. The execution time issue makes this algorithm useless even though it gives perfect results.