

11300A Bioinformatics

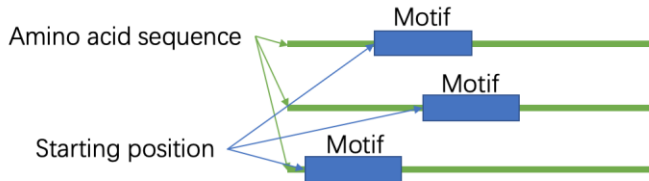
Homework 4 Dec. 19, 2023

1. Identifying protein motifs using Expectation Maximization (EM) algorithm.

Suppose we have 200 protein sequences of length 100, in the dataset. Each character represents an amino acid. The length of the motif is 10. The figure below shows the first 10 samples in the dataset.

1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25	26	27	28	29	30	31	32	33	34	35	36	37	38	39	40	41	42	43	44	45	46	47	48	49	50	51	52	53	54	55	56	57	58	59	60	61	62	63	64	65	66	67	68	69	70	71	72	73	74	75	76	77	78	79	80	81	82	83	84	85	86	87	88	89	90	91	92	93	94	95	96	97	98	99	100
M	C	K	R	E	C	F	N	G	N	C	W	H	M	S	K	H	N	N	N	K	G	L	N	T	F	M	V	T	L	C	W	S	C	I	T	V	S	C	S	V	L	T	C	G	P	D	V	I	A	M	G	M	V	N	H	G	S	N	C	I	D	P	V	G	E	G	T	V	W	W	H	A	L	N	H	H	L	E	I	V	F	H	S	E	F	R	D	R	A	S	D	S	E	F	N				
L	C	G	M	F	T	T	C	P	W	S	F	D	Y	T	A	R	N	E	C	S	N	N	K	R	A	N	D	T	K	E	G	D	I	F	T	C	D	F	A	K	A	C	K	N	A	V	G	L	F	E	L	C	R	G	E	L	R	A	I	A	L	G	W	T	Y	H	G	D	Y	W	E	M	M	E	P	H	Q	F	N	G	S	R	S	S	W	G	K	K	H	D	K								
T	D	F	H	S	E	S	E	V	T	H	T	M	P	N	A	D	R	W	S	P	T	R	T	L	L	Q	P	F	K	D	T	E	G	E	M	G	D	S	N	I	V	H	Q	D	L	I	S	R	S	C	P	V	L	T	C	G	P	S	E	D	N	P	Q	W	M	T	W	P	T	Y	Q	Y	M	N	T	P	I	T	A	I	D	D	C	K	H														
Y	T	T	V	M	L	Y	S	A	S	C	P	V	L	T	C	G	P	N	D	F	N	G	P	A	F	W	P	C	N	D	A	N	K	M	E	A	H	L	T	I	I	S	T	C	W	P	W	A	C	V	I	Q	Y	W	E	M	K	E	E	W	T	D	S	P	Q	H	A	L	R	C	K	P	Y	M	R	H	M	G	F	H	M	S	N	N	R	E	R	D	R	E	A	C	T	F	N	A	N	M	
E	F	D	N	F	E	D	G	V	L	A	H	N	N	F	T	G	K	A	E	N	F	Y	H	A	D	D	Y	S	L	O	S	E	T	I	E	P	R	S	N	E	K	I	K	G	A	R	G	S	V	E	K	R	G	H	T	Y	D	T	Y	K	P	N	G	W	Y	Q	Q	G	A	W	A	K	C	P	V	L	T	C	G	P	N	T	G	I	F	E	H	H	A	V	K	K	S	D	T	L	L		
E	P	F	K	F	Y	H	C	G	F	K	D	A	N	V	G	I	T	D	S	L	A	P	O	I	V	C	P	G	L	N	S	S	E	M	L	M	H	V	G	E	A	S	C	P	V	L	T	C	G	P	A	R	N	E	P	O	L	Y	M	H	L	N	R	Y	K	H	S	A	C	O	V	A	T	O	N	S	V	C	E	K	K	S	S	W	E	P	H	P											
R	N	M	T	A	T	A	R	S	R	V	C	H	F	K	S	R	C	A	V	T	T	I	M	K	D	I	V	P	E	T	E	R	W	P	F	G	K	Y	D	T	H	F	Y	P	O	K	C	I	V	C	D	H	R	I	D	D	B	A	H	I	L	O	G	F	C	S	D	D	Y	D	A	R	F	A	S	D	P	V	L	T	C	G	P	K	H	C	D	N											
R	C	S	M	E	S	K	V	I	F	C	L	Y	K	C	Y	E	A	D	E	F	R	D	T	F	N	C	E	A	S	C	P	V	L	T	C	G	P	T	T	G	A	S	F	S	F	A	Y	S	N	M	H	V	S	L	I	I	R	N	I	M	I	C	H	C	O	V	G	D	P	I	E	V	W	R	N	F	E	N	K	I	K	M	N	K	N	F	W	H	T	R	I	S	A	E	M				
W	D	L	P	L	A	F	P	W	S	N	A	F	F	M	V	F	I	N	S	W	W	F	P	K	M	D	T	P	L	K	I	T	G	W	E	F	H	K	K	P	Y	C	V	A	T	F	P	V	L	T	C	G	P	C	N	L	F	I	S	W	R	D	V	F	N	H	M	P	L	Y	G	S	P	D	I	C	Y	N	E	R	C	M	S	N	K	A	T	C	R	V	K	N	A	V	W	R			

Find the starting positions and the PSSM of the motif



Write a program to identify the starting position of the motif in each protein sequence. Write a report which includes: a) your understanding on EM algorithm, b) the core algorithm, c) the final result (starting positions and PSSM of the motif), d) the difficulties you met.

Hint:

- In short, we need to find the common (or similar) part shared by this group of proteins.
- How to represent a motif?
 - We first assume the motifs have a fixed length, and each protein has one motif.
 - Note that those similar parts are not necessarily to be exactly the same. We need a statistical model, such as the PSSM to describe them. (see PPT Multiple Sequence Alignment Page29-36 and PPT Motifs and domains Page34-36)
- How to apply EM?
 - Randomly initialize the starting positions for all the sequences.
 - Given the positions, you can extract the motif from each protein.

- c) Estimate the probability (or score) of each position to be the starting point given the motifs you extracted. (see PPT Motifs and domains Page36)
- d) Choose the positions that has the highest probability as the updated starting position.
- e) Go to b).

a) Expectation Maximization (EM) algorithm

■ Step 1: Initialization

First, I need to initialize your algorithm with a set of starting positions for the motif in each sequence. This can be done randomly or based on some heuristic.

■ Step 2: Expectation Step

Given the current estimate of the motif's starting positions, extract the motif instances from each sequence. A motif instance is the substring starting at the estimated position and having the specified length of the motif (in this case, 10 amino acids).

Using these instances, calculate a Position Specific Scoring Matrix (PSSM), which reflects the frequency of each amino acid at each position in the motif. This matrix represents the 'expected' motif based on the current starting positions.

■ Step 3: Maximization Step (M-step)

With the newly constructed PSSM, I can now evaluate every possible starting position in each sequence to determine how well it matches the expected motif. Calculate a score for each position that represents the likelihood of the motif starting there, given the PSSM.

Then update the starting position for each sequence to the one with the highest score. That is, for each protein sequence, slide the window of the motif length across the sequence, calculate the score based on the PSSM, and select the position with the highest score.

■ Step 4: Iteration

Repeat the E-step and M-step until the algorithm converges. Convergence can be determined when the changes between iterations are below a threshold, or the likelihood of the data given the model is no longer increasing significantly.

b) The Core Algorithm

```
#Expectation step
def expectation(pssm_matrix):
    expectation_matrix = []
    for i in range(sequences_coloumn_length):
        probability = 0
        for j in range(sequences_row_length - motif_length):
            for k in range(motif_length):
```

```

        probability +=
pssm_matrix[k][amino_acids.index(sequences[i][j + k])]
        probability = probability / motif_length
        expectation_matrix.append(probability)
    return expectation_matrix

```

#d) Implement the PSSM (Position-Specific Scoring Matrix) to describe the statistical model.

```

def pssm(motifs):
    pssm_matrix = []
    overall_frequency = [0] * amino_acids_length
    for i in range(motif_length):
        column = []
        for j in range(amino_acids_length):
            count = 0
            for k in range(sequences_coloumn_length):
                if amino_acids[j] == motifs[k][i]:
                    count += 1
            overall_frequency[j] += count
            column.append(count)

        column = [x / sequences_coloumn_length for x in column]
        pssm_matrix.append(column)

    overall_frequency = [x / (sequences_coloumn_length*motif_length) for x
in overall_frequency]
    pssm_matrix = [[math.log(x) if x > 0 else -10 for x in y] for y in
pssm_matrix]
    return pssm_matrix

```

#c) Implement the Expectation-Maximization algorithm to update the starting positions.

```

def maximization(expectation_matrix):
    updated_index = []
    for i in range(sequences_coloumn_length):
        x =
expectation_matrix.index(max(expectation_matrix[i*90:(i+1)*90])) % 90
        updated_index.append(x)
    return updated_index

```

c) the final result (starting positions and PSSM of the motif)

The final starting position is: [40, 23, 0, 57, 8, 73, 49, 85, 29, 46, 82, 4, 49, 79, 71, 23, 78, 41, 0, 43, 69, 82, 76, 3, 25, 5, 36, 82, 34, 74, 30, 89, 15, 54, 42, 49, 84, 66, 71, 22, 23, 11, 51, 31, 77, 45, 30, 83, 48, 7, 80, 45, 6, 0, 33, 84, 82, 86, 23, 79, 22, 21, 12, 68, 69, 81, 38, 68, 64, 89, 65, 47, 58, 53, 12, 56, 30, 12, 72, 16, 85, 57, 42, 86, 81, 30, 8, 6, 65, 32, 82, 24, 53, 11, 26, 15, 75, 44, 41, 56, 7, 27, 47, 34, 63, 46, 9, 32, 35, 13, 27, 88, 78, 8, 6, 76, 26, 26, 74, 44, 51, 7, 29, 87, 46, 64, 64, 72, 13, 5, 57, 19, 80, 10, 68, 1, 28, 64, 33, 82, 56, 88, 63, 41, 57, 31, 30, 29, 8, 51, 26, 29, 31, 58, 69, 84, 47, 1, 55, 88, 34, 31, 65, 66, 7, 7, 89, 76, 60, 58, 37, 73, 84, 87, 18, 6, 18, 14, 83, 1, 60, 16, 68, 42, 31, 10, 47, 8, 57, 32, 52, 51, 32, 59, 64, 2, 72, 26, 66, 58, 7, 5]

```
[[-0.12216763397420753, -4.605170185988091, -4.199705077879927, -5.298317366548036, -5.298317366548036, -4.199705077879927, -5.298317366548036, -10, -10, -10, -4.605170185988091, -10, -10, -5.298317366548036, -10, -4.605170185988091, -10, -4.605170185988091, -4.199705077879927, -4.605170185988091], [-5.298317366548036, -5.298317366548036, -5.298317366548036, -5.298317366548036, -5.298317366548036, -5.298317366548036, -10, -10, -4.605170185988091, -10, -4.605170185988091, -4.605170185988091, -5.298317366548036, -4.605170185988091, -10, -0.11093156070728166, -5.298317366548036, -10, -10, -4.605170185988091], [-4.605170185988091, -10, -4.605170185988091, -5.298317366548036, -0.11653381625595151, -5.298317366548036, -10, -4.605170185988091, -4.199705077879927, -4.605170185988091, -10, -10, -10, -5.298317366548036, -4.605170185988091, -10, -5.298317366548036, -5.298317366548036, -10, -3.912023005428146], [-5.298317366548036, -4.605170185988091, -5.298317366548036, -10, -10, -10, -4.605170185988091, -10, -10, -10, -3.912023005428146, -10, -0.07257069283483537, -5.298317366548036, -10, -10, -5.298317366548036, -5.298317366548036], [-10, -5.298317366548036, -10, -5.298317366548036, -5.298317366548036, -5.298317366548036, -10, -10, -10, -4.199705077879927, -10, -4.605170185988091, -5.298317366548036, -10, -10, -5.298317366548036, -4.605170185988091, -10, -0.06720874969344999], [-5.298317366548036, -5.298317366548036, -5.298317366548036, -5.298317366548036, -10, -4.199705077879927, -5.298317366548036, -5.298317366548036, -10, -10, -0.0779615414697118, -10, -10, -5.298317366548036, -5.298317366548036, -4.199705077879927, -10, -10, -5.298317366548036, -10], [-4.199705077879927, -5.298317366548036, -10, -4.605170185988091, -5.298317366548036, -10, -10, -5.298317366548036, -5.298317366548036, -4.199705077879927, -5.298317366548036, -10, -10, -5.298317366548036, -4.605170185988091, -0.11653381625595151, -10, -10, -4.605170185988091], [-4.605170185988091, -4.605170185988091, -5.298317366548036, -4.605170185988091, -0.11093156070728166, -5.298317366548036, -10, -4.605170185988091, -4.199705077879927, -4.605170185988091, -5.298317366548036, -10, -4.605170185988091, -4.199705077879927, -4.605170185988091, -5.298317366548036, -10, -10, -4.199705077879927, -10, -10, -10, -5.298317366548036], [-10, -10, -5.298317366548036, -4.199705077879927, -4.605170185988091, -4.605170185988091, -10, -0.11093156070728166, -5.298317366548036, -10, -10, -5.298317366548036, -5.298317366548036, -4.605170185988091, -5.298317366548036, -4.199705077879927, -5.298317366548036, -4.605170185988091, -10, -5.298317366548036], [-10, -5.298317366548036, -10, -4.199705077879927, -5.298317366548036, -4.199705077879927, -5.298317366548036, -10, -5.298317366548036, -5.298317366548036, -10, -5.298317366548036, -10, -5.298317366548036, -10, -10, -0.09431067947124129, -10, -3.912023005428146, -5.298317366548036, -5.298317366548036, -10]]
```

d) **the difficulties you met.**

- (ア) Although the dataset is small, it still is difficult for me to debug my program, especially involving several multi-dimensions matrices.
- (イ) I find it difficult to perform computations, using programming, involving multidimensional matrices, which often leads to errors.
- (ウ) Some problems occur during my programming, such as oversight of the boundaries of the matrices, normalizing the index and so on