

# 11300A Bioinformatics

## Homework 4

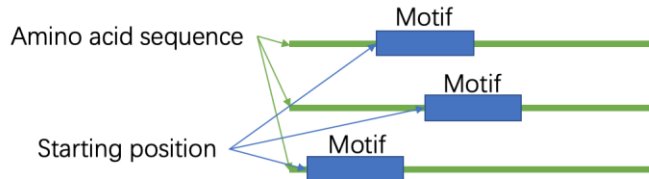
Jan. 10, 2023

### 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	V	I	A	M	G	M	V	N	H	G	S	N	C	I	D	P	V	G	E	D	T	V	W	W	H	A	L	N	H	H	L	E	I	V	F	H	S	E	I	R	D	R	A	S	D	S	E	F	N					
L	C	G	M	F	T	T	C	P	W	G	N	H	H	E	L	H	M	I	G	L	R	S	C	P	V	L	T	C	G	P	E	G	E	I	H	R	C	G	F	C	T	V	H	S	M	T	N	T	T	G	C	C	P	K	V	K	T	T	F	C	G	H	Y	S	V	Y	E	T	G	I	C	Q	F	A	W	F	V	W	G	E	N	H	T	H	R	E	E	N	P	G	A	V	M	K	I				
Y	S	C	M	L	T	C	G	P	T	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	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	I	A	H	L	T	I	I	S	T	C	W	P	W	A	C	V	I	Q	V	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	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	E	K	I	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	S	D	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	V	C	P	A	R	N	E	P	O	L	Y	M	H	L	N	W	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	H	M	T	A	N	K	E	R	V	C	H	F	K	R	W	C	A	V	T	I	M	K	D	I	V	P	E	T	E	R	W	P	F	G	K	V	D	T	H	F	Y	P	O	K	C	I	V	C	T	D	H	R	I	D	D	B	A	H	I	L	O	G	F	C	D	S	D	V	D	A	R	F	A	S	D	C	P	V	L	T	C	G	P	K	H	C	D	N											
R	C	S	M	E	G	V	I	F	C	L	Y	K	C	Y	E	A	D	E	F	R	D	R	T	F	N	C	I	A	S	C	P	V	L	T	C	C	P	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	G	N	A	F	F	M	V	F	I	N	S	W	W	P	F	K	M	D	T	P	L	K	T	G	W	E	F	H	K	K	P	Y	C	V	A	T	F	P	V	L	T	C	G	P	C	N	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.

Name: \_\_\_\_\_ Student ID: \_\_\_\_\_

- 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).