

CMSC423

Chapter 2 – Defining motifs

- Recap: Looking for something common in front of multiple genes
- How do we formalize what we look for?

Which one is a motif?

TTACCTTAAC
GATATCTGTC
ACGGCGTTCG
CCCTAAAGAG
CGTCAGAGGT

TTACCTTAAC
GTTATCTGTC
ACGTCGTTAG
CCCTCAAGAG
CTTCAGAGGC

CTACCTTAAC
CTTATCTGTC
ATGTCGTTAG
CTCTCGAGAG
CTTCCGTGTC

CTAACTTGAC
CTTATCTGTC
ATATCGTTAC
CTATCGTGAC
CTTACGTGTC

Option 1: count minority bases

	T	C	G	G	G	G	a	T	T	T	t	t
	c	C	G	G	t	G	A	c	T	T	a	C
	a	C	G	G	G	G	A	T	T	T	t	C
Motifs	T	t	G	G	G	G	A	c	T	T	t	t
	a	a	G	G	G	G	A	c	T	T	C	C
	T	t	G	G	G	G	A	c	T	T	C	C
	T	C	G	G	G	G	A	T	T	c	a	t
	T	C	G	G	G	G	A	T	T	c	C	t
	T	a	G	G	G	G	A	a	c	T	a	C
	T	C	G	G	G	t	A	T	a	a	C	C

SCORE(*Motifs*)

$$3 + 4 + 0 + 0 + 1 + 1 + 1 + 5 + 2 + 3 + 6 + 4 = 30$$

Option 2: compute entropy

$$H(p_1, p_2, \dots, p_n) = -\sum_{i=1}^N p_i \log_2(p_i)$$

	T	C	G	G	G	G	a	T	T	T	t	t	
Motifs	c	C	G	G	t	G	A	c	T	T	a	C	
	a	C	G	G	G	G	A	T	T	T	t	C	
	T	t	G	G	G	G	A	c	T	T	t	t	
	a	a	G	G	G	G	A	c	T	T	C	C	
	T	t	G	G	G	G	A	c	T	T	C	C	
	T	C	G	G	G	G	A	T	T	c	a	t	
	T	C	G	G	G	G	A	T	T	c	C	t	
	T	a	G	G	G	G	A	a	c	T	a	C	
	T	C	G	G	G	t	A	T	a	a	C	C	
PROFILE(Motifs)	A:	.2	.2	0	0	0	0	.9	.1	.1	.1	.3	0
	C:	.1	.6	0	0	0	0	0	.4	.1	.2	.4	.6
	G:	0	0	1	1	.9	.9	.1	0	0	0	0	0
	T:	.7	.2	0	0	.1	.1	0	.5	.8	.7	.3	.4

Entropy vs. counts

A Majority base = A
A SCORE(column) = 4
A
A Profile
C
A A – 0.6
C C – 0.2
G G – 0.1
T T – 0.1
A

entropy = 1.57

A Majority base = A
A SCORE(column) = 4
A
A Profile
C
A A – 0.6
C C – 0.4
C G – 0
C T – 0
A

entropy = 0.97

Next: Algorithms for finding motifs