Лабораторная работа № 5

Поиск регуляторного мотива

1 Мотив ДНК

Мотив ДНК определяется как подстрока строки ДНК, имеющая определенное биологическое значение. Мотив ДНК представляет собой небольшой кусок нуклеотидов А, Т, С или G длиной всего от 5 до 25 пар оснований. Известно, что регуляторный мотив слабо меняется в процессе эволюции,

Мотивы, присутствующие в экзонах (кодирующая часть генома), определяют структуру белка или метку белков, которые будут отправлены в определенные части клетки для осуществления таких процессов, как, например, фосфорилирование. Мотивы, присутствующие в интронах (которые составляют некодирующую часть генома), обычно представляют собой регуляторные последовательности, которые определяют степень экспрессии генов и сайты связывания белков.

Конкретные мотивы последовательностей могут функционировать как регуляторные последовательности, контролирующие биосинтез, или как сигнальные последовательности, которые направляют молекулу в конкретный сайт внутри клетки или регулируют ее созревание. Поскольку регуляторная функция этих последовательностей важна, считается, что они сохраняются в течение длительных периодов эволюции. В некоторых случаях эволюционное родство можно оценить по степени сохранности этих участков.

Мотив длиной k представляет собой k-мер. k-мер — это подстрока длины k — последовательность из k подряд идущих символов в строке (или нуклеотидов в последовательности ДНК, РНК или аминокислот в задачах биоинформатики).

Благодаря появлению доступных методов для массового прочтения последовательностей ДНК, стремительно растет объем прямых данных по ДНК-белковому узнаванию как in vivo, так и in vitro. Компьютерный анализ характерных ДНК- паттернов, мотивов, распознаваемых факторами транскрипции, потенциально позволяет изучать структуру регуляторных районов с однонуклеотидным разрешением. Однако, классические инструменты для анализа мотивов не справляются с возрастающими объемами данных и не учитывают специфику современных экспериментальных подходов. Поэтому требуются специальные программы для анализа регуляторных мотивов.

2 Мотивы последовательности и консенсусные последовательности

Оказывается, что даже выполняющие одну и туже функцию мотивы могут отличаться друг от друга. Наша задача найти консенсусную последовательность для этих мотивов.

Консенсусная последовательность (consensus sequence) — искусственная последовательность ДНК, содержащая в каждой позиции нуклеотид, наиболее часто встречающийся у нескольких гомологичных последовательностей. Обычно такие последовательности характерны для генов, кодирующих один и тот же белок у различных организмов.

Все реальные мотивы не должны отличаться от консенсусной последовательности более чем несколькими заменами.

Пример. Например, рассмотрим участок ДНК

Мы видим, что в этом участке 10 раз повторяется один и тот же 10-мер с точностью до 1-2 нуклеотидов. Эти 10-меры выделены заглавными буквами и жирным шрифтом, ошибочные нуклеотиды выделены красным. Расстояние Хэмминга, то есть количество разных символов в каждой позиции из двух строк одинаковой длины, в этом примере меньше или равно 2.

Учитывая набор мотивов, консенсусная последовательность представляет собой последовательность, полученную путем взятия наиболее частых нуклеотидов в каждом положении. Например, консенсусная последовательность

TAGATCTGAA

TGGATCCGAA

TAGACCCGAA

TAAATCCGAA

TAGGTCCAAA

TAGATTCGAA

CAGATCCGAA

TAGATCCGTA

TAGATCCAAA

TCGATCCGAA

представляет собой TAGATCCGAA.

Если мы напишем одну последовательность одну под другой и окрасим нуклеотиды в разный цвет, это будет легче увидеть: T — наиболее распространенный остаток в положении 1, A — наиболее распространенный остаток в положении 2 и так далее.

Таблица 1. Матрица Motifs

$N_{ar{0}}$	1	2	3	4	5	6	7	8	9	10
1	T	A	G	A	T	C	t	G	A	A
2	T	g	G	A	\mathbf{T}	\mathbf{C}	\mathbf{C}	G	A	A
3	T	A	G	A	\mathbf{c}	\mathbf{C}	\mathbf{C}	G	A	A
4	\mathbf{T}	A	a	A	\mathbf{T}	\mathbf{C}	\mathbf{C}	G	A	A
5	\mathbf{T}	A	G	g	\mathbf{T}	\mathbf{C}	\mathbf{C}	a	A	A
6	\mathbf{T}	A	G	A	\mathbf{T}	\mathbf{t}	\mathbf{C}	G	A	A
7	c	A	G	A	\mathbf{T}	\mathbf{C}	\mathbf{C}	G	A	A
8	\mathbf{T}	A	G	A	\mathbf{T}	\mathbf{C}	\mathbf{C}	G	\mathbf{t}	A
9	\mathbf{T}	A	G	A	\mathbf{T}	\mathbf{C}	\mathbf{C}	a	A	A
10	T	С	G	A	T	C	C	G	A	A

Следует отметить, что позиция 10 является наиболее консервативной (нуклеотид A полностью сохраняется в этих положениях), в то время как позиции 2 и 8 являются наименее консервативными.

В Python мы представим матрицу мотива как список строк Motifs. Будем назвать i-ю строку этой матрицы Motifs[i]; а j-й символ этой строки — Motifs[i][j].

Матрице Motifs поставим в соответствие число Score(Motifs) — сумму числа непопулярных (строчных) букв в матрице мотива.

Для нашего примера

$$Score(Motifs) = 1 + 1 + 1 + 1 + 1 + 1 + 1 + 1 + 2 + 1 + 0 = 10.$$

Построим теперь $4 \times k$ -матрицу, обозначаемую Count(Motifs) подсчета числа вхождений каждого нуклеотида в каждой столбец матрицы мотива. Элемент (i, j) матрицы Count(Motifs) содержит то число раз, которое нуклеотид с номером i появился в j-м столбце мотива.

Нуклеотид 1 2 $3 \ 4 \ 5$ 6 7 8 9 10 0 8 1 9 0 0 0 2 9 10 Α C 0 $0 \ 0 \ 1$ 9 9 $0 \quad 0$ 0 G 9 0 0 0 8 0 1 1 0 0 T 1 0 9 1 0

Таблица 2. Матрица Count(Motifs)

Один из способов представления такой матрицы в Python, состоит в том, чтобы создать список (list) для каждой строки матрицы, а затем организовать эти списки в больший словарь (тем самым создавая словарь, ключами которого являются нуклеотиды и значения — списки).

По матрице Count(Motifs) мы можем образовать консенсусную последовательность, обозначаемую Consensus(Motifs), из наиболее часто встречающихся нуклеотидов в каждом столбце матрицы Motifs. Если мы правильно выбрали мотивы, то Consensus(Motifs) и будет давать предполагаемый регуляторный мотив.

Например, на основе матрицы Count(Motifs) (см. таблицу 2) для матрицы Motifs (см. таблицу 1) консенсусной строкой является TAGATCCGAA.

На практике, для визуального представления мотивов часто используют логотип последовательностей — графическое представление консервативности каждой позиции в мотиве.

Приведем логотип последовательностей матрицы Motifs (см. таблицу 1 и рис. 1), построенный на сайте

https://weblogo.berkeley.edu/logo.cgi

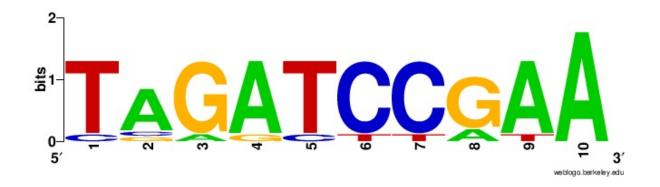


Рис. 1: Логотип последовательностей матрицы Motifs

Задания

- 1. Написать программу поиска консенсусной последовательности. Проиллюстрировать на примере из своего варианта. Вывести матрицу Motifs и матрицу Count(Motifs).
- **2.** Используя сайт https://weblogo.berkeley.edu/logo.cgi построить логотип последовательностей матрицы Motifs на примере из своего варианта.

Варианты

1. ATTATACTTA

AACCTGTGCT

CATTCACATT

TCAAAGCCAT

CTTCCAGAAA

GCTTAGCAAC

ACGAAAAAA

AAAAATTGTT

TAGATCTGAA

TATTTTTCAT

CTGAGAGCAT

TTAGAAAGGT

\circ	\mathbf{T}	α	A -	$\Gamma \Lambda$	Λ.	α	
Ζ.	١,	CA	Α.	IΑ	A	LΤ.	I C

AGTTCAACAA

ATATTTATCT

TGAACCAAGT

CCTTTTTTTC

TAAGCTTTTT

ACAACTTTGT

TCGGTTCAAA

AGAAATGCAA

TCCCTTATGT

TATTCCTTCA

TTTGTTTATA

3. AATCATTAAT

ATGATGGTTT

GAATGAGCTT

TTAAAGATGA

AAAAAAATCA

CCCTTTCATA

GTTGTTTAAA

GCAGCTTTTT

CTGTTAAATA

TGGAGGTCAT

TTTTTACAGA

AGAGAAACCG

4. AAATATACCA

AGATACTGTT

GCCTCCTAAT

TTGTACCATT

TGAAGTACCA

AAATTGACTA

AATACATTTT

CCTATTATGT

TTCTACATTT

TCAAGATGTA

CCTAAGATTG

CAGTGGGCCA

5. AGGCAAAGGG

GAAGTGAGCA

TGTCTTCACA

TGTCAGAGCA

GAAGACTGAG

AGTGAATGAG

GAAGTGCCAG

ACACTTTTAA

ACCATCAGGT

CTTGTTAGAA

CTCATTCATT

ACCATGAGAA

6. CAGCAAGGGG

GATATCCGCC

TCTATAATCA

ATCACCTTCC

ACTAATTGCC

TTCCTTAACA

CTGGGAATTA

ACAATTTTAC

ATGAGATTTG

GGTGGGGACA

CAGAGCTAAA

CTATATCAAC

7. TATCTACAAT

TTGCTTGATA

AGGAACAATT

TGTACACATT

CCTGAGTCCA

CCCAGAAGGA

GTGAAGTACA

GTCTGAATGG

AAAGTGTAGT

GGAAGAAGTT

TCTAGTAGCA

AAAAAGGGGA

8. CTATTCTCTG

AAAGCTTTGT

AATTAAACTA

AAGAGAAGAA

ACAAAGGTAC

AGCCAGAAAG

GAAGCCCTGG

TCATAGAAAG

TAAATAGAGT

GTTCACTCAG

AAGAACAAGG

TTGTGTACCA

9. CAACTATGGC

CAGTAAACTC

AGTGGAGAAG

AGGACTTAAA

CAATGGCTGG

GAAATGTAGA

GGGGTCACGA

ATCCTAAAAC

TCTGGACATG

AGATAGAATG

TTCACGTAAT

10. GCGGCCGCTG

CGTCCGCCAG

TAGCGGGTTG

CAGGCGCACC

CTCCCCTCCA

GGGCGGCCAC

GCAGCTGTCA

GTGCCGCCGC

CACTGCGAGG

CTGGAGCGGA

GCCCGGGTGG

CCGAGGGAGG

11. GGACCCCGCG

AGAGGGCCGC

GCGCCGGCCG

CCGCCGCCCC

GGCGCCCAGG

CTCGGTGCTG

GAGAGTCATG

CCTGTGAGCC

CTGGGCACCT

CCTGATGTCC

TGCGAGGTCA

CGGTGTTCCC

12. AAACCTCAGG

GTTGCCCTGC

CCCACTCCAG

AGGCTCTCAG

GCCCCACCCC

GGAGCCCTCT

GTGCGGAGCC

GCCTCCTCCT

GGCCAGTTCC

CCAGTAGTCC

TGAAGGGAGA

CCTGCTGTGT

13. GGAGCCTCTT

CTGGGACCCA

GCCATGAGTG

TGGAGCTGAG

CAACTGAACC

TGAAACTCTT

CCACTGTGAG

TCAAGGAGGC

TTTTCCGCAC

ATGAAGGACG

CTGAGCGGGA

AGGACTCCTC

14. TCTGCCTGCA

GTTGTAGCGA

GTGGACCAGC

ACCAGGGGCT

CTCTAGACTG

CCCCTCCTCC

ATCGCCTTCC

CTGCCTCTCC

AGGACAGAGC

AGCCACGTCT

GCACACCTCG

CCCTCTTTAC

15. ACTCAGTTTT

CAGAGCACGT

TTCTCCTATT

TCCTGCGGGT

TGCAGCGCCT

ACTTGAACTT

ACTCAGACCA

CCTACTTCTC

TAGCAGCACT

GGGCGTCCCT

TTCAGCAAGA

CGATGGCTGT

16. GCTCAGGCAG

CTGGCGCTCC

TCCTCTGGAA

GAACTACACC

CTGCAGAAGC

GGAAGGTCCT

GGTGACGGTC

CTGGAACTCT

TCCTGCCATT

GCTGTTTTCT

GGGATCCTCA

TCTGGCTCCG

17. ATAAAGAACT

ACCTGCGACT

GGGTAATTTA

TGAAGAAAAG

AGGTTTAATT

AACTCACAGT

TTCACAGGCT

GTACAGGTAG

CATGGCTGGG

AGACCTCAGG

AAACTTACAA

TCATGGCGGA

18. GCCCCAGGCT

CTGTCAACTG

CCCCATCTCA

AGTCTTCAAA

GCCTACTTCC

CATGAGTCAA

GAGTAACTCA

CTGTTCTGAA

GCTGTGTTAA

AGCATTGTAT

TACTCTGTTC

TCACACTGCT

19. AAGTCATTTT

GCAGCTAAGA

TTATTTGTCC

CCTGTTACAT

AACAGGGGGT

TATTCATAGA

GGTTATTCAT

AAACAGGTTA

TTCATCGAGT

TAAGTATAGA

ATGCAGGTCT

CCTGACCCTT

20. ATAAAAGGTC

TTACCTAATA

TATTTTACAA

TAAATACCCA

CATGGTACCC

TGAATCCATG

CTAAATTGAT

CTATCTCATA

TAGGCATAAT

GTGTAAAAGA

CTCTAGTGGT

CAAATGACTC

21. TAATGTATTA

 ${\bf ATGTTTTATC}$

TGGGGGTTCT

CTTTCTCCCT

TTACTCACTG

GCTAGAAGCC

TGTACAATTC

CTTCAAATTG

AGCTAATCTG

TAGGTACTCC

AAAATCACTG

CAGTTGATAG

22. AAAATATCTT

AAACAAAAAT

ATAGCCCATG

CCAGTCAAAC

GCTAGGGAAG

ATTATTTTCT

CAGAAAACCA

CATTGACAAT

ATCAGACACA

TGGAAGCAGA

ATATCTTAAA

TAATAGTTAA

23. GTTTATAGGC

AATAAGCTAT

TATCCTCTAG

GACTTAGCGA

ACTGGCTGTC

AGCCAACCTT

ATATTTGGAT

ACTCCTGGAA

GGATAAAATA

TGACAGTTAA

TTGGTTTATG

TGAGAAGAGT

24. ATTACTTGGG

CTACTCAGTG

AAGCTAATAC

CCAACATATG

ACCTTTCTTC

ACAGCCTTAT

TAATATCTTT

CACATCACAA

GTGGTGCCAG

AATTGAATAG

 ${\bf AGAGACTTTA}$

AAACAAGGCT

25. TTCAAAACTA

CTCAGCAAGT

CAGCTATGAA

AACCTTCAGA

CGTTTTCCAG

CAAATAGTGA

AAGGACATTA

AATGTTTTGG

TCATCTAAAG

AAACAAAGAG

CTGTTGAATG

GAAAAATAAA