

УДК 519.6+576.895.42+519.2

В. А. Гусев¹, Г. С. Лбов¹, Г. Л. Полякова^{1, 2}, В. С. Алтынцева², В. А. Габриэль²

¹ Институт математики СО РАН
пр. Акад. Коптюга, 4, Новосибирск, 630090, Россия

² Новосибирский государственный университет
ул. Пирогова, 2, Новосибирск, 630090, Россия

E-mail: vgus@math.nsc.ru; http://lbovgenady.narod2.ru
Polyakova@math.nsc.ru; Altynya@gmail.com, vitek-novosib@mail.ru

МЕТОДЫ ОБНАРУЖЕНИЯ ЛОГИЧЕСКИХ ЗАКОНОМЕРНОСТЕЙ В СТРУКТУРЕ ГЕНОМОВ *

Подтверждена гипотеза о наличии логических закономерностей в структуре геномов как микроорганизмов на примере *E. coli*, так и высших форм на примере X и Y хромосом человека. Для анализа рассматриваемых геномов был использован алгоритм полного перебора конъюнкций (*L*-грамм) с целью выявления логических закономерностей, обладающих высокой относительной частотой их встречаемости в бинарной последовательности. Приведено описание логико-вероятностных моделей для бинарных последовательностей и алгоритма обнаружения логических закономерностей в бинарной последовательности.

Ключевые слова: логико-вероятностная модель, бинарная последовательность, структура генома.

Введение

Обширная библиотека расшифрованных генетических последовательностей (GenBank) является в настоящее время объектом пристального внимания математиков. Как правило, работы по математическому анализу геномов посвящены применению различных математических методов для выявления регуляторных и кодирующих участков в структуре геномов [1–3]. Этот анализ основан на сопоставлении различных символьных последовательностей исследуемых геномов с паттернами ДНК, функции которых уже известны.

В работах [4–10] авторы проводили анализ структуры генетического кода, т. е. кодон-аминокислотного соответствия с использованием методов теоретико-группового анализа. Это позволило обнаружить неизвестные ранее молекулярным биологам и биохимикам закономерности в структуре кода.

Логично предположить, что соответствующие алгебраические и арифметические, т. е. символьные и числовые, закономерности генетического кода должны иметь отображения в нуклеотидной последовательности геномов. Для поиска числовых закономерностей необходимо представить стандартную символьную последовательность из *at* и *gc* пар в геноме в цифровом виде. Мы воспользовались данными работы [8], в которой показано, что молекулярные массы *at* и *gc* пар в составе двойной спирали равны соответственно 259 и 260 независимо от ориентации, т. е. принадлежности нуклеотида к конкретной нити ДНК. Таким образом, чередование *at* и *gc* пар в двойной спирали соответствует чередованию только двух

* Работа выполнена при финансовой поддержке РФФИ (проект № 10-01-00113-а).

чисел. Следовательно, для анализа рассматриваемых символьных последовательностей можно их представить в виде последовательностей нулей и единиц. Поиск возможных закономерностей в таком бинарном ряду чисел был проведен в рамках подхода к анализу эмпирической информации, изложенного в монографиях [11–13]. Подход сводится к построению логико-вероятностной модели объекта исследования. Под логико-вероятностной моделью понимается список логических закономерностей, обладающих достаточно большой прогнозирующей способностью (см. ниже). Целью данной работы является поиск логических закономерностей в бинарных последовательностях, сопоставленных кодирующим фрагментам генома прокариот *E. coli*, а также фрагментам X и Y хромосом эукариотического генома вида *Homo Sapiens*.

Обнаружение логических закономерностей

Метод поиска (обнаружения) логических закономерностей в бинарной последовательности состоит из выбора наилучшего разбиения, полученного для слова длины l и его окружения длины s , по некоторому критерию. Для каждого разбиения определяются частоты перехода из множества окружений слова x в множество слов y для всех пар $\langle x, y \rangle$ на основе анализа исходной последовательности; затем определяется порядок для всех пар по их частотам появления и выделяются пары с наибольшими частотами, которые и будут логическими закономерностями.

Пусть имеется последовательность *gene* длины L : $gene = (nucl_1, nucl_2, \dots, nucl_L)$, где $nucl_m \in D_{nucl}$, $m = 1, \dots, L$ и множество D_{nucl} представляет собой неупорядоченный набор значений элементов последовательности $D_{nucl} = \{a, t, c, g\}$.

Переходим к бинарной последовательности заменой в исходной последовательности символов a и t на 0, c и g на 1. Получаем последовательность $b = (b_1, b_2, \dots, b_L)$, где $b_m \in \{0, 1\}$, $m = 1, \dots, L$.

Под словом длины l в алфавите $\{0, 1\}$ понимается бинарная последовательность длины l (например, при $l = 4$ «слово» имеет вид 0110; заметим, что число возможных слов длины l в этом случае равно $2^4 = 16$). Рассматриваем слова длины l .

Предположим, что частота возникновения слова $w^k = (b_k, b_{k+1}, \dots, b_{k+l-1})$ длины l , начинающегося с k -й позиции в последовательности b , зависит только от некоторых из s_1 ближайших слева и s_2 ближайших справа элементов («окружения»), т. е. зависит от некоторых из элементов $b_s^k = (b_{k-s_1}, b_{k-s_1+1}, \dots, b_{k-1}, b_{k+1}, \dots, b_{k+l+s_2-1})$, где $s = s_1 + s_2$; s_1, s_2 – некоторые параметры, $k = s_1 + 1, \dots, L - l - s_2 + 1$. Например, при $s_1 = 3, s_2 = 2$ для указанного выше слова «окружение» из 3 ближайших слева и 2 ближайших справа элементов может иметь вид 110(0110)01.

Логической закономерностью для указанного слова называется логическое высказывание на символах окружения, которое с большой частотой характеризует данное слово. Например, для указанного слова 0110 высказывание «(слева в 1-й позиции стоит 1) и (справа в 1-й позиции стоит 0)» появляется в рассматриваемой бинарной последовательности с относительной частотой $\bar{P} \geq \delta = 0,95$, δ – параметр; а для всех других слов появляется с частотой, близкой к нулю.

Сопоставим каждому участку b_s^k набор значений некоторых бинарных переменных $X = X_1, X_2, \dots, X_n$, $n = s_1 + s_2$. $D_j = \{0, 1\}$ – область определения переменной X_j . Пусть $x^k = X(b_s^k) = (X_1(b_s^k), X_2(b_s^k), \dots, X_n(b_s^k))$; $X_j(b_s^k)$ – значение переменной X_j для участка b_s^k . Для указанного выше примера окружение имеет вид: $(x_1 = 1, x_2 = 1, x_3 = 0, x_4 = 0, x_5 = 1)$.

Частоту встречаемости данного слова w длины l в последовательности b определим как $\bar{P}(w) = \frac{N_w}{N}$, где $N = L - l - s + 1$ – число всех слов длины l (совпадающих и несовпадающих),

которые можно получить из последовательности b сдвигом на один символ с позиции $k = s_1 + 1$ в последовательности b до позиции $k = L - l - s_2 + 1$; N_w – число повторов слова w среди всех таких N слов. Для надежности метода необходимо выбрать среди всевозможных слов длины l слово w_* наиболее высокой частоты встречаемости.

Сопоставим слову $w^k = (b_k, b_{k+1}, \dots, b_{k+l-1})$, начинающегося с k -й позиции в последовательности b , $k = s_1 + 1, \dots, L - l - s_2 + 1$, значение целевой (прогнозируемой) переменной Y . Будем считать, что $y^k = Y(w^k) = 1$, если слово $w^k = w_*$; $y^k = Y(w^k) = 2$, если слово $w^k \neq w_*$.

Сопоставив каждому значению x^k значение y^k , получим таблицу данных $v = \{x^k, y^k\}$, размерностью $n \times N$, где $n = s_1 + s_2$, $N = L - l - s + 1$. Можно определить по таблице данных число $N_{(1)}$ объектов первого образа и число $N_{(2)}$ объектов второго образа.

Требуется по этим наблюдениям найти логические закономерности, обладающие большой прогнозирующей способностью, для предсказания значения y в зависимости от «окружения» x . Множество таких закономерностей представляет логико-вероятностную модель, отражающую причинно-следственные взаимосвязи между характеристиками. В процессе построения закономерностей автоматически отбираются наиболее информативные характеристики.

Задача обнаружения всех закономерностей является N_p – трудной задачей. Для обнаружения закономерностей используются алгоритмы класса ТЕМР [11–14], которые дают возможность значительно сократить время вычислений, учитывать разнотипность переменных, перебирать конъюнкции различной длины. Эти алгоритмы обнаруживают все логические закономерности на реальных таблицах за приемлемое время.

Обозначим $J(a, E_j)$ предикат, принимающий значения «истина» или «ложь». Предикат $J(a, E_j)$ эквивалентен утверждению: $X_j \in E_j$, $a \in \Gamma$ – объект из некоторой генеральной совокупности, описываемый характеристиками X_1, \dots, X_n, Y ; E_j является подмножеством множества значений D_j , $j = 1, \dots, n$.

Назовем $S(a, E) = J(a, E_{j_1}) \wedge \dots \wedge J(a, E_{j_d})$ конъюнкцией длины d . Областью истинности конъюнкции $S(a, E)$ является подмножество $E = \prod_{i=1}^d E_{j_i}$, $E_{j_i} \subset D_{j_i}$. Обозначим через μ нормированную меру подмножества E . Для любой конъюнкции $S(a, E)$ можно определить по таблице данных v число объектов первого образа $N_{(1,S)}$ и число объектов второго образа $N_{(2,S)}$, на которых указанная конъюнкция истинна.

Конъюнкцию $S(a, E)$ будем называть *логической закономерностью*, с большой вероятностью характеризующей первый образ, если выполняются неравенства: $\frac{N_{(1,S)}}{N_{(1)}} \geq \delta$, $\frac{N_{(2,S)}}{N_{(2)}} \leq \beta$, где δ и β – некоторые параметры; $0 \leq \beta < \delta \leq 1$. Чем больше δ и меньше β , тем сильнее логическая закономерность. Множество всех закономерностей обозначим через S^* .

Конъюнкцию $S(a, E)$ будем называть *потенциальной логической закономерностью* для первого образа (обозначим ее через S'), если выполняются неравенства: $\frac{N_{(1,S)}}{N_{(1)}} \geq \delta$, $\frac{N_{(2,S)}}{N_{(2)}} \leq \beta$.

Множество потенциальных закономерностей обозначим через S' . Очевидно, что из $S' \in S'$ можно получить закономерность S^* последовательным присоединением предикатов, т. е.

$S' \wedge J(a, E_j) \wedge \dots$; если для некоторой конъюнкции $S(a, E)$ выполняется неравенство $\frac{N_{(1,S)}}{N_{(1)}} \geq \delta$, то конъюнкция S по определению не является закономерностью и присоединение

к ней какого-либо предиката не даст закономерности (множество таких конъюнкций обозначим через \mathbf{S}). Таким образом, любая конъюнкция $S(a, E)$ может быть трех типов: \mathbf{S}^* , \mathbf{S}' , \mathbf{S} .

Алгоритм обнаружения логических закономерностей состоит в последовательном выполнении следующих шагов.

На *первом шаге* рассматриваются всевозможные конъюнкции длины один, т. е. конъюнкции вида $S(a, E) = J(a, E_j)$, E_j является подмножеством множества значений D_j , $j = 1, \dots, n$. Если $S(a, E) \in \mathbf{S}^*$, то она включается в список закономерностей и соответствующее подмножество E_j исключается из дальнейшего перебора; если $S(a, E) \in \mathbf{S}'$, то соответствующее подмножество E_j оставляется для дальнейшего перебора; если $S(a, E) \in \mathbf{S}$, то соответствующее подмножество E_j исключается из дальнейшего перебора. Обозначим через Q_j^1 множество подмножеств E_j , оставленных для дальнейшего перебора после выполнения первого шага алгоритма.

На *втором шаге* рассматриваются всевозможные конъюнкции длины два, т. е. конъюнкции вида $S(a, E) = J(a, E_i) \wedge J(a, E_j)$, $i \neq j$, $E_i \in Q_i^1$, $E_j \in Q_j^1$. Если $S(a, E) \in \mathbf{S}^*$, то соответствующие подмножества E_i и E_j исключаются из дальнейшего перебора и соответствующая конъюнкция включается в список закономерностей; если $S(a, E) \in \mathbf{S}'$, то соответствующие подмножества E_i и E_j оставляются для дальнейшего перебора; если $S(a, E) \in \mathbf{S}$, то соответствующие подмножества E_i и E_j исключаются из дальнейшего перебора. Аналогично обозначаем Q_j^2 множество подмножеств E_j , оставленных для дальнейшего перебора после выполнения второго шага алгоритма.

Далее, аналогично рассматриваются конъюнкции длины три, четыре, пять и т. д. В результате работы алгоритма получаются конъюнкции небольшой длины. Например, максимальная длина полученных конъюнкций в задаче, описанной ниже, не больше 6.

Закономерности в генетических последовательностях

В результате работы алгоритма было найдено несколько закономерностей. Вероятность образования таких закономерностей при условии равномерного распределения $P(S | H_0)$ различна. Из всех таких закономерностей естественно считать наилучшей ту, для которой эта вероятность минимальна.

Меру μ можно рассматривать как вероятность попадания в область E при равномерном распределении, $E \subset D$; $1 - \mu$ – как вероятность попадания в область $D \setminus E$. Следовательно, вероятность $P(S | H_0)$ образования закономерности заданной длины при известных $N_{(1,S)}$, $N_{(2,S)}$, $N_{(1)}$, $N_{(2)}$ и μ может быть вычислена следующим образом:

$$P(S | H_0) = C_{N_{(1)}}^{N_{(1,S)}} \mu^{N_{(1,S)}} (1 - \mu)^{N_{(1)} - N_{(1,S)}} \cdot C_{N_{(2)}}^{N_{(2,S)}} \mu^{N_{(2,S)}} (1 - \mu)^{N_{(2)} - N_{(2,S)}}.$$

Чем больше длина конъюнкции, тем меньше мера μ и меньше вероятность $P(S | H_0)$; следовательно, при одинаковых значениях $N_{(1,S)}$, $N_{(2,S)}$, $N_{(1)}$, $N_{(2)}$ предпочтительней будут конъюнкции большей длины.

Был проанализирован весь геном *E. coli* (4 266 генов, GenBank¹). Приведем в качестве примера только некоторые из полученных результатов.

Для кодирующей последовательности гена *thrL E.coli* (длина 63 символа) приведем гистограммы частот встречаемости слов в соответствующей бинарной последовательности. Примеры гистограмм, отражающих среднюю частоту встречаемости слов соответствующей длины, приведены на рис. 1, 2.

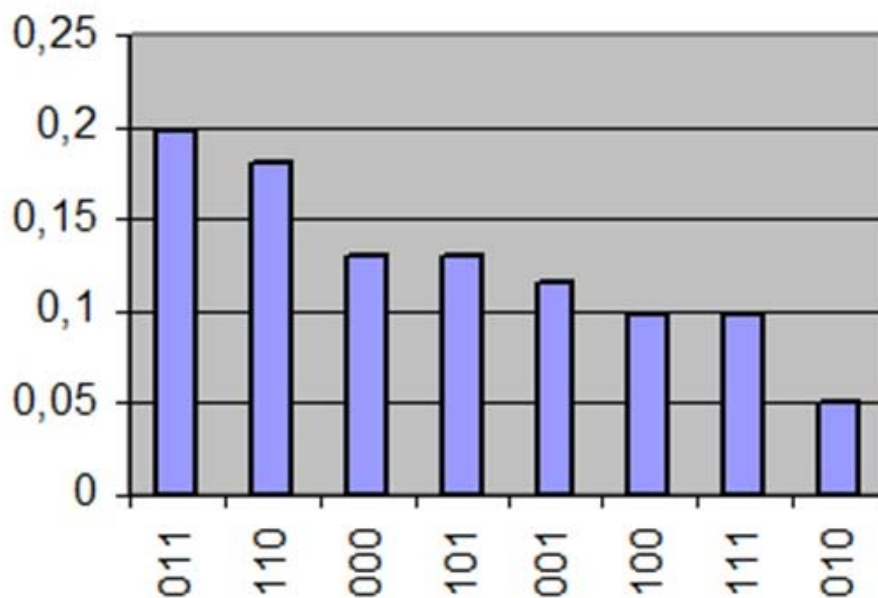


Рис. 1. Частоты встречаемости слов длины 3 в бинарной последовательности

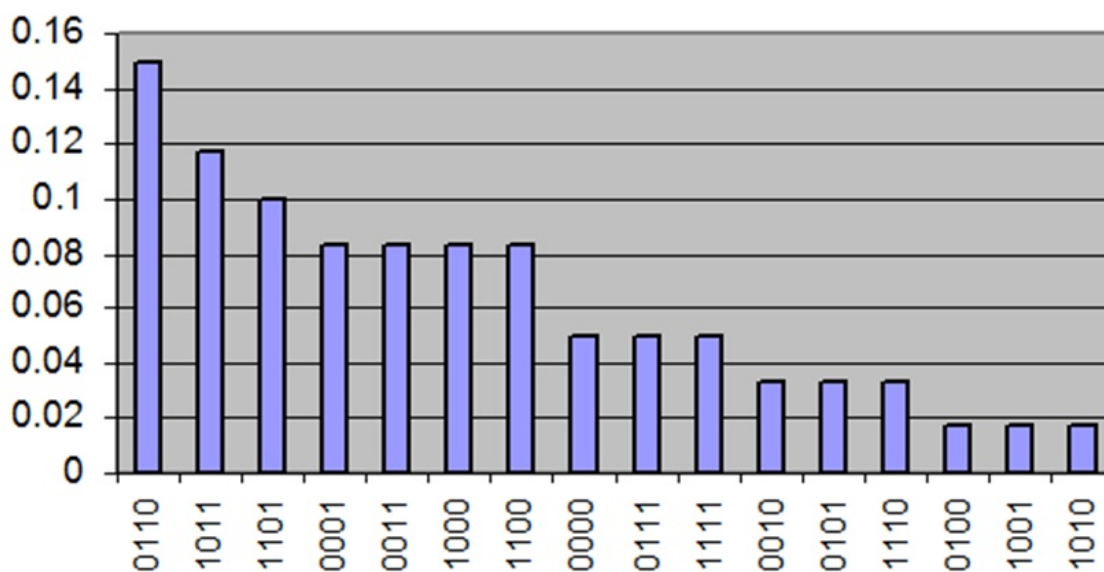


Рис. 2. Частоты встречаемости слов длины 4 в бинарной последовательности

¹ <http://www.ncbi.nlm.nih.gov/> (GenBank)



Рис. 4. Относительная частота встречаемости закономерности в хромосоме Y



Рис. 5. Относительная частота встречаемости закономерности в хромосоме X

рованные в соответствии с частотами ат и гс пар в исходных генетических последовательностях. Было проанализировано по 100 случайных последовательностей для каждого исследуемого гена. Ни в одной из них закономерностей не было найдено.

Следующая часть статьи посвящена проверке гипотезы о том, что в структуре генома человека, в частности в специальной паре хромосом, определяющих половую принадлежность, также содержатся логические закономерности.

Исследования проводились при длине слова от 3 до 7 символов. В результате были выделены наиболее часто встречающиеся закономерности в хромосомах X и Y. Затем в случайной бинарной последовательности определена относительная частота встречаемости для полученных закономерностей.

На рис. 4, 5

Рис. 1 приведены графики частот встречаемости выявленных закономерностей при $l = 5$, $s = 6$ и $N = 7000$. Отражены частоты появления слов, состоящих из пяти символов (по оси абсцисс средняя строка – пример выявленных закономерностей), с соответствующим окружением (по оси абсцисс первая и третья строки). Нижняя кривая соответствует встречаемости аналогичных комбинаций слов и окружения в рандомизированных последовательностях аналогичной длины.

Как видно из графиков, относительные частоты встречаемости закономерностей в хромосомах X и Y человека на порядок отличаются от встречаемости аналогичных комбинаций в случайных последовательностях. Это говорит о том, что полученные пары с большой частотой встречаемости практически не могут получиться из случайной последовательности. Поэтому в соответствии с гипотезой такие пары являются закономерностями в исходной последовательности.

Заключение

Анализ частот встречаемости слов различной длины в последовательностях показал, что частота слов в генетической последовательности значительно отличается от частот слов в случайной последовательности. Это дает основание утверждать, что наблюдаемые в геномах закономерности являются истинными. Следует особо подчеркнуть, что найденные закономерности в кодирующих последовательностях генома *E. coli*, а также во фрагментах X и Y хромосом человека имеют семантическую природу и непосредственно не связаны с триплетной структурой генома.

Благодарности

Авторы выражают признательность сотрудникам Института цитологии и генетики СО РАН В. А. Лихошваю и Ю. Г. Матушкину за конструктивную дискуссию в процессе постановки задачи.

Список литературы

1. Орлов Ю. Л. Анализ регуляторных геномных последовательностей с помощью компьютерных методов оценок сложности генетических текстов: Автореф. дис. ... канд. биол. наук. Новосибирск, 2004. 35 с.
2. Abnizova I., Schilstra M., te Boekhorst R., Nehaniv C. L. A Statistical Approach to Distinguish between Different DNA Functional Parts // WSEAS Transactions on Computational Methods. 2003. Vol. 2. Is. 4. P. 1188–1196.
3. Abe T., Kanaya S., Kinouchi M., Ichiba Y., Kozuki T., Ikemura T. Informatics for Unveiling Hidden Genome Signatures // Genome Res. 2003. Vol. 13 (4). P. 693–702.
4. Duplij D., Duplij S. Determinative degree and nucleotide content of DNA strands // Biophys. Bull. 2000. Vol. 497. P. 1–7.
5. Jimenez-Montano M. A., de la Mora-Basanez C. R., Poschel T. The Hypercube Structure of the Genetic Code Explains Conservative and Non-Conservative Aminoacid Substitutions in vivo and in vitro // BioSystems. 1996. Vol. 39. P. 117–125.
6. Jimenez-Montano M. A. Protein Evolution Drives the Evolution of the Genetic Code and Vice Versa // BioSystems. 1999. Vol. 54. P. 47–64.
7. Negadi T. Rumer's Transformation in Biology as the Negation in Classic Logic // Int. Journ. of Quant. Chem. 2003. Vol. 94. P. 65–82.
8. Shcherbak V. I. Arithmetic inside the universal genetic code // BioSystems 2003. Vol. 70. P. 187–209.
9. Карасев В. А. Генетический код: новые горизонты. СПб.: ТЕССА, 2003. 116 с.
10. Гусев В. А. Арифметика и алгебра в структуре генетического кода, логика в структуре генома и биохимическом цикле самовоспроизводства живых систем // Информационный вестник ВОГиС. 2005. Т. 9, № 2. С. 153–161.

11. Лбов Г. С. Методы обработки разнотипных экспериментальных данных. Новосибирск: Изд-во Наука, 1981. 160 с.
12. Лбов Г. С., Старцева Н. Г. Логические решающие функции и вопросы статистической устойчивости решений. Новосибирск: Изд-во Ин-та математики, 1999. 212 с.
13. Лбов Г. С., Бериков В. Б. Устойчивость решающих функций в задачах распознавания образов и анализа разнотипной информации. Новосибирск: Изд-во Ин-та математики, 2005. 218 с.
14. Лбов Г. С., Полякова Г. Л. Метод прогнозирования в классе логических решающих функций // Вестн. Сибирского государственного аэрокосмического университета имени академика М. Ф. Решетнева. 2010. Вып. 5 (31). С. 42–45.

Материал поступил в редколлегию 20.12.2011

V. A. Gusev, G. S. Lbov, G. L. Polyakova, V. S. Altynceva, V. A. Gabriel

**METHODS FOR THE DISCOVERY OF LOGICAL REGULARITIES
IN THE STRUCTURE OF THE GENOMES**

The hypothesis of the presence of logical regularities in the structure of the genomes of organisms as in the example *E.coli*, and the higher forms by the example of X and Y chromosomes of *Homo sapiens*. For the analysis of these genomes has been used exhaustive search algorithm conjunctions (*L*-gram) for the discovery of logical regularities, which have high relative frequency of their occurrence in a binary sequence. An algorithm for the discovery of logical regularities in a binary sequence is proposed. The examples of the obtained logical-and-probabilistic models are given.

Keywords: logical-and-probabilistic model, binary sequence, structure of the genome.