

1 УДК 621.391 : 519.218.5

2 © 20?? г. Г.А. Хазиев, О.А. Зверков, С.А. Пирогов, А.В. Селиверстов, В.А.
3 Любецкий

4 **МИНИМАЛЬНОЕ РАССТОЯНИЕ МЕЖДУ СЛУЧАЙНОЙ**
5 **ПОСЛЕДОВАТЕЛЬНОСТЬЮ И**
6 **ЕЁ СОВЕРШЕННЫМИ ПАЛИНДРОМАМИ**¹

7 Для произвольной последовательности x часто рассматривают нормированное
8 длиной x минимальное расстояние Левенштейна $\text{impr}(x)$ между x и переменным
9 словом вида $ws(w)$, где w – любой префикс в x , а $s(\cdot)$ – перестановка букв
10 в w в обратном порядке с заменой их на «двойственные» буквы. Например,
11 взаимозаменой А с Т и G с C. Свойства важной теоретически и в приложениях
12 функции $\text{impr}(x)$ остаются математически не изученными. Поэтому мы прове-
13 ли обширные компьютерные эксперименты, касающиеся поведения $\text{impr}(x)$ на
14 случайных последовательностях. В них наблюдается, что доля последовательностей
15 x , близких к совершенным палиндромам $ws(w)$, очень мала; медиана и
16 среднее значений для $\text{impr}(x)$ существенно зависят от вероятностей отдельных
17 букв. У длинных последовательностей x эти характеристики слабо зависят от
18 длины x , а для коротких, напротив, значимо зависят от длины. Но они весьма
19 зависят от GC-состава; и среднеквадратичное отклонение мало зависит от
20 него. Рассмотрены и другие характеристики. Результаты применимы, в частности,
21 в задачах биоинформатики. Например, при поиске сайтов кооперативного
22 связывания транскрипционных факторов и промоторов.

23 *Ключевые слова:* палиндром, случайная последовательность, GC-состав, биоин-
24 форматика.

25 **DOI:** 10.31857/S05552923??, **EDN:** ??

26 **§ 1. Введение**

27 В работе рассматривается минимальное нормированное длиной случайной после-
28 довательности x , редакционное расстояние между x и последовательностями специ-
29 ального вида (говорят «словами»), у которых префиксы те же самые, что у x . Эти

¹ Работа выполнена в рамках государственного задания ИППИ РАН, утвержденного Минобрнауки России.

30 слова однозначно вычисляются по каждому префиксу w в x и имеют вид $ws(w)$, где
31 функция (\cdot) определена ниже. Само расстояние обозначается $\text{imp}(x)$ как функция
32 от x .

33 В работе рассматриваются последовательности в четырёх буквенном алфавите.
34 В биоинформатическом контексте буквы называются нуклеотидами и обозначаются
35 A, C, G, T ; для них определяются так называемые комплементарные пары $\{A, T\}$ и
36 $\{G, C\}$. Точнее, определяется отношение комплементарности, которое записывается
37 как $c(A) = T, c(T) = A, c(C) = G$ и $c(G) = C$; и для любого слова $w = y_1 \dots y_k$ (в
38 роли префикса) $c(w)$ определяется как $c(w) = c(y_k) \dots c(y_1)$. Хотя результаты при-
39 водятся для нуклеотидных алфавита и комплементарности, предлагаемая методика
40 пригодна для любых алфавита и отношения комплементарности на нём; нуклео-
41 тидный случай функции $\text{imp}(x)$ выбран вследствие её широкой востребованности
42 в биоинформатике. Слово (так говорят, чтобы отличать его от последовательности
43 x) вида $ws(w)$, для любого подслова w , называется совершенным палиндромом. От-
44 куда следует, что $|w| = |c(w)|$ и такое слово чётной длины; где $|\cdot|$ – длина слова
45 или последовательности. В биоинформатике это определение выражают словами:
46 a совпадает с комплементарной к ней последовательностью, читаемой в обратном
47 направлении [1]. Несовершенным палиндромом назовём последовательность x близ-
48 кую по расстоянию к совершенному палиндрому $ws(w)$, где w – некоторый префикс
49 в x . Здесь подразумевается достаточная близость и $ws(w)$, но выбор границы свер-
50 ху (говорят: порога) этой близости зависит от прикладной задачи. Ниже предложен
51 вариант такого порога, но приведённые компьютерные эксперименты относятся к по-
52 следовательностям без ограничения на порог. Поиск несовершенных палиндромов
53 и выбор этого порога и весьма нетривиален и важен во многих прикладных зада-
54 чках, включая биофизические. В частности, несовершенные палиндромы, близкие по
55 последовательности к совершенным, могут служить сайтами кооперативного связы-
56 вания с ДНК транскрипционных факторов, образующих гомодимеры [2]; такие сай-
57 ты обычно несовершенные. Несовершенные палиндромы, близкие к совершенным,
58 входят как подпоследовательность в состав многих некодирующих РНК [3]; они рас-
59 пространены в геномах эукариот [4, 5] и прокариот [6]. Вместо префиксов аналогич-
60 ным образом рассматриваются суффиксы или даже слова, определённым образом
61 расположенные в x . Расстояние между произвольной последовательностью x и соот-
62 ветствующим ей совершенным палиндромом $ws(w)$ определяется как редакционное
63 расстояние (расстояние Левенштейна) [7], которое затем нормируется длиной после-
64 довательности x . Результаты приведены для случая, когда расстояние Левенштейна
65 не использует веса. Это расстояние между двумя последовательностями, как и их
66 наибольшая общая подпоследовательность, вычисляются за квадратичное время от
67 их длины. Известны алгоритмы вычисления редакционного расстояния за немного
68 меньшее время, чем квадратичное [8, 9]; математическое ожидание длины наиболь-

69 шей общей подпоследовательности двух случайных двоичных последовательностей
70 рассмотрено в [10]. GC-состав нуклеотидной последовательности определяется как
71 доля в ней букв G и C; можно аналогично рассматривать AT-состав. Напомним, что
72 GC-состав одинаков для обеих цепей ДНК. Например, GC-состав всего генома чело-
73 века 41%. Если в ДНК человека рассмотреть короткое окно, то в нём GC-состав ме-
74 няется в широких пределах: в окне длиной несколько сотен нуклеотидов GC-состав
75 меняется от 0% до 100%; а в окне длиной 100 тысяч нуклеотидов GC-состав меняет-
76 ся от 35% до 60%. Мы провели компьютерные эксперименты на основе описанного
77 ниже метода, выясняющие свойства минимального расстояния между случайной по-
78 следовательностью x и соответствующими ей совершенными палиндромами $ws(w)$,
79 которые частично приведены в разделе «Результаты».

80

§ 2. Метод

81 В [11] нами предложен алгоритм, который по последовательности x получает
82 список длин префиксов w в x (от начала x), для которых нормированное на длину
83 x редакционное расстояние dist между всей $x = wz$ и совершенным палиндромом
84 $ws(w)$ достигает минимума, вычисляемого по всем общим префиксам w в x и $ws(w)$.
85 Это расстояние, обозначается $\text{imp}(x)$. Наш алгоритм выдаёт и функцию $\text{imp}(x)$ –
86 значение минимума для последовательности x , равное

$$\text{imp}(x) = \frac{\min\{\text{dist}(x, ws(w)) \mid x = wz\}}{|x|}$$

87 Обычно последовательности x и $ws(w)$ разной длины. Отметим равенство

$$\text{imp}(x) = \text{imp}(c(x)).$$

88 Значение $\text{imp}(x)$ говорит, насколько произвольная последовательность x отличается
89 от ближайшего к ней совершенного палиндрама вида $ws(w)$, где w пробегает все пре-
90 фиксы в x . Для x с чётной длиной выполняется $\text{imp}(x) \leq \frac{1}{2}$; для x с нечётной длиной
91 выполняется $\text{imp}(x) \leq \frac{1}{2} + \frac{1}{2|x|}$, где $|x|$ – длина последовательности x . Например, для
92 x из одной буквы $\text{imp}(x) = 1$. В [12] получен алгоритм поиска (связной) строки u (в
93 последовательности x), которая по редакционному расстоянию достаточно близка к
94 какому-то совершенному палиндрому; и порог этой близости включён в квадратич-
95 ную оценку сложности алгоритма. Эта и наша задачи важны в биоинформатике и в
96 других приложениях; они различаются как глобальное, у нас, и локальное, в [12], вы-
97 равнивания последовательностей. В прикладных задачах минимальное отклонение
98 совершенного палиндрама $ws(w)$ от всей последовательности x с общим префиксом
99 w может быть велико, и величина отклонения имеет, в частности, биофизическое
100 значение, как показано, например, в [2]. Совершенный палиндром $ws(w)$, на кото-
101 ром достигается минимальное отклонение от x с произвольным общим префиксом w

102 в x , назовём оптимальным палиндромом для x . Нами доказано, что за квадратичное
103 время от длины исходной последовательности x наш алгоритм строит указанный в
104 начале этого раздела список длин префиксов, в то время как полный перебор всех
105 префиксов в x с вычислением редакционного расстояния для каждого из них тре-
106 бует кубического времени. Компьютерная реализация на языке Python алгоритма,
107 описанного в [11], доступна по адресу <http://lab6.iitp.ru/-/pali> вместе с примерами
108 её использования.

109 Компьютерная генерация выборки независимых случайных последовательностей
110 x с данной длиной n выполнялась по случайной величине, у которой вероятности
111 букв G и C равны $\frac{gc}{2}$, а букв A и T равны $\frac{1-gc}{2}$, где gc – параметр, $0 \leq gc \leq 1$. С ростом
112 объёма выборки среднее GC-составов отдельных последовательностей x из выборки
113 стремится к gc , а среднеквадратичное отклонение GC-составов выборки стремит-
114 ся к $\sqrt{\frac{gc \cdot (1-gc)}{n}}$. В этом смысле параметр gc характеризует GC-состав случайной
115 последовательности x с данной длиной n . Разработан скрипт, который генериру-
116 ет независимые выборки одинакового заданного объёма из независимых случайных
117 последовательностей x с заданной длиной n и данным gc . Ниже приводятся резуль-
118 таты, усреднённые по выборкам объёма сто тысяч при длине последовательностей
119 $n = 1000$, и объёма миллион при длине последовательностей $n = 100$ или $n = 50$.
120 Увеличение объёма и длины $n \geq 50$ приводят к качественно сходным результатам.

121 По выборке (переменная x пробегает её элементы) образуется числовая после-
122 довательность значений функции $\text{impr}(x)$ или другой заданной функции $y = f(x)$,
123 у которой одинаковым значениям приписывается кратность $f^1(y)$, т.е. доли вычис-
124 ляются с учётом кратности, по соответствующему числу x -ов. В этом смысле ниже
125 говорится о характеристиках функции $\text{impr}(x)$ или другой функции $f(x)$.

126 § 3. Результаты

127 На выборках объёма в миллион последовательностей для многих фиксированных
128 длин в интервале 50-100-200 найдены эмпирические зависимости медианы $Me(n)$ и
129 среднеквадратичного отклонения $\sigma(n)$ для функции $\text{impr}(x)$ от длины $n = |x|$ после-
130 довательностей x . Типичный вид этих зависимостей показан на рис. 1–2. Мы пред-
131 полагаем, что $Me(n)$ при фиксированном GC-составе стремится к пилообразному
132 графику с некоторой периодичностью и убывающей амплитудой; строгое описание
133 предельного вида графика представляло бы значительный интерес.

134 При $gc = 0,5$ среднеквадратичное отклонение $\sigma(n)$ функции $\text{impr}(x)$ хорошо ап-
135 проксимируется (с коэффициентом детерминации $R^2 = 0,999$) быстро убывающей
136 функцией от длины $n = |x|$ последовательностей x :

$$\sigma_{50\%} = \frac{(n - 0,6)^{0,35}}{2,58n}$$

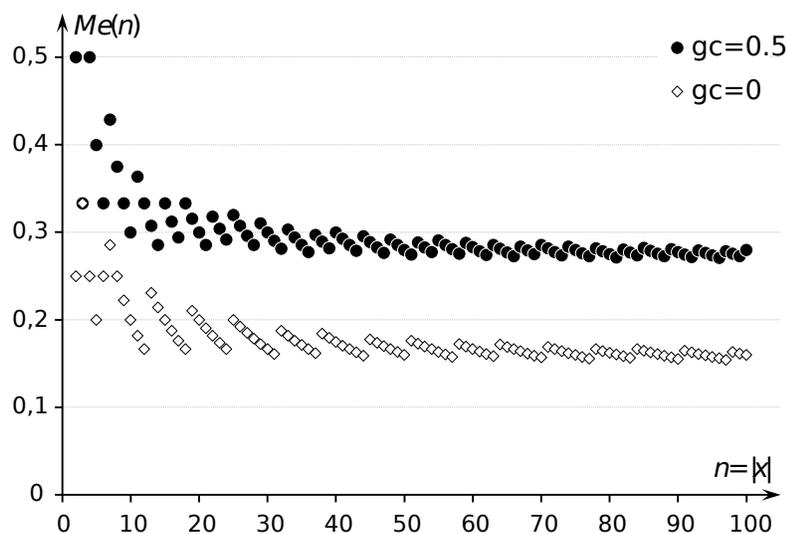


Рис. 1. Зависимость медианы $Me(n)$ (ордината) от данной длины $n = |x|$ (абсцисса) последовательностей x для функции $\text{imp}(x)$ при значениях параметра $gc = 0,5$ и $gc = 0$.

137 При $gc = 0$ среднеквадратичное отклонение $\sigma(n)$ аппроксимируется ($R^2 = 0,99$)
 138 другой быстро убывающей функцией от длины $n = |x|$ последовательностей x :

$$\sigma_{0\%} = \frac{(n - 0,1)^{0,36}}{2,65n}$$

139 Обе зависимости показаны на рис. 2а.

140 Далее, при $gc = 0,5$ на выборках объёма миллион последовательностей с теми
 141 же длинами найдена эмпирическая зависимость доли последовательностей x в за-
 142 висимости от редакционного расстояния dist между x и оптимальным палиндромом
 143 $wc(w)$. На рис. 3 показана зависимость числа N последовательностей x в зависимо-
 144 сти от редакционного расстояния dist при $n = 100$. В частности, доля последова-
 145 тельностей, близких к таким совершенным палиндромам, очень мала: не нашлось
 146 ни одной последовательности на расстояниях меньше 16 и больше 37.

147 Колоколообразный вид этой зависимости сохраняется и для других длин и вы-
 148 борок.

149 Далее, вычислены эмпирические среднее, медиана и другие квантили для функ-
 150 ции $\text{imp}(x)$ и случайных нуклеотидных последовательностей в зависимости от дли-
 151 ны и параметра gc . Например, для независимых случайных последовательностей x
 152 с длиной $n = |x| = 1000$, независимыми позициями и данным gc среднее $\mu(n, gc)$ и
 153 среднеквадратичное отклонение $\sigma(n, gc)$ для $\text{imp}(x)$, приведены в табл.1. Для длин
 154 свыше 1000 среднее μ почти не зависит от длины n последовательности, рис. 4.

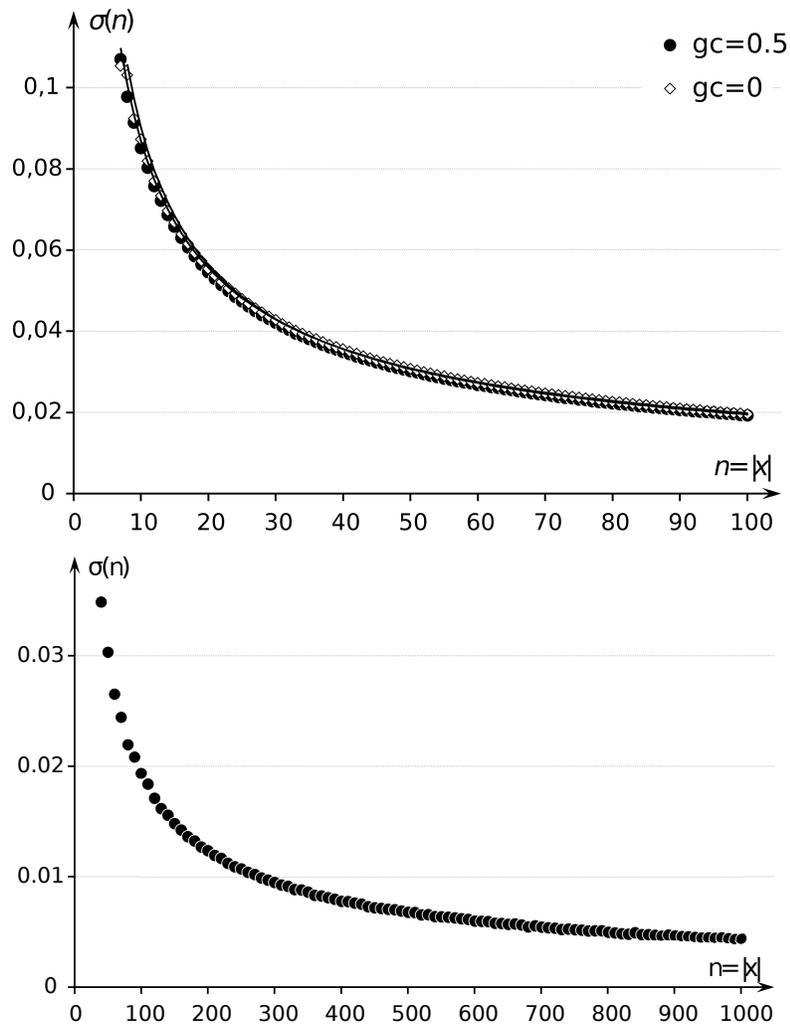


Рис. 2. (а) Зависимость среднеквадратичного отклонения $\sigma(n)$ (ордината) от длины $n = |x|$ (абсцисса) последовательностей x для функции $\text{imp}(x)$ при $gc = 0,5$ и $gc = 0$. На рис. 2а показаны ещё соответствующие регрессии, их формулы приведены ниже. (б) Первая из этих зависимостей для большего интервала длин последовательностей.

155 Ниже “ $k\%$ -квантиль равен q ” означает, что k процентов выборки последователь-
 156 ностей x удовлетворяет неравенству $\text{imp}(x) \leq q$ с наименьшим возможным q ; а 50%-
 157 квантиль совпадает с медианой. Рассмотрены выборки случайных последователь-
 158 ностей x с длинами $n = |x| = 100$ и $n = |x| = 50$ при различных GC-составах; анало-
 159 гичные соотношения наблюдаются и при $n = |x| = 1000$. В частности, из последних
 160 столбцов табл. 2-3 видно, что при любом GC-составе для почти всех последователь-
 161 ностей x значения $\text{imp}(x)$ близки к медиане.

162 Далее, наблюдается следующее. При $gc = 0,5$ и всех длинах $6 < n \leq 100$ последо-
 163 вательностей x медиана числовой последовательности, составленной из количеств

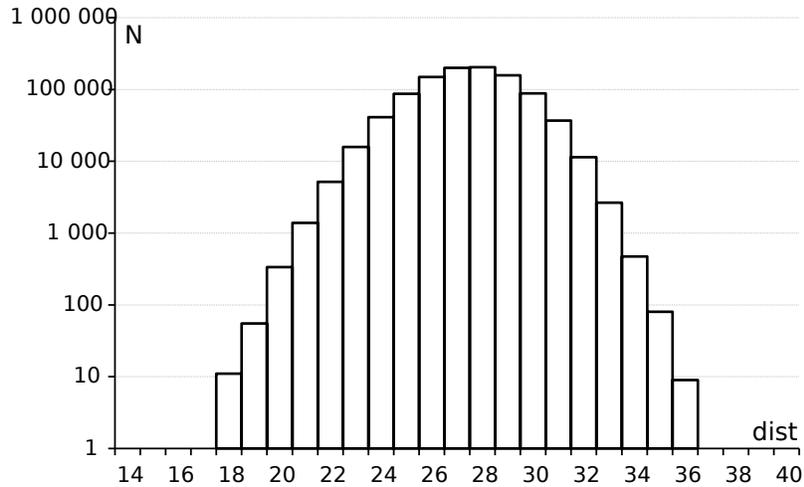


Рис. 3. Число N указано в логарифмическом масштабе. При $gc = 0,5$ показана зависимость числа N (ордината) последовательностей x (с длиной $n = 100$) в выборке объёма миллион от минимального редакционного расстояния dist для функции $\text{imp}(x)$ (абсцисса) между x и оптимальным палиндромом $ws(w)$.

Таблица 1. Среднее значение $\mu(n, gc)$ и среднеквадратичное отклонение $\sigma(n, gc)$ для функции $\text{imp}(x)$ и длинных последовательностей с $n = |x| = 1000$ в зависимости от GC-состава.

Параметр gc	Среднее $\mu(1000, gc)$ для $\text{imp}(x)$	Отклонение $\sigma(1000, gc)$ для $\text{imp}(x)$
0,00	0,147	0,0044
0,05	0,172	0,0054
0,10	0,193	0,0056
0,15	0,211	0,0056
0,20	0,225	0,0054
0,25	0,237	0,0052
0,30	0,246	0,0049
0,35	0,253	0,0048
0,40	0,258	0,0046
0,45	0,261	0,0044
0,50	0,262	0,0044

164 оптимальных палиндромов в последовательности x по выборкам с данным n , равна
 165 двум; т.е. примерно в половине по x имеется ≥ 2 оптимальных палиндромов. Мода
 166 этой последовательности равна 1. При любом $gc \in [0, 0,5]$ и $n = 100$ 90%-й кван-
 167 тиль количества оптимальных палиндромов равен четырём, то есть примерно 10%
 168 последовательностей x соответствует ≥ 4 по оптимальных палиндромов; а 99%-й
 169 квантиль равен 7. При этих gc среднее значение по выборкам этой числовой после-
 170 довательности лежит в интервале $(2, 14, 2, 28)$. Среднее число оптимальных палин-
 171 дромов при $n = |x| = 100$ и выборках объёма миллион слабо зависит от параметра
 172 gc , рис. 5.

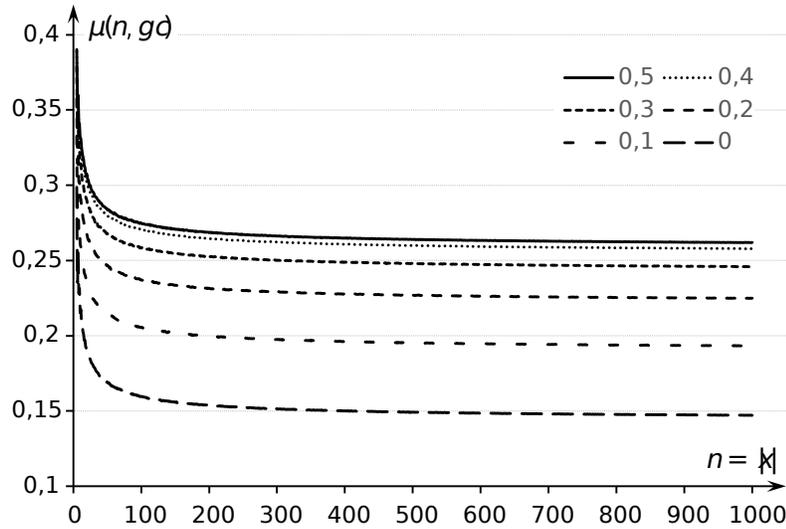


Рис. 4. Зависимости среднего значения $\mu(n, gc)$ для $\text{imp}(x)$ от длины $n = |x|$ последовательностей x при GC-составах 20%, 30%, 40% и 50%.

Таблица 2. Квантили для функции $\text{imp}(x)$ и среднее квадратичное отклонение $\sigma(n, gc)$ для коротких последовательностей x с длиной $n = |x| = 100$

Параметр gc	0,001%	0,01%	0,1%	1%	10%	50%	σ
0,00	0,08	0,09	0,10	0,12	0,14	0,16	0,020
0,05	0,09	0,11	0,12	0,14	0,16	0,18	0,022
0,10	0,11	0,12	0,14	0,15	0,18	0,21	0,022
0,15	0,13	0,14	0,15	0,17	0,19	0,22	0,022
0,20	0,14	0,15	0,17	0,19	0,21	0,24	0,022
0,25	0,15	0,17	0,18	0,20	0,22	0,25	0,021
0,30	0,17	0,18	0,19	0,21	0,23	0,26	0,021
0,35	0,17	0,18	0,20	0,22	0,24	0,27	0,020
0,40	0,18	0,19	0,21	0,22	0,25	0,27	0,020
0,45	0,19	0,19	0,21	0,23	0,25	0,27	0,019
0,50	0,18	0,20	0,21	0,23	0,25	0,28	0,019

173 Для всех последовательностей x небольшой длины вычислено редакционное рас-
174 стояние

$$\text{dist}(x) = |x| \text{imp}(x)$$

175 между x и оптимальным палиндромом. Например, в таблице 4 для всех последова-
176 тельностей x длины 15 в алфавите $\{A, C, G, T\}$ (миллиард с небольшим) приведено их
177 распределение в зависимости от значения $\text{dist}(x)$.

178 Найдены средние значения для функции $\text{imp}(x)$ по всем последовательностям
179 фиксированной длины n , которое назовём «точным» значением. В частности, для
180 длин $n = 15$ и $gc = 0,5$ и длин $n = 30$ и $gc = 0$ отличие этого точного среднего от

Таблица 3. Квантили и среднеквадратичное отклонение $\sigma(n, gc)$ для функции $\text{inp}(x)$ и ещё более коротких последовательностей x с длиной $n = |x| = 50$.

Параметр gc	0,001%	0,01%	0,1%	1%	10%	50%	σ
0,00	0,04	0,06	0,08	0,1	0,14	0,16	0,031
0,05	0,06	0,08	0,1	0,12	0,16	0,2	0,033
0,10	0,08	0,1	0,12	0,14	0,18	0,22	0,034
0,15	0,08	0,1	0,12	0,16	0,18	0,24	0,034
0,20	0,1	0,12	0,14	0,16	0,2	0,24	0,033
0,25	0,12	0,14	0,16	0,18	0,22	0,26	0,032
0,30	0,12	0,14	0,16	0,2	0,22	0,26	0,032
0,35	0,12	0,16	0,18	0,2	0,24	0,28	0,031
0,40	0,14	0,16	0,18	0,2	0,24	0,28	0,031
0,45	0,14	0,16	0,18	0,22	0,24	0,28	0,030
0,50	0,14	0,16	0,18	0,22	0,24	0,28	0,030

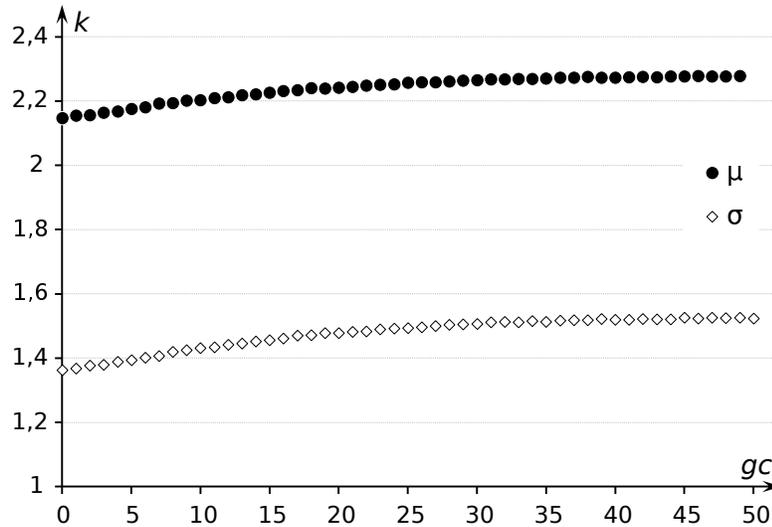


Рис. 5. Среднее число и стандартное отклонение числа оптимальных палиндромов при $n = |x| = 100$ и выборках объёма миллион в зависимости от параметра gc .

181 среднего по выборкам не превышает 0,2%. Аналогичная точная медиана совпадает
 182 с медианой по выборке, если медиана редакционного расстояния между x и опти-
 183 мальной палиндромом целое число. Например, это условие нарушается при $gc = 0$,
 184 то есть для последовательностей в алфавите $\{A, T\}$ и длинах $n = 2$ или $n = 6$, когда
 185 медианы расстояния равны 0,5 и 1,5 соответственно.

186 § 4. Заключение

187 Нами выполнены обширные компьютерные эксперименты, которые позволили
 188 предположить следующие утверждения. Доля случайных последовательностей x с
 189 длиной $n \geq 100$ и оценкой GC-состава параметром gc , которые близки к оптималь-

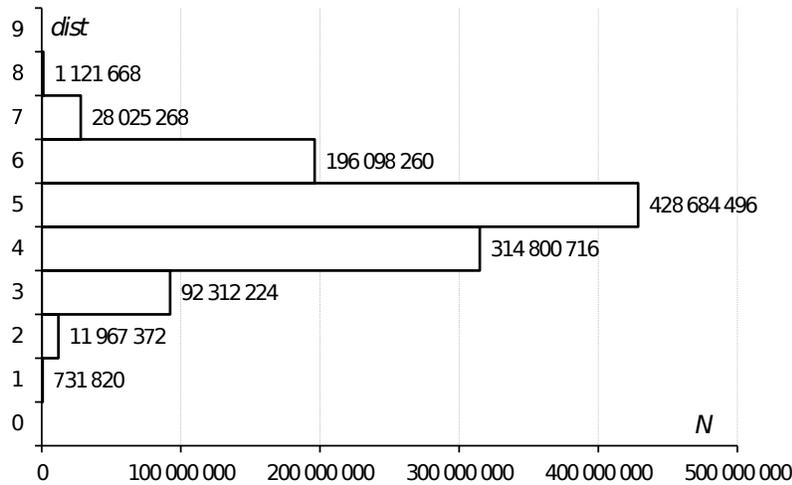


Рис. 6. Для всех последовательностей x длины 15 в алфавите $\{A, C, G, T\}$ приведено распределение их числа в зависимости от значения $\text{dist}(x)$.

190 ным палиндромам, очень мала. Распределение числа таких последовательностей x
 191 в зависимости от расстояния до оптимальных палиндромов имеет колоколообраз-
 192 ную форму, рис. 3. Для почти всех последовательностей x значение $\text{imp}(x)$ близко к
 193 среднему значению $\mu(n, gc)$ по всем x . Среднее значение $\mu(n)$ существенно зависит
 194 от параметра gc , достигая максимума при $gc = 0,5$, а симметрично расположенные
 195 минимумы этого среднего достигаются при $gc = 0$ и $gc = 1$ (с заменой в после-
 196 довательностях x из выборки букв А на G и Т на С). Квантили и среднее значение
 197 для функции $\text{imp}(x)$ зависят от длины последовательности x . Значения квантилей
 198 стабилизируются при длинах от 1000 букв. Для коротких последовательностей, дли-
 199 ны которых типичны для сайтов связывания транскрипционных факторов с ДНК,
 200 среднее значение $\mu(n, gc)$ оказывается гораздо выше, рис. 4, чем его предельное
 201 значение. Мы предполагаем, что это справедливо для любого GC-состава, хотя ско-
 202 рости сходимости к предельному значению различные. Квантили числовой после-
 203 довательности значений $\text{imp}(x)$ для случайных последовательностей x указывают
 204 на зависимость квантилей от длины x и её GC-состава, что позволяет отсеять по-
 205 следовательности, неотличимые от типичного случая. Иными словами, позволяет
 206 по экспериментальным длинам и GC-составу выбрать верхнюю границу $\text{imp}(x)$, при
 207 которой экспериментальные последовательности неслучайно близки к оптимальным
 208 палиндромам. В частности, при поиске несовершенных палиндромов в высококон-
 209 сервативных элементах генома [13], которые часто имеют небольшую длину, а их
 210 GC-состав значительно отличается от среднего значения по геному. Здесь важно
 211 учитывать длину и GC-состав индивидуальных последовательностей в геноме, и
 212 трудно подобрать универсальный порог для отсека типичных случаев. Тем не ме-
 213 нее, проведённые компьютерные эксперименты предлагают в биоинформатическом
 214 контексте грубую оценку верхней границы при поиске несовершенных палиндромов

215 x в форме неравенства $0 < \text{imp}(x) < Me(|x|)$. Многие авторы искали алгоритмы для
216 выделения короткого совершенного или почти совершенного палиндрома в длин-
217 ной последовательности x . Мы предполагаем, что алгоритм низкой полиномиальной
218 сложности для решения такой задачи возможен лишь при ограничениях на $\text{imp}(x)$,
219 которые превращают x практически в совершенный палиндром. Полученные выше
220 оценки также могут найти применение, например, в теории кодирования в связи
221 с q -ичными кодами, устойчивыми относительно появления вставок определённого
222 вида. Коды, включая ДНК-коды, устойчивые к появлению коротких tandemных по-
223 второв, рассмотрены в [14]. Аналогичными вставками являются инвертированные
224 повторы, которые образуют совершенный палиндром, когда инвертированный уча-
225 сток встраивается без зазора, или несовершенный палиндром, когда вставка отде-
226 лена некоторым словом.

227 СПИСОК ЛИТЕРАТУРЫ

- 228 1. T. Mieno, M. Funakoshi, Y. Nakashima, S. Inenaga, H. Bannai, and M. Takeda. Computing
229 maximal palindromes in non-standard matching models. *Information and Computation*,
230 304(105283):1–20, 2025.
- 231 2. R.R. Datta and J. Rister. The power of the (imperfect) palindrome: sequence-specific roles
232 of palindromic motifs in gene regulation. *Bioessays*, 44(4):1–22, 2022.
- 233 3. Н.Н. Назипова. Разнообразие некодирующих РНК в геномах эукариот. *Математиче-*
234 *ская биология и биоинформатика*, 16(2):256–298, 2021.
- 235 4. K.V. Mikhailov, B.D. Efeykin, A.Y. Panchin, D.A. Knorre, M.D. Logacheva, A.A. Penin,
236 M.S. Muntyan, M.A. Nikitin, O.V. Popova, O.N. Zanegina, M.Y. Vyssokikh, S.E. Spiridonov,
237 V.V. Aleoshin, and Y.V. Panchin. Coding palindromes in mitochondrial genes of
238 nematomorpha. *Nucleic Acids Research*, 47(13):6858–6870, 2019.
- 239 5. O.V. Nikolaeva, A.M. Beregova, B.D. Efeykin, T.S. Miroljubova, A.Yu. Zhuravlev, A.Yu.
240 Ivantsov, K.V. Mikhailov, S.E. Spiridonov, and V.V. Aleoshin. Expression of hairpin-
241 enriched mitochondrial dna in two hairworm species (nematomorpha). *International Journal*
242 *of Molecular Sciences*, 24(14):1–17, 2023.
- 243 6. Л.А. Мирошниченко, Н.А. Арефьева, Ю.П. Джиоев, В.Д. Гусев, А.Ю. Борисенко, С.В.
244 Эрдынеев, and Ю.С. Букин. Структура повторов в геномах сальмонелл. *Математи-*
245 *ческая биология и биоинформатика*, 18(2):602–620, 2023.
- 246 7. В.И. Левенштейн. Двоичные коды с исправлением выпадений, вставок и замеще-
247 ний символов. *Доклады АН СССР*, 163(4):845–848, 1965. Перевод: Levenshtein, V.I.
248 (1966). Binary codes capable of correcting deletions, insertions, and reversals. *Soviet Physics*
249 *Doklady*, 10(8), 707–710.
- 250 8. A. Tiskin. Semi-local longest common subsequences in subquadratic time. *Journal of*
251 *Discrete Algorithms*, 6(4):570–581, 2008.
- 252 9. S. Grabowski. New tabulation and sparse dynamic programming based techniques for
253 sequence similarity problems. *Discrete Applied Mathematics*, 212:96–103, 2016.

- 254 10. A. Tiskin. The chvátal–sankoff problem as a problem of symbolic dynamics. *Zapiski*
255 *Nauchnykh Seminarov POMI*, 528:214–237, 2023.
- 256 11. О. Зверков, А. Селиверстов, and Г. Шиловский. Выравнивание скрытого палиндрома.
257 *Математическая биология и биоинформатика*, 19(2):427–438, 2024.
- 258 12. M. Alzamel, C. Hampson, C.S. Iliopoulos, Z. Lim, S. Pissis, D. Vlachakis, and S. Watts.
259 Maximal degenerate palindromes with gaps and mismatches. *Theoretical Computer Science*,
260 978(114182):1–16, 2023.
- 261 13. L.I. Rubanov, A.V. Seliverstov, O.A. Zverkov, and V.A. Lyubetsky. A method for
262 identification of highly conserved elements and evolutionary analysis of superphylum
263 alveolata. *BMC Bioinformatics*, 17(385):1–16, 2016.
- 264 14. М. Ковачевич. О максимальном числе различных строк под действием коротких
265 тандемных дупликаций. *Проблемы передачи информации*, 58(2):12–23, 2022.

Хазиев Георгий Андреевич

Зверков Олег Анатольевич

Пирогов Сергей Анатольевич

Селиверстов Александр Владиславович

Любецкий Василий Александрович

Институт проблем передачи информации

им. А.А. Харкевича Российской академии наук, Москва

khaziev@iitp.ru

zverkov@iitp.ru

pirogov@iitp.ru

slvstv@iitp.ru

lyubetsk@iitp.ru

Поступила в редакцию

26.06.2023

После доработки

26.09.2023

Принята к публикации

20.10.2023