

Analysis of imperfect palindromes

Georgii A. Khaziev¹

¹PhD student at IITP RAS, Moscow, 1.5.8 Mathematical biology and bioinformatics,
physical-mathematical sciences

December 17, 2025

Today is a special day

I would like to formally congratulate my scientific supervisor Vassily Alexandrovich Lyubetsky with his 80-th birthday!!!



Palindromes in Bioinformatics.

Let us consider nucleotide sequences. The nucleotides $\{A, C, G, T\}$ form two complementary pairs:

$$c(A) = T, c(T) = A, c(C) = G, \text{ and } c(G) = C.$$

Next, $c()$ denotes the reverse complement, i. e., $c(xy) = c(y)c(x)$. For example, $c(AACG) = CGTT$. In DNA, a perfect palindrome is an inverted sequence repeat, i. e., reverse complement of itself. Let us omit the concatenation symbol. So, all perfect palindromes are of the type $xc(x)$, where x denotes a sequence. In particular, a sequence of odd length cannot be any perfect palindrome.

Needleman-Wunsch algorithm

let x and y be nucleotide sequences. Then to find edit distance between x and y we compute

$$f(j, k) = \begin{cases} j, & j \geq 0, k = 0 \\ k, & j = 0, k \geq 0 \\ \min\{f(j, k-1) + 1, f(j-1, k) + 1, f(j-1, k-1) + s(j, k)\}, & j > 0, k > 0 \end{cases}$$

where $s(j, k) \in \{0, 1\}$ depends on elements of x and y :

$$s(j, k) = \begin{cases} 0, & x_j = y_k \\ 1, & x_j \neq y_k \end{cases}$$

Needleman-Wunsch example

Let $x = \text{ACCAT}$ and $y = \text{AGGA}$, then we compute alignment matrix.

		A	C	C	A	T
	0	1	2	3	4	5
A	1	0	1	2	3	4
G	2	1	1	2	3	4
G	3	2	2	2	2	4
A	4	3	3	3	3	3

Edit distance between x and y is at bottom-right corner and equals 3.

Theorem 1

Theorem (Zverkov et al., 2024)

There exists an algorithm that for input sequences x and y computes optimal partition of y to such concatenation $y = wz$ that edit distance between x and a palindrome $wc(w)$ is minimal. In the same time edit distance between x and optimal palindrome $wc(w)$ is also computed. Overall algebraic complexity is $O(|x||y|)$.

Proof of Theorem 1

To compute edit distance between x and y it is enough to compute $f(|x|, |y|)$ of function f , recurrently defined as

$$f(j, k) = \begin{cases} j, & j \geq 0, k = 0 \\ k, & j = 0, k \geq 0 \\ \min\{f(j, k-1) + 1, f(j-1, k) + 1, f(j-1, k-1) + s(j, k)\}, & j > 0, k > 0 \end{cases}$$

where $s(j, k) \in \{0, 1\}$ depends on elements of x and y :

$$s(j, k) = \begin{cases} 0, & x_j = y_k \\ 1, & x_j \neq y_k \end{cases}$$

$f(j, k)$ allow us to compute edit distance between prefixes of x and y with lengths j and k .

Proof of Theorem 1

Now let's compute edit distance between x and $c(y)$. Let us define recursive function $g(j, k)$ as follows

$$g(j, k) = \begin{cases} j, & j \geq 0, k = 0 \\ k, & j = 0, k \geq 0 \\ \min\{g(j, k-1) + 1, g(j-1, k) + 1, \\ g(j-1, k-1) + r(j, k)\}, & j > 0, k > 0, \end{cases}$$

where $r(j, k) \in \{0, 1\}$ depends on elements of x and y :

$$r(j, k) = \begin{cases} 0, & x_{m-j+1} = c(y_k) \\ 1, & x_{m-j+1} \neq c(y_k) \end{cases}$$

Edit distance between x and $c(y)$ equals $g(|x|, |y|)$.

Proof of Theorem 1

Let h be defined as

$$h(j, k) = f(j, k) + g(|x| - j, k).$$

A value of $h(j, k)$ equals sum of two edit distances. First distance is between prefixes of x and y , and a second one is between corresponding suffixes of x and $c(y)$. yuEdit distance between x and $wc(w)$ equals minimal value of h for $0 \leq j \leq |x|$ and $0 \leq k \leq |y|$. At the same time corresponding value of k , defines the partition $y = wz$.

imp function

Let us denote by $|x|$ the length of x . Let us denote by $\text{imp}(x)$ the ratio of the minimum edit distance between sequence x and optimal palindrome to the length of the sequence:

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

Example

For $x = \text{ATATGT}$, $\text{imp}(x) = \frac{1}{6}$

The ratio shows how imperfect the palindrome is. $\text{imp}(x)$ could be easily computed using algorithm from Theorem 1¹.

¹If $y = x$, then $f(j, k)$ could be easily obtained as $f(j, k) = |j - k|$.

Theorem 2

Theorem

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

It is enough to show that for any $x = wz$ equalities

$$\text{dist}(x, wc(w)) = \text{dist}(x, c(z)z) = \text{dist}(c(x), c(z)z)$$

are satisfied. Indeed,

$$\text{dist}(x, wc(w)) = \text{dist}(z, c(w)),$$

since prefixes are matching and

$$\text{dist}(x, c(z)z) = \text{dist}(w, c(z)),$$

since suffixes are also matched. Also, since dist is symmetrical function it is true that

$$\text{dist}(z, c(w)) = \text{dist}(c(w), z).$$

Proof of Theorem 2

And because of invariance of dist under simultaneous replacement of both arguments to their reverse complements, it is true that

$$\text{dist}(c(w), z) = \text{dist}(w, c(z)).$$

Thus, first equality is satisfied

$$\text{dist}(x, wc(w)) = \text{dist}(z, c(w)) = \text{dist}(c(w), z) = \text{dist}(w, c(z)) = \text{dist}(x, c(z)z).$$

The second equality

$$\text{dist}(x, c(z)z) = \text{dist}(c(x), c(z)z)$$

is true for any z , since $c(z)z$ is a palindrome (i.e. $c(c(z)z) = c(z)z$).

Theorem 3

Theorem

For all even-length sequences it is true that $\text{imp}(x) \leq 1/2$. For all odd-length sequences it is true that $0 < \text{imp}(x) \leq (1 + 1/|x|)/2$.

For x of even length and prefix w which is equal half of x ,

$$\text{dist}(x, w) = \frac{|x|}{2}.$$

For odd-length sequence x and prefix w , which is limited by half of x , but don't contain the middle symbol,

$$\text{dist}(x, w) = \frac{|x| + 1}{2}.$$

In any of the cases it's true that

$$\text{dist}(x, wc(w)) \leq \text{dist}(x, w).$$

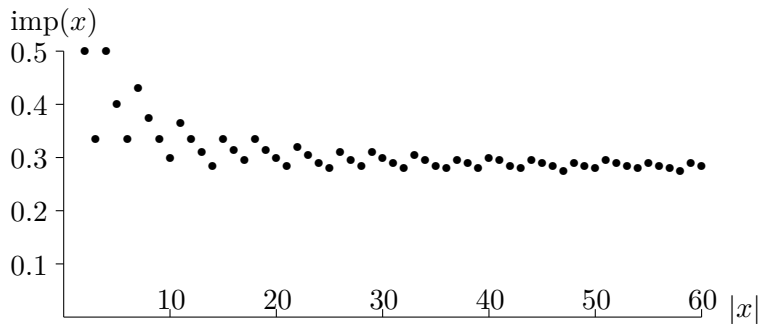


Figure: Empirical dependency between median of $\text{imp}(x)$ and $|x|$ for random nucleotide sequences for $gc = 0.5$.

median of imp

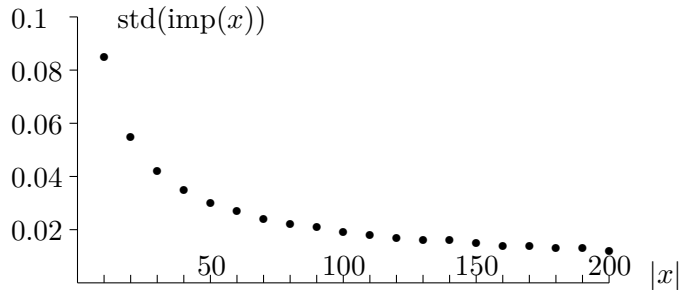


Figure: Empirical dependency between standard deviation of $\text{imp}(x)$ and $|x|$ for random nucleotide sequences.

trimmers

Substring is a contiguous sequence of characters within a string. The main idea behind the algorithm to select an imperfect palindrome is checking whether one of the optimal lengths of the prefix w of the input sequence x differs significantly from $|x|/2$. If such case occurs, then the algorithm deletes either prefix or suffix by difference between $|w|$ and $|x|/2$.

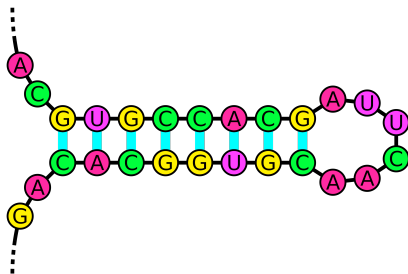


Figure: Hairpin example.

trimmers

All three functions `pref_trimmer`, `suff_trimmer`, and `double_trimmer` take as input nucleotide sequence x , sorted list `optimal_lengths` of optimal prefix lengths $|w|$, and floating-point number `cutoff_condition`. Note that `min(optimal_lengths)` and `max(optimal_lengths)` are the first and last elements of `optimal_lengths`, respectively.

The `pref_trimmer` function trims first rd symbols of x , where

$$rd = \max(\text{optimal_lengths}) - \left\lfloor \frac{|x|}{2} \right\rfloor,$$

when $rd \geq \text{length}(x) \cdot \text{cutoff_condition}$ is satisfied.

The `suff_trimmer` function trims last ld symbols of x , where

$$ld = \left\lfloor \frac{|x|}{2} \right\rfloor - \min(\text{optimal_lengths}),$$

when $ld \geq |x| \cdot \text{cutoff_condition}$ is satisfied.

double_trimmer

The `double_trimmer` function initially computes

$$rd = \max(\text{optimal_lengths}) - \left\lfloor \frac{|x|}{2} \right\rfloor$$

and

$$ld = \left\lfloor \frac{|x|}{2} \right\rfloor - \min(\text{optimal_lengths}).$$

Subsequently it checks if

$$\begin{cases} rd \geq |x| \cdot \text{cutoff_condition} \\ ld \geq |x| \cdot \text{cutoff_condition} \end{cases}$$

If the first inequality is satisfied, the function trims the first rd symbols from the string x . Similarly, if the second inequality is satisfied, the function trims the last ld symbols from x .

Theorem

For perfect palindrome x for any `cutoff_condition` > 0 no trimming would be performed.

Let x be a perfect palindrome. Then its only optimal partition would be exactly at

$$\left\lfloor \frac{|x|}{2} \right\rfloor$$

Then

$$ld = \left\lfloor \frac{|x|}{2} \right\rfloor - \left\lfloor \frac{|x|}{2} \right\rfloor = 0$$

and

$$rd = \left\lfloor \frac{|x|}{2} \right\rfloor - \left\lfloor \frac{|x|}{2} \right\rfloor = 0.$$

For a non-zero `cutoff_condition` both rd and ld would be less than $|x| \cdot \text{cutoff_condition}$ so no cutoff condition would be met.

Theorem (Theorem 5)

For any $n \geq 3$, there exist $\text{cutoff_condition} > 0$ and x with length of at least n , which satisfies both prefix and suffix trimming conditions such that

$$\text{pref_trimmer}(x_{\text{suffix}}) \neq \text{suff_trimmer}(x_{\text{prefix}}),$$

where x_{suffix} and x_{prefix} are results of suffix and prefix trimming of x , respectively.

Let

$$x = \underbrace{\text{ATA} \dots \text{ATA}}_{\alpha},$$

where $\alpha = 2 \lfloor \frac{n}{2} \rfloor + 1$, i.e. minimal odd number, not smaller than n . For such x optimal partitions would be precisely at $\lfloor \frac{n}{2} \rfloor - 1$ and $\lfloor \frac{n}{2} \rfloor + 1$. Minimal index of optimal partition would correspond to the perfect palindrome

$$\underbrace{\text{AT} \dots \text{AT}}_{\alpha-1}$$

Proof of Theorem 5

and maximal index would correspond to perfect palindrome

$$\underbrace{AT\dots AT}_{\alpha+1}.$$

Let $\text{cutoff_condition} = \frac{1}{10\alpha}$. Since this value satisfies both suffix and prefix trimming conditions, x_{suffix} and x_{pref} could be computed and would be equal to

$$\underbrace{AT\dots AT}_{\alpha-1}$$

and

$$\underbrace{T\dots ATA}_{\alpha-1}$$

respectively. Both of those sequences are perfect palindromes and would not be trimmed due to the previous theorem.

Another trimmers

We've also developed trimmer algorithms based on another approach. `pref_GRT`, `suff_GRT`, and `double_GRT` recursively trim $\text{cutoff_value} \cdot |x|$ symbols from x for a given `recursion_depth` number of iterations. On each iteration they check, whether value of `imp` decreased, if not, they decrease `cutoff_value` by half and return substring x' with minimal value of `imp` that was obtained.

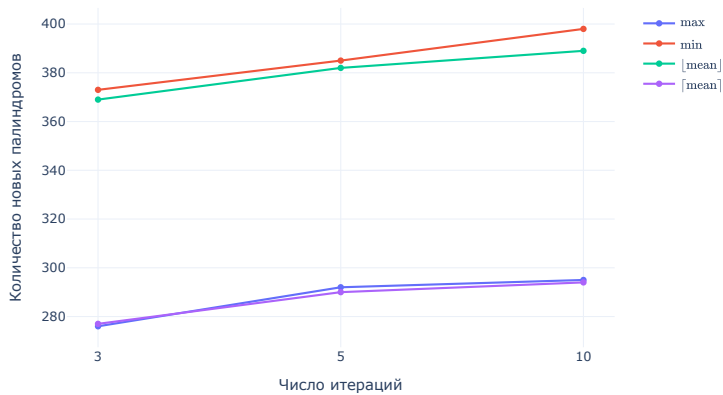
de_shapker

To find a noncomplementary substring(loop) inside x we've developed new algorithm called `de_shapker`. It takes as an input a string x , a floating point value `imp_given`, which equals $\text{imp}(x)$, matrix H of values of function $h(j, k)$ from `palindrome_self_alignment(x)` algorithm, integers `iteration_counter`, `window_delta` and `window_size`, and a Python function `strategy`. Let $x_{result} = x$. On each iteration `de_shapker` algorithm computes l_1 norms of all continuous submatrices $S \in \mathbb{R}^{\text{window_size} \times \text{window_size}}$ of H defined as

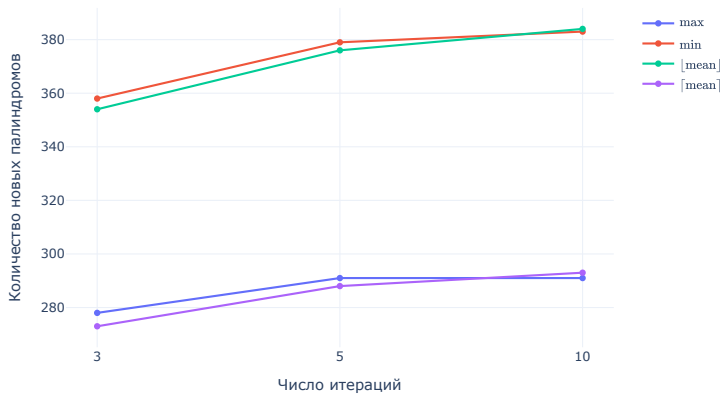
$$l_1(S) = \sum_{i=1}^{\text{window_size}} \sum_{j=1}^{\text{window_size}} (S)_{ij} \quad (1)$$

After that, the submatrix S_{\min} with minimal norm is chosen. Then, the index i of the beginning of potential loop in x is computed as $i = \text{strategy}(u, v)$, where u and v are coordinates of top-left element of S_{\min} in H . We recommend to choose as a strategy a function biased towards minimal value between u and v (e.g. $\min(u, v)$ or $\lfloor \text{mean}(u, v) \rfloor$) as they tend to capture loops more precisely. Subsequently, x_b is computed as x with cutted symbols from i to $i + \text{window_size}$ positions. If $\text{imp}(x_b) < \text{imp_given}$, then x_{result} will be overwritten with x_b , and imp_given will be overwritten with $\text{imp}(x_b)$. Next, the window_size is increased by window_delta and the next iteration starts. If $\text{imp}(x_b) \geq \text{imp_given}$, then algorithm returns x_{result} . After all iterations, algorithm returns x_{result} .

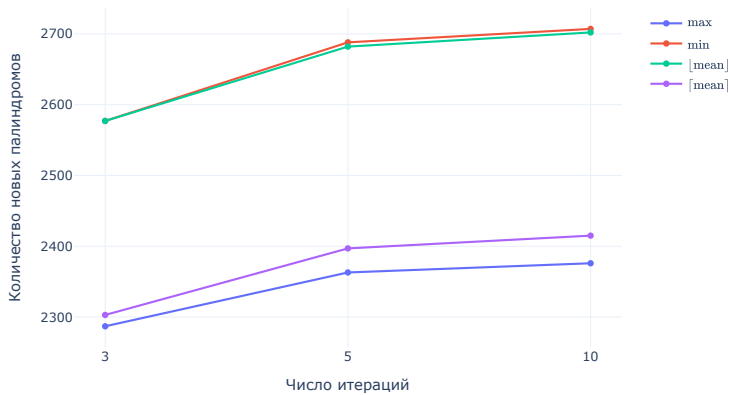
On real data



On real data

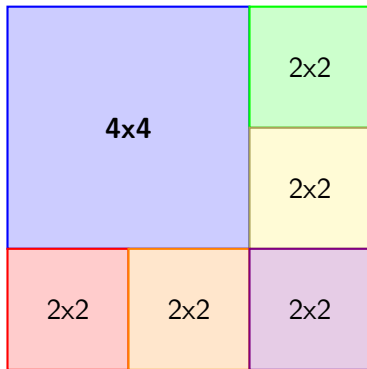


On real data



Boosting de_shapker

Note that on each iteration we recompute $l_1(S)$ increasing size of S each time by `window_delta`. But if we will use `window_size = window_delta` for the first iteration, then new l_1 norm could be computed from previous with addition of l_1 norms of the first iteration.



Boosting de_shapker

To implement this we need to store 3 square matrices of length $n - \text{window_size} + 1$, where $|x| = n$ and by `window_size` we mean input value of `window_size`. The first matrix would be used to store initial norms, second and third matrices would be previous and next step norms by rotation. After computing new norms matrix that stores previous ones would be multiplied by infinity.

Wobbly palindromes

In RNA hairpin is sometimes formed with the use of non-complementary pair GU². Such bond is called wobbly. To address this me and my student S. Mankevich implemented wobbly version of algorithm from theorem 1 with changed value of function

$$r(j, k) = \begin{cases} 0, & x_{m-j+1} = c(y_k) \vee (c(y_k), x_{m-j+1}) \in \{(C, U), (U, C), (A, G), (G, A)\} \\ 1, & otherwise \end{cases}$$

Such modification allows us to compute more weak distance and imp_w functions, but with accounting of wobbly pair.

We've tested this algorithm on mirBase dataset and from 38.5K RNA hairpins for less than 30 we couldn't find substrings with value $\text{imp}_w \leq 0.2$ after trimming. While on the same dataset for $\text{imp} \leq 0.2$ couldn't be found for about 8K of sequences.

²in RNA T is replaced by U

Conclusion

1. New algorithms for picking a close-to-palindrome subsequence of x are introduced.
2. Modification of algorithm for finding distance to perfect palindrome with respect to RNA wobbly pairs were made.
3. Those algorithms were tested on synthetic and real DNA and RNA data.

References



Zhao X.X., Wang Z.X., Tang D. et al. *The expectation and the variance of the weights of de Bruijn sequences*, Des. Codes Cryptogr. 2025.
<https://doi.org/10.1007/s10623-025-01729-2>



Nazipova N.N. *Application of suffix arrays to detect repeats in genomic sequences*, Mathematical Biology and Bioinformatics, 2025. V. 20 I. 2. P. 348-362.
<https://doi.org/10.17537/2025.20.348>



Zverkov O., Seliverstov A., Shilovsky G. *Alignment of a Hidden Palindrome*, Mathematical biology and bioinformatics. 2024. V.19. I.2. P. 427-438.
<https://doi.org/10.17537/2024.19.427>



Khaziev G.A., Seliverstov A.V., Zverkov O.A. *Searching for an imperfect palindrome*, Computer algebra: 6th International Conference Materials, Moscow, Russia, June 23–25 2025, RUDN University, 2025, P. 62–65.