# On repetition-free binary words of minimal density

Roman Kolpakov[*]        Gregory Kucherov[†]        Yuri Tarannikov[‡]

## Abstract

We study the minimal proportion (density) of one letter in $n$-th power-free binary words. First, we introduce and analyse a general notion of minimal letter density for any infinite set of words which don't contain a specified set of "prohibited" subwords. We then prove that for $n$-th power-free binary words the density function is $\frac{1}{n} + \frac{1}{n^3} + \frac{1}{n^4} + \mathcal{O}(\frac{1}{n^5})$. We also consider a generalization of $n$-th power-free words for fractional powers (exponents): a word is $x$-th power-free for a real $x$, if it does not contain subwords of exponent $x$ or more. We study the minimal proportion of one letter in $x$-th power-free binary words as a function of $x$ and prove, in particular, that this function is discontinuous at $\frac{7}{3}$ as well as at all integer points $n \geq 3$. Finally, we give an estimate of the size of the jumps.

**Keywords:** Unavoidable patterns, power-free words, exponent, minimal density.

## 1 Introduction

One of classical topics of formal language theory and word combinatorics is the construction of infinite words verifying certain restrictions. A typical restriction is the requirement that the word does not contain a subword of the form specified by some general pattern. Results of this kind find their applications in different areas such as algebra, number theory, game theory (see [15, 21]).

The oldest results of this kind, dating back to the beginning of the century, are Thue's famous constructions of infinite square-free and (strongly) cube-free words over alphabets of three and two letters respectively [22, 23] (see also [4]). A word is *square-free* (respectively *cube-free*, *strongly cube-free*) if it does not contain a subword $uu$ (respectively $uuu$, $uua$), where $u$ is a non-empty word and $a$ is the first letter of $u$.

During the last two decades, different generalizations of Thue's results have been studied. A natural generalization is to consider, instead of squares or cubes, any $n$-th power, or, yet more generally, any *pattern* (a word over some alphabet of variables). Works [3, 1] introduce a general property of *avoidability* of a pattern and propose an algorithm to test it. A pattern is avoidable iff for some $k$, there is an infinite word over $k$ letters that does not contain a subword which is an instance of the pattern. If $k$ is fixed, the pattern is called $k$-avoidable.[1] In this terminology, Thue's

[1]The difference between avoidability and $k$-avoidability is important. While avoidability was shown to be decidable in [3, 1], decidability of $k$-avoidability is a long-standing open problem (see [9]).

results state that pattern $xx$ is 3-avoidable, and pattern $xxx$ and, more strongly, the pattern $xyxyx$ are 2-avoidable. We refer the reader to [6, 7] for a survey of the area of pattern avoidability.

Many results on avoidability establish some threshold values or some "borderline conditions". As an example, let us mention the result of Roth [19] showing that every pattern over two variables of length six is 2-avoidable. Six is the best possible value, as there are patterns of length five that are not 2-avoidable (e.g. $xxyxx$).

As another example, Dejean [11] strengthens the Thue construction of a square-free word by constructing an infinite word over three letters such that any two occurrences of a non-empty word $u$ are separated by at least $|u|/3$ letters, and she shows that this bound is optimal. There is another formulation of this result: There exists an infinite word over three letters that not only avoids repetitions (subwords $uu$), but does not admit subwords $uv$, where $v$ is a prefix of $u$ of length more than $3|u|/4$. Generalizations of this result for bigger alphabets have been obtained (see [4] for more references; see also [8] for a related result).

In this paper, that fits into this general research direction, we address the following general problem. Assume that each letter has some weight, and we try to minimize the total weight of a word of given length avoiding the pattern. For example, if one letter is much "heavier" than the others, this leads to the following problem: Assume that a pattern is $k$-avoidable but not $(k-1)$-avoidable, then what is the minimal proportion of the $k$-th letter in an infinite word avoiding the pattern?

In this paper we solve this problem for the case of binary alphabet ($k = 2$), and patterns $x^n$ ($n$-th power) for $n > 2$. Specifically, we show that the minimal proportion $\rho(n)$ of one letter in an $n$-th power-free binary word is $\frac{1}{n} + \frac{1}{n^3} + \frac{1}{n^4} + \mathcal{O}(\frac{1}{n^5})$. As for strongly cube-free words, this proportion is asymptotically $1/2$, i.e. it is not possible in this case to reduce the number of occurrences of one letter with respect to the other. Both these results can be expressed uniformly through the generalized minimal frequency function based on the notion of exponent (fractional power) of a word. In this way, we consider $x$-th power-free words (words without subwords of exponent $x$ or more) for any real $x > 2$ (cf [16, 17, 10]). In the second part, we study properties of the generalized minimal frequency function $\rho(x)$ and prove, in particular, that it is discontinuous to the right of $\frac{7}{3}$ as well as to the right of all integer points $n \geq 3$. We then estimate the values $\rho(n+0)$ – the right limit of $\rho(x)$ at integer points $n \geq 3$ – and prove that $\rho(n+0) = \frac{1}{n} - \frac{1}{n^2} + \frac{2}{n^3} - \frac{2}{n^4} + \mathcal{O}(\frac{1}{n^5})$.

To the best of our knowledge, minimal density has been first studied in a related paper [14]. However, some work has been done on counting limit densities of subwords in words defined by DOL-systems (cf e.g. [12]).

The paper is organized as follows. We introduce the notion of minimal density in Section 2. In particular, we prove that two natural definitions of this notion lead to the same quantity. Section 3 is devoted to the estimate of $\rho(n)$. In Section 4 we study the generalized minimal density function. We conclude in Section 5 with possible directions for future work.

As usual, $A^*$ denotes the free monoid over an alphabet $A$. $u \in A^*$ is a *subword* of $w \in A^*$ if $w$ can be written as $u_1 u u_2$ for some $u_1 u_2 \in A^*$. $|u|$ stands for the length of $u \in A^*$. $A^\omega$ stands for the set of *one-way infinite* words, often called $\omega$-words, over $A$, that are defined as mappings $\mathbb{N} \to A$. For $n \in \mathbb{N}$, the word $w$ obtained by concatenating $n$ copies of a word $v$ is called the *$n$-th power* of $v$ and denoted $v^n$. A word $v$ is a *period* of $w$ iff $w$ is a subword of $v^n$ for some $n \in \mathbb{N}$.

# 2   Minimal density: general definition and properties

In this section we analyse, in a general context, the notion of minimal limit density of a letter in words of an infinite set.

Assume we have specified a set of words $P \subseteq A^*$, and consider the set $F \subseteq A^*$ of words that don't contain any word of $P$ as subword. For example, $P$ can be the set of instances of a given pattern, and $F$ the set of words avoiding this pattern (cf. Introduction). Note for any set $P$ of prohibited subwords, the set $F$ of avoiding words is *closed under subwords*, that is if a word $w$ is in $F$, then any subword of $w$ belongs to $F$ too. Moreover, any set $F$ closed under subwords is the set of avoiding words for some $P$ (just take $P = A^* \setminus F$). Therefore, being closed under subwords can be considered as a characterization for the sets of words that can be specified by means of prohibited subwords.

Assume we have an infinite set $F \subseteq A^*$ which is closed under subwords. Then there exist an infinite word from $A^\omega$ such that its every finite subword belongs to $F$. With interpretation of pattern avoidance, this allows to speak about infinite words avoiding the set of patterns. We denote by $F^\omega \subseteq A^\omega$ the set of all infinite words with every finite subword belonging to $F$.

Let $a \in A$ be a distinguished letter, and we are interested in the minimal limit proportion of $a$'s in words of $F$ of unbounded length. For $w \in F$, define $c_a(w)$ to be the number of occurrences of $a$ in $w$ and $\rho_a(w) = \frac{c(w)}{|w|}$. Denote $F(l) = \{w \in F \mid |w| = l\}$.

**Definition 1** *For every $l \in \mathbb{N}$, let $\rho_a(F, l) = \frac{1}{l} \min_{w \in F(l)} c_a(w)$ and $\rho_a(F) = \underline{\lim}_{l \to \infty} \rho_a(F, l)$. $\rho_a(F)$ is called the* minimal (limit) density *of $a$ in $F$.*

Note that the type of argument of $\rho_a$ will always make it clear if the density of an individual word, or the minimal density is meant.

Obviously, all numbers $\rho_a(F, l)$ belong to $[0, 1]$ and therefore $\rho_a(F)$ belongs to $[0, 1]$ too. The following two Lemmas clarify the behaviour of the sequence $\{\rho_a(F, l)\}_{l=1}^\infty$ with respect to $\rho_a(F)$.

**Lemma 1** *For every $l \in \mathbb{N}$, $\rho_a(F, l) \leq \rho_a(F)$.*

**Proof:** Take any $l \in \mathbb{N}$ and assume that $\{\rho_a(F, l_i)\}_{i=1}^\infty$ is a subsequence converging to $\rho_a(F)$. Take some $l_i > l$. By definition of $\rho_a(F, l_i)$, there exists a word $w_i \in F(l_i)$ such that $c_a(w_i) = l_i \rho_a(F, l_i)$. Consider $\lfloor l_i/l \rfloor$ non-overlapping subwords of $w_i$ of length $l$. Since $F$ is closed under subwords, each of these subwords belongs to $F(l)$ and then contains at least $l\rho_a(F, l)$ $a$'s. Therefore, $w_i$ contains at least $\lfloor l_i/l \rfloor l \rho_a(F, l)$ $a$'s, that is $c_a(w_i) \geq \lfloor l_i/l \rfloor l \rho_a(F, l)$. We obtain that $\rho_a(F, l_i) \geq \lfloor l_i/l \rfloor l \rho_a(F, l)/l_i > ((l_i/l) - 1) l \rho_a(F, l)/l_i = (1 - (l/l_i)) \rho_a(F, l)$. By taking the limit for $i \to \infty$, we conclude that $\rho_a(F) = \lim_{i \to \infty} \rho_a(F, l_i) \geq \lim_{i \to \infty} (1 - (l/l_i)) \rho_a(F, l) = \rho_a(F, l)$. $\square$

**Lemma 2** $\rho_a(F) = \lim_{l \to \infty} \rho_a(F, l) = \sup_{l \geq 1} \rho_a(F, l)$.

**Proof:** By Lemma 1, $\rho_a(F, l) \leq \rho_a(F)$ for every $l$, and then $\overline{\lim}_{i \to \infty} \rho_a(F, l) \leq \rho_a(F) = \underline{\lim}_{i \to \infty} \rho_a(F, l)$. Thus, $\rho_a(F, l)$ converge to $\rho_a(F)$ from below, and then $\rho_a(F) = \sup_{l \geq 1} \rho_a(F, l)$. $\square$

By Lemma 2, the lower limit in Definition 1 can be replaced by the simple limit. Thus, the definition $\rho_a(F) = \lim_{l \to \infty} \min_{w \in F(l)} \frac{c_a(w)}{l}$ is correct and seems to capture in a right way the notion of minimal density. However, there is another natural way to define the minimal limit density directly in terms of infinite words $F^\omega$, and not as the limit density value for finite words. One may ask if this approach always leads to the same density value or may lead to a different one.

For a word $w \in F \cup F^\omega$, let $w[1 : j]$ denotes the prefix of $w$ of length $j$. The density of letter $a$ in an infinite word $v \in F^\omega$ is naturally defined as the limit $\lim_{j \to \infty} \rho_a(v[1 : j])$. Obviously, this limit may not exist. However, below we show that among all words for which this limit exists, there is

one that realizes the minimum of these limits, which is equal to $\rho_a(F)$. This confirms that $\rho_a(F)$ is the right quantity caracterizing the limit density.

We define an auxiliary measure $\sigma_a(F,l) = \min_{w \in F(l)} \max_{1 \le j \le l} \rho_a(w[1:j])$. The following lemma gives a key argument.

**Lemma 3** *For every $l \in \mathbb{N}$, $\rho_a(F,l) \le \sigma_a(F,l) \le \rho_a(F)$.*

**Proof:** It is easily seen that $\sigma_a(F,l) \ge \rho_a(F,l)$. Let us prove that $\sigma_a(F,l) \le \rho_a(F)$ for all $l \in \mathbb{N}$. Assume that $\sigma_a(F,L) > \rho_a(F)$ for some $L \in \mathbb{N}$. This means that every word $v \in F$ of length at least $L$ has a pefix $v[1:j]$ with $\rho_a(v[1:j]) > \rho_a(F)$. Let $\varepsilon = \min\{\rho_a(v[1:j]) - \rho_a(F)\}$ where minimum is taken over all such prefixes. Take any word $w \in F(N)$ with $N > \frac{2L}{\varepsilon}(\rho_a(F)+\varepsilon)$. Find a decomposition $w = w_1 w_2 \ldots w_m$ such that $|w_j| \le L$ and $\rho_a(w_j) > \rho_a(F)$ for every $j$, $1 \le j \le m-1$, and $|w_m| < L$. Then $c_a(w) \ge (\rho_a(F)+\varepsilon)(|w|-L)$ and $\rho_a(w) \ge \rho_a(F)+\varepsilon - L(\frac{\rho_a(F)+\varepsilon}{|w|}) \ge \rho_a(F)+\frac{\varepsilon}{2}$. Since $w$ was chosen arbitrarily, this contradicts to $\rho_a(F,N) \le \rho_a(F)$ (Lemma 1). $\square$

**Corollary 1** *The limit $\lim_{l \to \infty} \sigma_a(n,l)$ exists and is equal to $\rho_a(F)$.*

**Lemma 4** *There exists a word $v \in F^\omega$ such that $\lim_{j \to \infty} \rho_a(v[1:j])$ exists and is equal to $\rho_a(F)$.*

**Proof:** From Lemma 3 it follows that for every $l \in \mathbb{N}$, there exists a word $w \in F(l)$ with $\rho(w) = \sigma_a(n,l) \le \rho_a(F)$, that is $\max_{1 \le j \le |w|} \rho_a(w[1:j]) \le \rho_a(F)$. Moreover, every prefix of $w$ verifies the same inequality. Therefore, the set of words $w$ verifying the inequality forms an infinite tree with respect to the prefix relation such that the parent of a word $w$ in the tree is its immediate prefix, obtained by removing the rightmost letter. Since the alphabet $A$ is finite, the tree is finitely branching. By König's Lemma, there exists an infinite path in this tree which defines the infinite word $v$ with $\rho_a(v[1:j]) \le \rho_a(F)$ for all $j \in \mathbb{N}$. Since $\rho_a(F,j) \le \rho(v[1:j]) \le \rho_a(F)$, the result follows from Lemma 2. $\square$

**Lemma 5** $\min_{v \in F^\omega} \lim_{j \to \infty} \rho_a(v[1:j]) = \rho_a(F)$, *where minimum is taken over $v \in F^\omega$ for which the limit exists.*

**Proof:** By Lemma 4, there exists a word $v \in F^\omega$ such that $\lim_{j \to \infty} \rho_a(v[1:j]) = \rho_a(F)$. Therefore, $\inf_{v \in F^\omega} \lim_{j \to \infty} \rho_a(v[1:j]) \le \rho_a(F)$. On the other hand, since $v[1:j] \in F(j)$, then $\rho_a(v[1:j]) \ge \rho_a(F,j)$, then $\lim_{j \to \infty} \rho_a(v[1:j]) \ge \lim_{j \to \infty} \rho_a(F,j) = \rho_a(F)$ and $\inf_{v \in F^\omega} \lim_{j \to \infty} \rho_a(v[1:j]) \ge \rho_a(F)$. The lemma follows. $\square$

Lemmas 4 and 5 imply that there exists a word $v \in F^\omega$ that realizes the minimal limit $\lim_{j \to \infty} \rho_a(v[1:j])$ among all words of $F^\omega$ for which the limit exists. Moreover, this minimum is equal $\rho_a(F)$. To avoid the problem of existence of the limit, we could replace it by the lower limit and define the quanity $\inf_{v \in F^\omega} \underline{\lim}_{j \to \infty} \rho_a(v[1:j])$ where the infimum is taken over *all* words $v \in F^\omega$. The proof of Lemma 5 shows that this value is also equal to $\rho_a(F)$, and the infimum is reached on some word $v \in F^\omega$.

Finally, note that one might suggest yet another, though less natural definition of minimal letter density as the value $\underline{\lim}_{j \to \infty} \min_{v \in F^\omega} \rho_a(v[1:j])$. Using Lemma 4 and arguments similar to the proof of Lemma 5, it is easily shown that the lower limit here can be replaced by the simple limit which is again equal to $\rho_a(F)$.

The equvalence of different definitions gives a strong evidence that $\rho_a(F)$ is an interesting quantity to study. In this paper, we undertake this study for a particular family of sets $F$ – the sets of $n$-th power-free binary words.

# 3  Minimal letter density in $n$-th power-free binary words

Consider an alphabet $A$. For a natural $n \geq 2$, a word $w \in A^*$ is called *n-th power-free* iff it does not contain a subword which is the $n$-th power of some non-empty word. We denote $PF(n) \subseteq A^*$ the set of $n$-th power-free finite words. Words from $PF(2)$ are called *square-free*, and words from $PF(3)$ are called *cube-free*. If $w \in A^*$ does not contain a subword $uua$, where $u$ is a non-empty word and $a$ is the first letter of $u$, then $w$ is called *strongly cube-free*. An equivalent property (see [20]) is overlap-freeness – $w$ is *overlap-free* if it does not contain two overlapping occurrences of a non-empty word $u$. Well known Thue's results [22, 23] state that there exist square-free words of unbounded length on the 3-letter alphabet, and strongly cube-free words of unbounded length on the 2-letter alphabet. An infinite sequence of strongly cube-free words can be constructed by iterating the morphism $h(0) = 01$, $h(1) = 10$, known as Thue-Morse morphism. Note that the existence of infinite strongly cube-free words on the 2-letter alphabet implies that for that alphabet the set $PF(n)$ is infinite for every $n \geq 3$.

From now on we fix on the binary alphabet $A = \{0, 1\}$. Our goal is to compute, for all $n > 2$, the value $\rho_1(PF(n))$ – minimal density of 1 in the words $PF(n)$. Note that by symmetry, $\rho_1(PF(n)) = \rho_0(PF(n))$, and to simplify the notation, we denote $\rho_1(PF(n))$ (respectively $\rho_1(PF(n), l)$) by $\rho(n)$ (respectively $\rho(n, l)$) in the sequel. Similarly, we will drop the index in $c_1(w)$ and $\rho_1(w)$, and will write $c(w)$ and $\rho(w)$ instead.

In [14] it has been proved that $\rho(n) = \frac{1}{n} + \mathcal{O}(\frac{1}{n^2})$. Here, using a different method, we prove the following more precise estimation, that corresponds to the first four terms in the asymptotic expansion of $\rho(n)$.

**Theorem 1** $\rho(n) = \frac{1}{n} + \frac{1}{n^3} + \frac{1}{n^4} + \mathcal{O}(\frac{1}{n^5})$.

We first establish the upper bound

$$\rho(n) \leq \frac{1}{n} + \frac{1}{n^3} + \frac{1}{n^4} + \frac{C}{n^5}, \tag{1}$$

for all $n \geq 3$ and some positive constant $C$. The proof is based on the following lemma.

Denote by $\alpha_i$ the word $0^i 1$.

**Lemma 6** Let $k \geq 3$. For $i, j$, $0 \leq i, j \leq k$ and $i \neq j$, consider a morphism $h : \{0, 1\}^* \rightarrow \{0, 1\}^*$ defined by $h(0) = \alpha_i$, $h(1) = \alpha_j$. For a word $w \in \{0, 1\}^*$, if $w \in PF(k)$ then $h(w) \in PF(k + 1)$.

**Proof:** First observe that $\{h(0), h(1)\}$ is a prefix code, i.e. the inverse image $w$ of any word $h(w)$ is unique. Furthermore, for any $u \in \{0, 1\}^*$, the occurrences of 1 in $h(u)$ delimit the images of individual letters of $w$. This means that any subword of $h(w)$ which ends with 1 and is preceeded by 1 (or starts at the beginning of $h(w)$) is the image of some subword of $w$.

To prove the lemma, assume by contradiction that for some $w \in PF(k)$, $h(w)$ contains a subword $v^{k+1}$. Proceed by case analysis on the number of 1's in $v$. If $v$ contains no 1's, then $v^{k+1}$ contains at least $k+1$ consecutive 0's which is impossible as $h(w)$ is a concatenation of words $\alpha_i, \alpha_j$. If $v$ contains one 1, then $v = 0^l 1 0^m$, and $v^{k+1} = 0^l 1 (0^{l+m} 1)^k 0^m$. Since $h(w) \in \{\alpha_i, \alpha_j\}^*$, we conclude that $l + m \in \{i, j\}$ and $w$ must contain $k$ consecutive occurrences of the letter $h^{-1}(0^{l+m} 1)$. Finally, if $v$ contains $s$ 1's, then $v = 0^l 1 \alpha_{i_1} \dots \alpha_{i_{s-1}} 0^m$, and $v^{k+1} = 0^l 1 (\alpha_{i_1} \dots \alpha_{i_{s-1}} 0^{l+m} 1)^k 0^m$. Again, $l + m \in \{i, j\}$ and $w$ contains the $k$-th power of the inverse image $h^{-1}(\alpha_{i_1} \dots \alpha_{i_{s-1}} 0^{l+m} 1)$. $\square$

**Lemma 7** *For every $n \geq 4$,*

$$\rho(n) \leq \frac{1}{n - \rho(n-1)} \tag{2}$$

**Proof:** For $l \in \mathbb{N}$, take a word $w \in PF(n-1)$ with $|w| = l$ and $\rho(w) = \rho(n-1, l)$. Denote by $h$ the morphism defined by $h(0) = \alpha_{n-1}$, $h(1) = \alpha_{n-2}$. Let $u = h(w)$. By Lemma 6, $u \in PF(n)$. Since $c(u) = |w|$, and $|u| = (n-1)c(w) + n(|w| - c(w)) = n|w| - c(w)$, we have $\rho(n, |u|) \leq \rho(u) = \frac{c(u)}{|u|} = \frac{1}{n - \rho(w)} = \frac{1}{n - \rho(n-1,l)}$. Taking the limit for $l \to \infty$, and then $|u| \to \infty$, we have $\rho(n) \leq \frac{1}{n - \rho(n-1)}$. $\square$

Upper bound (1) is now proved by simple induction on $n \geq 3$. Using the trivial inequality $\rho(3) \leq 1/2$, the base case $n = 3$ can be satisfied by choosing any constant $C \geq 57/2$. To prove the inductive step, we apply Lemma 7. This leads to the inequality

$$\frac{1}{n - \left( \frac{1}{n-1} + \frac{1}{(n-1)^3} + \frac{1}{(n-1)^4} + \frac{C}{(n-1)^5} \right)} \leq \frac{1}{n} + \frac{1}{n^3} + \frac{1}{n^4} + \frac{C}{n^5}$$

for $n \geq 4$, which reduces to the polynomial inequality

$$(-3 + C)n^6 + (-5C + 8)n^5 + (8C - 9)n^4 + (2 - 6C)n^3 + (-3C + 3)n^2 + (3C - 1)n - (C^2 + C) \geq 0$$

After substituting $C = 30$, the routine check shows that the inequality holds for all $n \geq 4$ (substitute $n - 4$ for $n$, expand and notice that all coefficients get positive). This proves that upper bound (1) holds for $C = 30$.

Note that constant $C$ can be reduced if we take into consideration the next term in the asymptotic expansion. Using a similar argument, it can be shown that

$$\rho(n) \leq \frac{1}{n} + \frac{1}{n^3} + \frac{1}{n^4} + \frac{3}{n^5} + \frac{C_1}{n^6}$$

for $C_1 = 90$.

Now we turn to bounding $\rho(n)$ from below, and prove the following lower bound.

$$\rho(n) \geq \frac{n-1}{n^2 - n - 1} \tag{3}$$

for all $n \geq 3$.

Consider an arbitrary finite $n$-th power-free word $w$. First, group its letters into blocks $\alpha_i = 0^i 1$, $0 \leq i \leq n - 1$. For a technical reason we assume that $w$ does not start with $\alpha_{n-1}$. If it does, we temporarily remove the first symbol 0. $w$ is uniquely decomposed into a concatenation of $\alpha_i$'s and a suffix of at most $n - 1$ 0's. Then, we group occurrences of $\alpha_i$'s into larger blocks $\beta(m, k) = (\alpha_{n-1})^m \alpha_k$, $0 \leq m \leq n - 1$, $0 \leq k \leq n - 2$. Informally, blocks $\beta$ are delimited by occurrences of $\alpha_i$ with $i \leq n - 2$. Again, $w$ is uniquely decomposed into blocks $\beta$ and the remaining suffix $Q$ of length at most $n^2 - 1$ ($n - 1$ occurrences of $\alpha_{n-1}$ followed by $n - 1$ 0's). We proceed by grouping blocks $\beta$ into yet more large blocks. Let

$$\gamma(l, k_0, k_1, \ldots, k_s) = \beta(l, k_0)\beta(n-1, k_1) \ldots \beta(n-1, k_s) = (\alpha_{n-1})^l \alpha_{k_0}(\alpha_{n-1})^{n-1}\alpha_{k_1} \ldots (\alpha_{n-1})^{n-1}\alpha_{k_s},$$

where $0 \leq l \leq n - 2$, $s \geq 0$, $0 \leq k_0, k_1, \ldots, k_s \leq n - 2$. Blocks $\gamma$ are delimited by each occurrence of $\beta(l, k)$ with $l \leq n - 2$. Note that since $w$ starts with $\alpha_k$, $k \leq n - 2$, it starts with $\beta(0, k)$ and

6

therefore the first block $\gamma$ starts at the beginning of $w$. Thus, the decomposition of $w$ is uniquely defined with a possibly remaining suffix $Q$ of length up to $n^2 - 1$. Taking into account the first possibly removed 0, we have $w = Pw'Q$, where $|P| \leq 1, |Q| \leq n^2 - 1$, and $w'$ is uniquely decomposed into blocks $\gamma$.

Let us now compute the minimal possible ratio of 1's in blocks $\gamma$. Consider a block $\gamma(l, k_0, k_1, \ldots, k_s)$. We distingish two cases:

**Case $s \geq 1$:** We show that $k_j + k_{j+1} \leq n - 2$ for every $j$, $0 \leq j \leq s - 1$. Indeed, consider the subword $\alpha_{k_j}(\alpha_{n-1})^{n-1}\alpha_{k_{j+1}}$ of $\gamma(l, k_0, k_1, \ldots, k_s)$. If $k_j + k_{j+1} \geq n - 1$ then it has the prefix $(0^{k_j}10^{n-1-k_j})^n$ which contradicts to the $n$-th power-freeness of $w$.

Using this observation, we can bound

$$\sum_{j=0}^{s} |\alpha_{k_j}| \leq \begin{cases} \frac{s+1}{2}n & s \text{ impair,} \\ \frac{s}{2}n + (n-1) & s \text{ pair.} \end{cases}$$

Then $|\gamma(l, k_0, k_1, \ldots, k_s)| \leq \frac{s}{2}n + (n-1) + sn(n-1) + ln = s(n^2 - \frac{n}{2}) + nl + n - 1$. Since the number of 1's in $\gamma(l, k_0, k_1, \ldots, k_s)$ is $ns + l + 1$, we have $\rho(\gamma(l, k_0, k_1, \ldots, k_s)) \geq \frac{ns+l+1}{s(n^2-\frac{n}{2})+nl+n-1}$. The right-hand side minimizes when $l$ is maximal ($l = n - 2$) and $s$ is minimal ($s = 1$). We then obtain $\rho(\gamma(l, k_0, k_1, \ldots, k_s)) \geq \frac{2n-1}{2n^2 - \frac{3n}{2} - 1}$.

**Case $s = 0$:** In this case $\gamma(l, k_0) = \beta(l, k_0)$, $|\gamma(l, k_0)| = ln + k_0 + 1$, and $\rho(\gamma(l, k_0)) = \frac{l+1}{ln+k_0+1}$. The right-hand mininizes when both $l$ and $k_0$ are maximal ($l = k_0 = n - 2$), which gives $\rho(\gamma(l, k_0)) \geq \frac{n-1}{n^2-n-1}$.

The second case gives a smaller bound for all $n \geq 3$ and we conclude that $\rho(\gamma(l, k_0, \ldots, k_s)) \geq \frac{n-1}{n^2-n-1}$. Since $w'$ is a concatenation of blocks $\gamma$, this implies $\rho(w') \geq \frac{n-1}{n^2-n-1}$. Returning to $w$, we have $c(w) \geq c(w') \geq \frac{n-1}{n^2-n-1}|w'| \geq \frac{n-1}{n^2-n-1}(|w| - n^2)$, and then $\rho(w) = \frac{c(w)}{|w|} \geq \frac{n-1}{n^2-n-1}(1 - \frac{n^2}{|w|})$. As $w$ is an arbitrary $n$-th power-free word, we have $\rho(n, l) \geq \frac{n-1}{n^2-n-1}(1 - \frac{n^2}{l})$ for all $l$. Taking the limit for $l$ going to infinity, we obtain $\rho(n) \geq \frac{n-1}{n^2-n-1}$. This implies in particular that

$$\rho(n) \geq \frac{1}{n} + \frac{1}{n^3} + \frac{1}{n^4} + \frac{2}{n^5} \tag{4}$$

Lower bound (4) together with upper bound (1) implies Theorem 1.

## 4 Generalized minimal density function

As a function on integer argument $n \geq 3$, $\rho(n)$ admits an interesting extension to real argument. The extension is achieved through the notion of *exponent* (see [11, 5, 8]) that generalizes the notion of $n$-th power. The exponent of a word $w$ is the ratio $\frac{|w|}{\min|v|}$, where the minimum is taken over all periods $v$ of $w$. Instead of $n$-th power-free words, we can now consider words which don't contain subwords of exponent $x$ or more, where $x$ is a real number (see, e.g., [16, 17, 10]).

Formally, for a real number $x$, define $PF(x)$ (resp. $PF(x + \varepsilon)$) to be the set of binary words that do not contain a subword of exponent greater than or equal to (resp. strictly greater than) $x$. Note that $PF(2 + \varepsilon)$ is precisely the class of strongly cube-free words. For the binary alphabet, the

existence of infinite cube-free words implies that $PF(x)$ (resp. $PF(x+\varepsilon)$) is infinite for $x > 2$ (resp. for $x \geq 2$). Using the results of Section 2, values $\rho_1(PF(x))$ and $\rho_1(PF(x+\varepsilon))$ are well-defined for $x > 2$ and $x \geq 2$ respectively. Similar to the previous section, we denote them respectively by $\rho(x)$ and $\rho(x+\varepsilon)$. Notation $\rho(x,l)$ and $\rho(x+\varepsilon,l)$ is defined accordingly. Note that for natural values of $x > 2$, $\rho(x)$ coincides with $\rho(n)$ studied in the previous section.

## 4.1 Discontinuity of $\rho(x)$

Now we study the generalized function $\rho(x)$ and prove, in particular, that it has an infinite number of discontinuity points.

Functions $\rho(x), \rho(x+\varepsilon)$ are non-increasing with values from $[0, \frac{1}{2}]$. Therefore, at every $x > 2$ there exists a right limit, denoted $\rho(x+0)$, and $\rho(x+0) = \sup_{y>x} \rho(y)$. The following lemma is useful.

**Lemma 8** *For every $x > 2$, $\rho(x+0) = \rho(x+\varepsilon)$.*

**Proof:** Clearly, for every $y > x$, $\rho(y) \leq \rho(x+\varepsilon)$, and therefore $\rho(x+0) = \sup_{y>x} \rho(y) \leq \rho(x+\varepsilon)$. Assume that $\rho(x+0) < \rho(x+\varepsilon)$. Then by Lemma 2, for some $l$, $\rho(x+\varepsilon,l) > \rho(x+0)$. The exponents of subwords of words of length $l$ can take finitely many possible values. Let $\hat{x}_l$ be the smallest such value strictly greater than $x$. Then $\rho(x+\varepsilon,l) = \rho(\hat{x}_l,l) > \rho(x+0)$. By Lemma 1, $\rho(\hat{x}_l) \geq \rho(\hat{x}_l,l)$ and then $\rho(\hat{x}_l) > \rho(x+0)$. This contradicts to the fact that $\rho(x+0) = \sup_{y>x} \rho(y)$. $\square$

Let us now compute the value $\rho(2+\varepsilon)$. The class $PF(2+\varepsilon)$ of *strongly cube-free* (overlap-free) binary words has been extensively studied (see [18, 13]) and the structure of these words has been thoroughly characterized. In particular, it is known that every strongly cube-free word can be written as $v_1 v v_2$, where $|v_1| \leq 2$, $|v_2| \leq 2$ and $v \in \{01, 10\}^*$ (see Lemma 2.2 of [13]). This implies immediately that $\rho(2+\varepsilon) = \frac{1}{2}$. However, a stronger result can be stated.

**Lemma 9** *For all $x \in (2, \frac{7}{3}]$, $\rho(x) = \frac{1}{2}$.*

**Proof:** To prove that a word $w$ can be decomposed as above, it is sufficient to assume that $w$ does not contain subwords $vva$ where $|v| \leq 3$ and $a$ is the first letter of $v$. We refer the reader to the proof of Lemma 2.1 in [13] to check this out. $\square$

We now show that $\rho(x)$ is discontinuous at $x = 7/3$. Specifically, we prove

**Theorem 2**

$$\rho\left(\frac{7}{3} + \varepsilon\right) \leq \frac{10}{21}.$$

Together with $\rho(\frac{7}{3}) = \frac{1}{2}$ (Lemma 9), this proves a jump of $\rho(x)$ to the right of $x = \frac{7}{3}$.

Consider the morphism $h$ defined by

$$h(0) = 011010011001001101001,$$
$$h(1) = 100101100100110010110.$$

We call $h(0), h(1)$ *coding words*. Note that a coding word is uniquely determined by its first (or last) letter. Consider a word $w \in \{0,1\}^*$ and its image $h(w)$. An occurrence of $h(a)$, $a \in \{0,1\}$, in $h(w)$, corresponding to the image of $a$ in $w$, is called a *coding occurrence*. Consider an occurrence

in $h(w)$ of some subword $u$, that is $h(w) = u_1 u u_2$ for some $u_1, u_2 \in \{0,1\}^*$. Consider the minimal subword $w'$ of $w$ such that $u$ is covered by $h(w')$, that is $h(w) = h(w_1)h(w')h(w_2)$, $h(w') = \delta_1 u \delta_2$, $h(w_1)\delta_1 = u_1$, $\delta_2 h(w_2) = u_2$, and $\delta_1$ (resp. $\delta_2$) is a proper prefix (resp. suffix) of a coding occurrence. We call $\delta_1$ the *precursor* of this occurrence of $u$.

We show that $h$ preserves the property of absence of subwords of exponent greater than $7/3$.

**Lemma 10** *For every $w \in \{0,1\}^*$, if $w$ does not contain subwords of exponent greater than $7/3$, then neither does $h(w)$.*

**Proof:** Assume that $w$ does not contain subwords of exponent greater than $7/3$. First show that $h(w)$ does not contain a subword of exponent greater than $7/3$ and with a period less than or equal to 15. If such a subword exists, there exists another one, say $v$, of length at most 36, with the same period and of exponent greater than $7/3$. Since $|h(0)| = |h(1)| = 21$, $v$ is covered by three contiguous coding occurrences. Since $w$ does not contain 000 or 111, $v$ occurs in one of $h(001)$, $h(010)$, $h(011)$, $h(100)$, $h(101)$, $h(110)$. A direct exhaustive check shows that none of these words contains a subword of exponent greater than $7/3$ and with a period at most 15.

Now assume that $h(w)$ contains a subword $v = v_1 v_1 v_2$, where $|v_1| \geq 16$, $v_2$ is a prefix of $v_1$, and $|v_2| > |v_1|/3$. Let $w'$ be the shortest subword of $w$ such that $h(w')$ contains $v$, that is $h(w') = \delta_1 v \delta_2$, where $\delta_1, \delta_2 \in \{0,1\}^*$, and $\delta_1$ is the precursor of the considered occurrence of $v$.

We now observe that if $u$ is a subword of $h(w)$ of length 16 or more, then the precursor of $u$ is uniquely defined for all occurrences of $u$. Since every subword of length 16 is located within two coding occurrences, this can be shown by checking this property for all subwords of length 16 occurring in words $h(00), h(01), h(10), h(11)$. [2]

By applying this argument to the two occurrences of $v_1$ in $v$ and by using properties of $h$, we can rewrite $h(w') = h(a)v_1' h(a)v_1' h(a)v_2' \delta_2$, $a \in \{0,1\}$. Now observe that
(1) $v_1'$ is non-empty (otherwise $w$ would contain $aaa$),
(2) $|h(a)v_1'| = |v_1|$, $|h(a)v_2'| \geq |v_2|$,
(3) $v_2' \delta_2$ is a prefix of $v_1' h(a)$,
(4) $v_1' = h(w_1)$ and $v_2' \delta_2 = h(w_2)$ for some $w_1, w_2 \in \{0,1\}$.
By taking the inverse image of $h(w')$, we get $w' = a w_1 a w_1 a w_2$, where $w_2$ is a prefix of $w_1 a$, and

$$|aw_2| = \frac{|h(a)v_2'| + |\delta_2|}{21} \geq \frac{|v_2|}{21} > \frac{1}{3} \cdot \frac{|v_1|}{21} = \frac{1}{3}|aw_1|.$$

We conclude that $w'$ is a word of exponent greater than $7/3$, which is a contradiction. $\qquad\square$

Theorem 2 now follows from Lemma 10. Consider words $h(0), h^2(0), \ldots, h^k(0), \ldots$. By Lemma 10, these words don't contain subwords of exponent greater than $7/3$. On the other hand, since both $h(0), h(1)$ are of length 21 and contain ten 1's, then $\rho(h^k(0)) = \frac{10}{21}|h^k(0)|$. By Lemma 2, we conclude that $\rho(7/3 + \varepsilon) \leq \frac{10}{21}$.

By Lemma 8, $\rho(7/3 + 0) \leq \frac{10}{21}$ which proves that $\rho(x)$ has a jump at $x = 7/3$.

Now we show that, besides $x = 7/3$, the generalized function $\rho(x)$ is discontinuous to the right at all integer points $x \geq 3$. We use the following lemma which is somewhat similar to Lemma 6. Recall that $\alpha_i = 0^i 1$.

**Lemma 11** *Let $A = \{a_1, \ldots, a_k\}$ and $n \geq 3$. Let $h : A \to \{0,1\}$ be a morphism such that $h(a_i) = \alpha_{m_i}$, where $m_i \leq n$ for all $1 \leq i \leq k$, and $m_i \neq m_j$ for all $i \neq j$. Then for every $(n-1)$-th power-free word $w \in A^*$, $h(w)$ is $(n + \varepsilon)$-th power-free.*

---

[2]This does not hold for subwords of length 15. For example, 010011001011001 occurs in $h(01)$ as well as in $h(10)$, and these occurrences have different precursors.

**Proof:** Similar to Lemma 6, morphism $h$ is injective, and every subword of $h(w)$ ending with 1 and preceeding by 1 is the image of a subword of $w$.

Assume that $h(w)$ is not $(n+\varepsilon)$-th power-free. Then it contains a subword $u^n a$ for a non-empty word $u \in \{0,1\}^*$ and $a$ the first letter of $u$. If $u$ contains no 1's, then $u^n a$ contains at least $n+1$ consecutive 0's, which is impossible as $h(w)$ is a concatenation of $\alpha_{m_i}$'s, and $m_i \leq n$. Assume that $u$ contains at least one 1, that is $u = 0^p 1 u'$, $p \geq 0$, $u' \in \{0,1\}^*$. Then $u^n = (0^p 1 u')^n = 0^p 1 v^{n-1} u'$ for $v = u' 0^p 1$. By properties of morphism $h$, each occurrence of $v$ is the image of some subword of $w$ under morphism $h$. Since this subword is the same for all occurrences of $v$, then $w$ contains a subword $(h^{-1}(v))^{n-1}$ which contradicts to $n$-th power-freeness of $w$. $\qquad\square$

**Lemma 12** *For every* $n \geq 4$,
$$\rho(n+\varepsilon) \leq \frac{1}{n+1-\rho(n-1)} \tag{5}$$

**Proof:** Denote $h_n : \{0,1\}^* \to \{0,1\}^*$ the morphism defined by $h_n(0) = \alpha_n$, $h_n(1) = \alpha_{n-1}$. Let $w_l$ be an $(n-1)$-th power-free word of length $l$ with minimal number of 1's ($\rho(w_l) = \rho(n-1,l)$). Clearly, $|h_n(w_l)| = (n+1)(l - c(w_l)) + nc(w_l) = (n+1)l - c(w_l)$, and $c(h_n(w_l)) = l$. By Lemma 11, $h_n(w_l)$ is $(n+\varepsilon)$-th power-free, and we have

$$\rho(n+\varepsilon, |h_n(w_l)|) \leq \rho(h_n(w_l)) = \frac{l}{(n+1)l - c(w_l)} = \frac{1}{n+1-\rho(n-1,l)}$$

By taking the limit for $l \to \infty$ (see Lemma 2), inequality (5) follows. $\qquad\square$

Inequality (5) together with the trivial inequality $\rho(n-1) \leq 1/2$ gives $\rho(n+\varepsilon) \leq \frac{1}{n-1/2} < \frac{1}{n}$ for $n \geq 4$. On the other hand, from lower bound (3) it follows that $\rho(n) \geq \frac{n-1}{n^2-n-1} > \frac{1}{n}$. This implies that $\rho(n+0) = \rho(n+\varepsilon) < \rho(n)$, that is $\rho(x)$ has a jump to the right of all integer points $n \geq 4$.

For $n = 3$, inequality (5) does not make sense ($\rho(2)$ is not defined). Therefore, the case $n = 3$ should be analysed separately.

**Lemma 13** $\rho(3+\varepsilon) \leq \frac{1}{3}$.

**Proof:** Take a 3-letter alphabet $A = \{1,2,3\}$. For $w \in A^*$, let $c_i(w)$ ($i = 1,2,3$) denote the number of occurrences of $i$ in $w$. For any $l \in \mathbb{N}$, choose a square-free word $w_l \in A^*$ of length $l$ such that $c_1(w) \leq c_2(w) \leq c_3(w)$. Note that for all $l \in \mathbb{N}$, $w_l$ is well-defined, which follows from the existence of infinite square-free words on the 3-letter alphabet. Consider the morphism $h : A^* \to \{0,1\}^*$ defined by $h(1) = 01$, $h(2) = 001$, $h(3) = 0001$. Then $|h(w_l)| = 2c_1(w_l) + 3c_2(w_l) + 4c_3(w_l) = 3l + (c_3(w_l) - c_1(w_l)) \geq 3l$, and $\rho(h(w_l)) \leq \frac{l}{3l} = \frac{1}{3}$. By Lemma 11, word $w_l$ is $(3+\varepsilon)$-th power-free, and then $\rho(3+\varepsilon, |h(w_l)|) \leq \frac{1}{3}$. Taking the limit for $l \to \infty$ and using Lemma 2, we get $\rho(3+\varepsilon) \leq \frac{1}{3}$. $\square$

On the other hand, from lower bound (3) it follows that $\rho(3) \geq \frac{2}{5}$. Therefore, $\rho(x)$ has a jump to the right of $x = 3$.

Putting all together, we obtain

**Theorem 3** $\rho(x)$ *is discontinuous to the right of* $x = \frac{7}{3}$ *as well as to the right of all natural points* $n \geq 3$.

## 4.2   Estimating $\rho(n + \varepsilon)$

In Section 3 we obtained an estimate of $\rho(n)$, for natural $n \geq 3$ (Theorem 1). Theorem 3 says that $\rho(x)$, considered as function on real argument, has a jump on the right of all these points. In this final part of the paper, we estimate the size of these jumps by estimating the values $\rho(n + \varepsilon)$ for natural $n \geq 3$. Recall that $\rho(n + \varepsilon) = \rho(n + 0)$ by Lemma 8.

We start with proving the lower bound

$$\rho(n + \varepsilon) \geq \frac{n - 1}{n^2 - 2} \tag{6}$$

for all $n \geq 3$. The proof follows closely the proof of lower bound (3) from Section 3. Therefore, we only give a sketch of it, underlining the differences with the proof of Section 3.

Consider a finite $(n + \varepsilon)$-th power-free word $w$. As in Section 3, we group its letters into blocks $\alpha_i = 0^i 1$, where $0 \leq i \leq n$ ($w$ may contain $n$-th powers). Again, we assume that $w$ does not start with $\alpha_n$, otherwise we remove the first 0 into a separate prefix. We now note that under this assumption, $w$ cannot contain a subword $(\alpha_n)^n$. Indeed, since $w$ does not start with $\alpha_n$, the occurrence of $(\alpha_n)^n$ is preceeded by at least one letter. This latter cannot be 0, as $w$ would then have a subword $0^{n+1}$ which contradicts to the fact that $w$ does not contain subwords of exponent greater than $n$. This letter cannot be 1 either, as this would give the subword $(10^n)^n 1$ which again contradicts to the fact that $w$ is $(n + \varepsilon)$-th power-free. Thus, no occurrence $(\alpha_n)^n$ exists.

We then group $\alpha_i$'s into blocks $\beta(m, k) = (\alpha_n)^m \alpha_k$, $0 \leq m \leq n - 1$, $0 \leq k \leq n - 1$, and then further into blocks

$$\gamma(l, k_0, k_1, \ldots, k_s) = \beta(l, k_0)\beta(n - 1, k_1) \ldots \beta(n - 1, k_s) = (\alpha_n)^l \alpha_{k_0} (\alpha_n)^{n-1} \alpha_{k_1} \ldots (\alpha_n)^{n-1} \alpha_{k_s},$$

where $0 \leq l \leq n - 2$, $s \geq 0$, $0 \leq k_0, k_1, \ldots, k_s \leq n - 2$.

We now compute the minimal value of $\rho(\gamma(l, k_0, k_1, \ldots, k_s))$. Consider a block $\gamma(l, k_0, k_1, \ldots, k_s)$. As in Section 3, we distingish two cases:

**Case $s \geq 1$:**   Here we show that $k_j + k_{j+1} \leq n$ for every $j$, $0 \leq j \leq s - 1$. By contradiction, assume that $k_j + k_{j+1} > n$. If $k_j = 0$ then $k_{j+1} > n$ which is a contradiction. Assume $k_j > 0$, and consider the subword $\alpha_{k_j}(\alpha_{n-1})^{n-1}\alpha_{k_{j+1}}$. If $k_j + k_{j+1} > n$, then it has the prefix $(0^{k_j}10^{n-k_j})^n$ followed by at least one 0. This gives a subword of exponent greater than $n$ which is a contradiction.

Now we can bound $\sum_{j=0}^{s} |\alpha_{k_j}| \leq (\frac{s}{2} + 1)n + s$, and then $|\gamma(l, k_0, k_1, \ldots, k_s)| \leq l(n+1) + s(n^2 + \frac{n}{2}) + n$. Then $\rho(\gamma(l, k_0, k_1, \ldots, k_s)) \geq \frac{l + ns + 1}{l(n+1) + s(n^2 + \frac{n}{2}) + n}$. Again, the right-hand side minimizes when $l$ is maximal ($l = n - 2$) and $s$ is minimal ($s = 1$). Finally for this case, $\rho(\gamma(l, k_0, k_1, \ldots, k_s)) \geq \frac{2n - 1}{2n^2 + \frac{n}{2} - 2}$.

**Case $s = 0$:**   In this case $\gamma(l, k_0) = \beta(l, k_0)$, $|\gamma(l, k_0)| = l(n + 1) + k_0 + 1$, and $\rho(\gamma(l, k_0)) = \frac{l+1}{l(n+1)+k_0+1}$. The right-hand mininizes when both $l$ and $k_0$ are maximal ($l = n - 2$, $k_0 = n - 1$), which gives $\rho(\gamma(l, k_0)) \geq \frac{n-1}{n^2-2}$.

The second case gives a smaller or equal bound for all $n \geq 3$ and we conclude that $\rho(\gamma(l, k_0, \ldots, k_s)) \geq \frac{n-1}{n^2-2}$. Since $w$ is a concatenation of blocks $\gamma$ (with possibly remaining prefix and suffix of bounded length), this implies inequation (6).

Turning to asymptotic expansion of (6), we have

$$\rho(n + \varepsilon) \geq \frac{1}{n} - \frac{1}{n^2} + \frac{2}{n^3} - \frac{2}{n^4} + \mathcal{O}(\frac{1}{n^5}). \tag{7}$$

To obtain the lower bound of $\rho(n + \varepsilon)$ that matches upper bound (7), it suffices to substitute into inequality (5) the upper bound of $\rho(n-1)$ implied by (1) (instead of trivial upper bound $\rho(n-1) \le \frac{1}{2}$).

We then get

$$\rho(n + \varepsilon) \le \frac{1}{n + 1 - \left(\frac{1}{n-1} + \frac{1}{(n-1)^3} + \frac{1}{(n-1)^4} + \frac{C}{(n-1)^5}\right)} = \frac{1}{n} - \frac{1}{n^2} + \frac{2}{n^3} - \frac{2}{n^4} + \mathcal{O}(\frac{1}{n^5})$$

Together with (7), this gives

**Theorem 4** $\rho(n + \varepsilon) = \frac{1}{n} - \frac{1}{n^2} + \frac{2}{n^3} - \frac{2}{n^4} + \mathcal{O}(\frac{1}{n^5})$.

# 5  Concluding remarks

In this paper we initiated the study of minimal density function for the words avoiding a set of patterns. In particular, we analysed the minimal density $\rho(x)$ of a letter in binary words that don't contain an exponent greater than or equal to $x$. We proved $\rho(x)$ to be discontinuous to the right of point $7/3$ as well as of all integer points starting from 3, and we gave an estimate of values $\rho(n)$ and $\rho(n + \varepsilon)$.

Many questions about minimal density function $\rho(x)$ remain open. Does it have other discontinuities? What are they? Is this function piece-wise constant? All these questions are still to be answered.

Another direction of generalizing the results of this paper is to consider the general notion of $k$-avoidability of a pattern (see Introduction). The general question is: If a pattern $p$ is not $k$-avoidable but is $(k+1)$-avoidable, what is the minimal frequency of a letter in an infinite word over $(k+1)$ letters avoiding $p$? For example, what is the minimal frequency of a letter in an infinite ternary square-free word? A pattern which is 4-avoidable but not 3-avoidable is given in [2]. What is the minimal proportion of the 4th letter needed to avoid that pattern?

# References

[1] А.И. Зимин. Блокирующие множества термов. *Математический Сборник*, 119(3):363–375, 1982. English Translation: A.I.Zimin, Blocking sets of terms, Math. USSR Sbornik 47 (1984), 353-364.

[2] K. Baker, G. McNulty, and W. Taylor. Growth problems for avoidable words. *Theoret. Comp. Sci.*, 69:319–345, 1989.

[3] D. Bean, A. Ehrenfeucht, and G. McNulty. Avoidable patterns in strings of symbols. *Pacific J. Math.*, 85(2):261–294, 1979.

[4] J. Berstel. Axel thue's work on repetitions in words. Invited Lecture at the 4th Conference on Formal Power Series and Algebraic Combinatorics, Montreal, 1992, June 1992. accessible at http://www-igm.univ-mlv.fr/~berstel/index.html.

[5] J. Berstel and D. Perrin. *Theory of codes*. Academic Press, 1985.

[6] J. Cassaigne. *Motifs évitables et régularités dans les mots*. Thèse de doctorat, Université Paris VI, 1994.

[7] C. Choffrut and J. Karhumäki. Combinatorics of words. In G. Rozenberg and A. Salomaa, editors, *Handbook on Formal Languages*, volume I. Springer, Berlin-Heidelberg-New York, 1997.

[8] M. Crochemore and P. Goralcik. Mutually avoiding ternary words of small exponent. *International Journal of Algebra and Computation*, 1(4):407–410, 1991.

[9] J. Currie. Open problems in pattern avoidance. *American Mathematical Monthly*, 100:790–793, 1993.

[10] J. Currie and R. Shelton. Cantor sets and Dejean's conjecture. *Journal of Automata, Languages and Combinatorics*, 1(2):113–128, 1996.

[11] F. Dejean. Sur un théorème de Thue. *J. Combinatorial Th. (A)*, 13:90–99, 1972.

[12] M. Dekking. On the Thue-Morse measure. *Acta Univ. Carolin. Math. Phis*, 33(2):35–40, 1992.

[13] A. Kfoury. A linear time algorithm testing whether a word contains an overlap. *RAIRO Inf. Th.*, 22:135–145, 1988.

[14] R. Kolpakov and G. Kucherov. Minimal letter frequency in $n$-power-free binary words. In I. Privara and P. Ružička, editors, *Proceedings of the 22nd International Symposium on Mathematical Foundations of Computer Science (MFCS), Bratislava (Slovakia)*, volume 1295 of *Lecture Notes in Computer Science*, pages 347–357. Springer Verlag, 1997.

[15] M. Lothaire. *Combinatorics on Words*, volume 17 of *Encyclopedia of Mathematics and Its Applications*. Addison Wesley, 1983.

[16] F. Mignosi and G. Pirillo. Repetitions in the Fibonacci infinite word. *RAIRO Theoretical Informatics and Applications*, 26(3):199–204, 1992.

[17] F. Mignosi, A. Restivo, and S. Salemi. A periodicity theorem on words and applications. In *Proceedings of the 20th International Symposium on Mathematical Foundations of Computer Science (MFCS)*, volume 969 of *Lecture Notes in Computer Science*, pages 337–348. Springer Verlag, 1995.

[18] A. Restivo and S. Salemi. On weakly square free words. *Bull. of the EATCS*, 21:49–56, 1983.

[19] P. Roth. Every binary pattern of length six is avoidable on the two-letter alphabet. *Acta Informatica*, 29:95–106, 1992.

[20] A. Salomaa. *Jewels of formal language theory*. Computer Science Press, 1986.

[21] M. Sapir. Combinatorics on words with applications, December 1993. accessible at http://www.math.unl.edu/~msapir/ftp/course.

[22] A. Thue. Über unendliche Zeichenreihen. *Norske Vid. Selsk. Skr. I. Mat. Nat. Kl. Christiania*, 7:1–22, 1906.

[23] A. Thue. Über die gegenseitige Lage gleicher Teile gewisser Zeichenreihen. *Norske Vid. Selsk. Skr. I. Mat. Nat. Kl. Christiania*, 10:1–67, 1912.