

Grasshopper avoidance of patterns

Michał Dębski Urszula Pastwa Krzysztof Węsek

Faculty of Mathematics and Information Science
Warsaw University of Technology
Warsaw, Poland

michaldebski87@gmail.com, urszula@pastwa.pl, k.wesek@mini.pw.edu.pl

Submitted: Jun 10, 2016; Accepted: Oct 15, 2016; Published: Oct 28, 2016

Mathematics Subject Classifications: 68R15

Abstract

Motivated by a geometrical Thue-type problem, we introduce a new variant of the classical pattern avoidance in words, where jumping over a letter in the pattern occurrence is allowed. We say that pattern $p \in E^+$ *occurs with jumps* in a word $w = a_1a_2 \dots a_k \in A^+$, if there exist a non-erasing morphism f from E^* to A^* and a sequence (i_1, i_2, \dots, i_l) satisfying $i_{j+1} \in \{i_j + 1, i_j + 2\}$ for $j = 1, 2, \dots, l - 1$, such that $f(p) = a_{i_1}a_{i_2} \dots a_{i_l}$. For example, a pattern xx occurs with jumps in a word $abdcadbc$ (for $x \mapsto abc$). A pattern p is *grasshopper k -avoidable* if there exists an alphabet A of k elements, such that there exist arbitrarily long words over A in which p does not occur with jumps. The minimal such k is the *grasshopper avoidability index* of p . It appears that this notion is related to two other problems: pattern avoidance on graphs and pattern-free colorings of the Euclidean plane. In particular, we show that a sequence avoiding a pattern p with jumps can be a tool to construct a line p -free coloring of \mathbb{R}^2 .

In our work, we determine the grasshopper avoidability index of patterns α^n for all n except $n = 5$. We also show that every doubled pattern is grasshopper $(2^7 + 1)$ -avoidable, every pattern on k variables of length at least 2^k is grasshopper 37-avoidable, and there exists a constant c such that every pattern of length at least c on 2 variables is grasshopper 3-avoidable (those results are proved using the entropy compression method).

Keywords: Thue sequence; Avoidable pattern; Entropy compression

1 Introduction

1.1 Avoidance of patterns

A celebrated result by Axel Thue asserts that there exists an infinite sequence on 3 symbols without a square (or repetition), i.e. two consecutive identical blocks of any

positive length. This theorem can be proved constructively using an infinite sequence on 2 symbols without a cube (or 3-repetition), i.e. three consecutive identical blocks – it was independently discovered by Prouhet [23] and by Morse [18], and is generally known as Thue-Morse sequence. The work of Thue [28, 29] inspired a whole branch of combinatorics called combinatorics on words - an area with many challenging questions and variety of applications in other fields of mathematics and science in general [17]. For a meticulous survey on historical roots of the field see the article of Berstel and Perrin [3].

Let A be a finite set of letters and E be a set of variables. We say that a pattern $p \in E^+$ occurs in a word $w \in A^+$, if there exists a non-erasing morphism f from E^* to A^* such that $f(p)$ is a factor of w . For example, a pattern $xyxy$ occurs in a word $abbabbcde$ (for $x \mapsto ab$ and $y \mapsto b$). Otherwise, we say that word w avoids pattern p . A pattern p is said to be k -avoidable if for any alphabet A with k elements, A^* contains infinitely many words in which p does not occur. Minimal such k , if it exists, is the avoidability index of p , denoted by $\mu(p)$. A pattern p is avoidable if it is k -avoidable for some k , otherwise it is unavoidable. The classic Thue's theorem can be restated in this language in the following way: the pattern xx is 3-avoidable. Clearly, the pattern xx is not 2-avoidable - xx occurs in every binary word of length 4. Thus the avoidability index of xx is 3.

Not all patterns are avoidable, for example, xyx . Bea, Ehrenfeucht, McNulty [1] and Zimin [31] provided a complete characterization of avoidable patterns (it can be checked using Zimin's algorithm by reductions of patterns). However, the characterization of k -avoidability seems to be very hard to establish and there is no known characterization for an arbitrary k . In fact, it is not known if there exists a constant C such that every avoidable pattern is C -avoidable. Up to our knowledge, the biggest known avoidability index of a pattern is 5, which was provided by Clark [7]. Moreover, on the algorithmic side, the following question is open: Is it decidable, given pattern p and integer k , whether p is k -avoidable?

Nevertheless, there are some partial results on k -avoidability of patterns. For example, all binary patterns are completely classified: a finite number of them is unavoidable, a finite number of them has the avoidability index equal to 3, the rest of binary patterns (including all of length at least 6) has the avoidability index 2. The classification was established in parts by Schmidt [25, 26], Roth [24] and Cassaigne [4]. All ternary patterns are also classified: it was started by Cassaigne [5] and finished by Ochem [20].

Other "simple" class of patterns worth mentioning is the class of *doubled* patterns, i.e. in which every variable occurs at least twice. What if some pattern p has a variable used only once? Assume we are trying to find an occurrence of p in some word w . Then, since such variable can be mapped to any nonempty subword (without any demanded relation to the rest of the pattern occurrence), this variable gives us much freedom. Such variable stands for "jumping" over any positive number of letters in the word when we are trying to find p in w . Restriction to doubled patterns makes avoiding easier: every such pattern is 3-avoidable (proved in parts by Bell and Goh [2], Cassaigne [5], Ochem [20, 21]). In this case, it is also possible to impose 2-avoidability by a stronger condition on the length of the pattern in terms of the number of variables, see work by Zydr  n [32].

Another direction of research is to consider patterns which are sufficiently long in

terms of the number of variables. Intuitively, it should be easier to avoid such patterns, and this intuition is confirmed. An important observation in this case is that if a pattern is long enough, then it contains a doubled pattern as a block. The best results of this form were presented by Ochem and Pinlou [22]: Assume a pattern p has l variables. If p has length at least 2^l then it is 3-avoidable, and if p has length at least $3 \cdot 2^{l-1}$ then it is 2-avoidable.

Furthermore, there are some general bounds on the avoidability index of an avoidable pattern p in terms of the number of variables of p , although probably far from being optimal. If l is the number of variables in p , then $\mu(p) \leq 4^{\lceil \frac{l+1}{2} \rceil} \leq 2l + 4$ [6]. In fact, for any l there exists a single infinite word over an alphabet of size $2l + 4$ that avoids every avoidable pattern with the number of variables not greater than l .

For a more detailed introduction to the topic, see [16, Chapter “On avoidable patterns” by Cassaigne]. For a survey concentrated on interesting open questions, see the article by Currie [8].

1.2 Grasshopper avoidance of patterns

In our work we consider a new type of pattern avoidability - we allow jumping over letters in the pattern occurrence. Let us imagine a grasshopper going through some word to the right - the grasshopper can jump between two consecutive letters or can jump over exactly one letter. Assume p is a pattern that the grasshopper likes very much. The grasshopper reads every letter it stands on and tries to choose a sequence of letters admitting an occurrence of p . Intuitively, pattern p occurs with jumps in a word w if the grasshopper can choose such a path. Let us define this notion formally.

Definition 1. We say that pattern $p \in E^+$ occurs with jumps in a word $w = a_1 a_2 \dots a_k \in A^+$, if there exist a non-erasing morphism f from E^* to A^* and a sequence (i_1, i_2, \dots, i_d) satisfying $i_{j+1} \in \{i_j + 1, i_j + 2\}$ for $j = 1, 2, \dots, d - 1$, such that

$$f(p) = a_{i_1} a_{i_2} \dots a_{i_d}.$$

We say, that the word $w = a_{i_1} a_{i_1+1} a_{i_1+2} \dots a_{i_d}$ is an *occurrence with jumps* of a pattern p . If $i_{j+1} = i_j + 2$ for some j , we say that index $i_j + 1$ and symbol a_{i_j+1} are *skipped* in the occurrence of p .

For example, a pattern xx occurs with jumps in a word $abdcadbcb$ (for $x \mapsto abc$).

Definition 2. Pattern p is said to be grasshopper k -avoidable, if for a k -element alphabet A the set A^* contains infinitely many words in which p does not occur with jumps.

It is equivalent, by König’s Lemma, to the existence of one infinite word (we can also require this word to be doubly infinite) with letters from A avoiding p with jumps. We say that this word *avoids p with jumps*. We use all three definitions depending on context.

Definition 3. Pattern p is grasshopper avoidable, if it is grasshopper k -avoidable for some finite k . For every such p we define grasshopper avoidability index

$$\mu'(p) = \min\{k : p \text{ is grasshopper } k\text{-avoidable}\}.$$

Creating a word avoiding p with grasshopper jumps is more difficult than in the classic sense - if a pattern does not occur with jumps, then the pattern does not occur in the classic sense. As an example consider any infinite sequence on 3 symbols. This sequence cannot avoid xx with jumps. Indeed, if there is the same symbol on some positions i and $i + 1$ or some positions i and $i + 2$, then the grasshopper can produce a simple repetition of the same letter. Otherwise, the sequence would be of the form $\dots 012012012012\dots$ which has many repetitions. On the other hand, recall that it is possible to avoid xx with 3 symbols.

1.3 Connections to other problems

The idea of jumps was inspired by a geometrical problem: so-called line nonrepetitive colorings of the Euclidean plane. Grytczuk, Kosiński and Zmarz in their work [14] implicitly constructed a sequence on 6 symbols that avoids the pattern xx with jumps and used it to construct a line nonrepetitive 36-coloring of the plane. The reader interested in more Thue type problems with graph-theoretic, geometrical or number-theoretic flavor can be referred to the survey of Grytczuk [13].

As mentioned before, grasshopper avoidance was partially inspired by research on avoiding repetitions (or other patterns) in geometrical structures. Let us take a closer look at this context. Recently, Grytczuk, Kosiński and Zmarz [14] introduced a Thue-type problem related to the famous Hadwiger–Nelson problem of coloring the Euclidean plane (see [27]). A *line path* is a sequence of distinct collinear points in \mathbb{R}^2 with consecutive distances equal to 1. Grytczuk, Kosiński and Zmarz presented a 36-coloring of \mathbb{R}^2 in which every sequence of colors of a line path avoids pattern x^2 . This concept was investigated further for various patterns by Wenus and Węsek [30] (including improving the result for repetitions to 18 colors) and by Dębski, Grytczuk, Pastwa, Pilat, Sokół, Tuczyński, Wenus and Węsek [10] (studies concentrated on 2-colorings). For an arbitrary pattern p , as one can guess, the goal is to construct a coloring of \mathbb{R}^2 which avoids p on every sequence of colors of a line path. A coloring satisfying this condition is called *p -free*. Let us call a pattern p *line avoidable* if there exists a p -free coloring of \mathbb{R}^2 . It is easy to observe that every line avoidable pattern has to be avoidable. In [30], authors suggested that the converse implication may also be true (with a confirmation for doubled patterns and patterns of length at least 2^l , where l stands for the number of variables).

Conjecture 4. [30] A pattern is avoidable if and only if it is line avoidable.

Examples of grasshopper avoidance were used to construct a line xx -free colorings in [14] and [30]. It appears that this is a more general property: for a given pattern p , sequence avoiding p with jumps can be used as a tool to prove line avoidability of p .

Proposition 5. *For every grasshopper k -avoidable pattern p , there exists a line p -free k^2 -coloring of \mathbb{R}^2 .*

Proof. Let $(a_n)_{n \in \mathbb{Z}}$ be a sequence over a k -elemental alphabet A in which p does not occur with jumps. Let $f : \mathbb{R}^2 \rightarrow A^2$ be the following coloring:

$$f((x, y)) = (a_{\lfloor x\sqrt{2} \rfloor}, a_{\lfloor y\sqrt{2} \rfloor}).$$

Let $P = (x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$ be a line path in \mathbb{R}^2 and L be the straight line containing P . Without loss of generality we assume that the counterclockwise angle α between L and x -axis is in the set $[-\frac{\pi}{4}, \frac{\pi}{4}]$. Otherwise, we can reverse the order in P or switch the names of x and y coordinates. For any natural i , we have $x_{i+1} - x_i = \cos \alpha \in [\frac{1}{\sqrt{2}}, 1]$. Hence, $\lfloor x_{i+1}\sqrt{2} \rfloor - \lfloor x_i\sqrt{2} \rfloor \in \{1, 2\}$. Since there is no occurrence with jumps of p in the sequence $(a_{\lfloor x_i\sqrt{2} \rfloor})_{i=1}^n$, we conclude that the sequence $(f(x_i))_{i=1}^n$ avoids p . \square

We find it also worth pointing out that a bit different concept of jumps can be found in the literature. Currie and Simpson [9] introduced the following generalization of Thue sequences, also originating in geometrical problems. Let a sequence a be k -Thue if every j -subsequence avoids repetitions, for $1 \leq j \leq k$, where j -subsequence is any sequence of the form $a_i a_{i+j} a_{i+2j} \dots$. Grytczuk [11] conjectured that for any k it is possible to construct a k -Thue sequence with $k + 2$ symbols. The best effort in this matter is the recently proved sufficiency of $2k$ symbols due to Kranjc, Luzar, Mockovciakova and Sotak [15]. Note the relation with our notion: In the language of grasshopper, the definition of 2-Thue sequence takes into account only two extremal possible behaviors of the grasshopper - that means a repetition can occur with jumps in a 2-Thue sequence.

1.4 Major open problems

Consider an arbitrary pattern p . If p can be avoided with jumps in a sequence over k symbols, then the same sequence is a witness of k -avoidability of p - although, it is possible that for the classic avoidance less than k symbols may be sufficient. On the other hand, for a k -avoidable p , it may be reasonable to take into account that k symbols can be far from being enough to construct a sequence avoiding p with jumps. But are there any avoidable patterns that are not grasshopper avoidable? Does the possibility of jumping actually changes the property of being avoidable over a sufficiently large alphabet? We believe that the answer is negative and hence we state the following conjecture.

Conjecture 6. A pattern is avoidable if and only if it is grasshopper avoidable.

It follows from Proposition 5 that Conjecture 6 would imply Conjecture 4.

Thue-type problems have been considered also in the area of graph colorings. A lot of studies concentrate on vertex colorings satisfying the condition that the sequence of colors of every simple path avoids repetitions or, more generally, a given pattern. For a pattern p , such a coloring is called p -free. In his work, Grytczuk [12] considered the following concept of pattern avoidance on graphs. We say that a pattern p is *avoidable on graphs* if for any Δ there exists a constant k such that every graph with the maximal degree at most Δ has a p -free k -coloring. Grytczuk stated the following conjecture (again, with a confirmation for doubled patterns and patterns of length at least 2^d , where d stands for the number of variables).

Conjecture 7. [12] A pattern is avoidable if and only if it is avoidable on graphs.

Consider an infinite graph G defined on the set \mathbb{Z} with two integers a, b joined by an edge if $0 < |a - b| \leq 2$. Suppose that for a given avoidable pattern p we have a p -free

k -coloring of G . Then, this coloring produces a sequence avoiding p with jumps (since any jumping subsequence corresponds to a simple path in G). Therefore, Conjecture 7 would imply Conjecture 6. We can summarize the relation between conjectures stated in this section:

$$\text{Conjecture 7} \Rightarrow \text{Conjecture 6} \Rightarrow \text{Conjecture 4}$$

We believe that Conjecture 6, besides being interesting in itself, is also a good approach to make a progress in Conjecture 4. On the other hand, Conjecture 7 seems to be much stronger - but any step forward in Conjecture 6 may be an important suggestion for Conjecture 7.

1.5 Our contributions

We investigate grasshopper avoidability of patterns known to be avoidable in the classic sense, starting with patterns using one variable (that is, xx, xxx, \dots); we determine the grasshopper avoidability index of x^n for all values of n except $n = 5$ (Theorem 8). In particular, x^6 can be avoided with jumps on 2 symbols.

We also show that doubled patterns (that is, patterns where every variable is used at least twice) are grasshopper avoidable on $2^7 + 1$ symbols (Theorem 18). Note that it proves Conjecture 6 for this class of patterns. However, $2^7 + 1$ symbols is probably far from optimal (recall that doubled patterns are 3-avoidable).

Our further investigations concern patterns with low grasshopper avoidability index. We show that patterns on 2 variables that are sufficiently long are 3-grasshopper avoidable (Theorem 22). A similar, but weaker statement is given for patterns with more variables: if a pattern on k variables is longer than 2^k , then it is 37-grasshopper avoidable (Theorem 21).

Note that Theorems 18 and 21 are best possible, as far as grasshopper avoidability is concerned. For every k there exists an unavoidable pattern Z_k on k variables such that the length of Z_k is $2^k - 1$ and exactly one variable appears in Z_k once. Those patterns are known as Zimin words (or: sesquipowers), and are constructed as follows: $Z_1 = x_1$ and $Z_{k+1} = Z_k x_{k+1} Z_k$.

Theorem 8 is proved in Section 2, using elementary arguments. A part of it (that $\mu'(x^2) \geq 6$) is obtained by an extensive computer search, and we do not know any simple argument. The bounds $\mu'(x^6) \leq 2$, $\mu'(x^3) \leq 3$ and $\mu'(x^2) \leq 6$ are obtained using the Thue and Thue-Morse sequence, and $\mu'(x^4) > 2$ follows from a direct construction of an occurrence of x^4 in any sufficiently large binary word.

Theorems 18, 22 and 21 are proved in Section 3 using the entropy compression method. The proofs are somewhat similar to earlier applications of this method in pattern avoidance (in [19] and [22]), but involve essential modifications in order to handle the letters skipped in occurrences with jumps of considered patterns.

2 Avoiding x^n

The simplest class of patterns is clearly the class of patterns with just one variable. In this section we prove Theorem 8, that give the exact value of grasshopper avoidability index of x^n for $n = 2, 3, 4$ and $n > 5$. It follows that $2 \leq \mu'(x^5) \leq 3$.

Theorem 8. *For any natural number n we have*

$$\mu'(x^n) = \begin{cases} 2 & \text{for } n \geq 6 \\ 3 & \text{for } 4 \geq n \geq 3 \\ 6 & \text{for } n = 2 \end{cases}$$

Proof. We divide the proof into four lemmas (Lemma 9, 10, 11 and 13).

By Lemma 9 we have $\mu'(x^6) = 2$, and it follows that $\mu'(x^n) = 2$ for $n \geq 6$.

Lemma 10 says that $\mu'(x^4) > 2$, and by Lemma 11 we have $\mu'(x^3) \leq 3$. Since $\mu'(x^4) \leq \mu'(x^3)$, we have $\mu'(x^3) = \mu'(x^4) = 3$.

It was shown in [14] that $\mu'(x^2) \leq 6$. By Lemma 13 we get $\mu'(x^2) = 6$, which completes the proof. \square

Lemma 9. *Let w be the Thue sequence over the alphabet $\{0, 1, 2\}$ and let h be a morphism defined as*

$$\begin{aligned} h(0) &= a^5b^5 \\ h(1) &= a^5b^5a^2b^2 \\ h(2) &= a^5b^5a^2b^2a^2b^2. \end{aligned}$$

Then x^6 doesn't occur with jumps in $h(w)$.

Proof. Suppose that the grasshopper can read v^6 from $h(w)$ for some $v \in \{a, b\}^*$. Let $v = v_0v_1v_2 \dots v_k$, where v_1, v_2, \dots, v_{k-1} and v_kv_0 are of the form $(a^2a^*b^3b^* + a^3a^*b^2b^*)((a + a^2)(b + b^2))^*$ (we can make this division by starting new v_i before every $a^2a^*b^3b^*$ or $a^3a^*b^2b^*$). Now, each of the words v_1, v_2, \dots, v_{k-1} and v_kv_0 is read by the grasshopper from an image of a single symbol: $h(0)$, $h(1)$ or $h(2)$, and we can easily recognize from which one. Since v^6 contains $(v_1v_2 \dots v_{k-1}(v_kv_0))^5$, we can find x^5 in w , which contradicts the properties of Thue sequence. \square

Lemma 10. *The pattern x^4 is grasshopper unavoidable on 2 symbols.*

Proof. Suppose that there exists an infinite word w on $\{a, b\}$ in which x^4 does not occur with jumps. Note that w must contain an infinite number of occurrences of aa and bb (otherwise it would contain $(ab)^4$). We can find a subword u in w satisfying $|u| > 32$ and $u = (AB)^*$, where $A = aa(a^*)(ba^+)^*$ and $B = bb(b^*)(ab^+)^*$ (we find first occurrence of aa in w and read it to the next bb , then to the next aa and so on).

Note that if a word from A has length at least 5, then the grasshopper can read a^4 from it, and if a word from B has length at least 5, then the grasshopper can read b^4 from it. Therefore, we may assume that $u = (A'B')^*$, where $A' = \{aaa, aaba, aa\}$ and $B' = \{bbb, bbab, bb\}$.

The grasshopper can read aa from every word from A' and bb from every word from B' (read the first or the second letter and then read the last letter of the word). Because u consists of more than 8 words from A' and B' , the grasshopper can read $(aabb)^4$ from it, contradicting the choice of w . \square

Lemma 11. *Let w be the Thue-Morse sequence on the alphabet $\{0, 1\}$ and let h be a morphism defined as*

$$\begin{aligned} h(0) &= c^2a^2 \\ h(1) &= c^2b^2. \end{aligned}$$

Then x^3 does not occur with jumps in $h(w)$.

Proof. Denote $h(w) = a_0a_1a_2a_3\dots$. Let us suppose that x^3 occurs with jumps in $h(w)$ and let $(i_1, i_2, \dots, i_{3l})$ be a sequence satisfying $i_{j+1} \in \{i_j + 1, i_j + 2\}$ for $i = 1, 2, \dots, 3l - 1$ such that $a_{i_1}a_{i_2}\dots a_{i_{3l}} = v^3$, where v is a word from $\{a, b, c\}^*$.

Without loss of generality we may assume that v starts with c and ends with another letter. Indeed, if v ends with c or begins with a or b , let

$$v' = \begin{cases} c^2u & \text{for } v = cuc, \text{ where } u \text{ starts and ends with } a \text{ or } b \text{ (1a)} \\ c^2u & \text{for } v = uc^2, \text{ where } u \text{ starts and ends with } a \text{ or } b \text{ (1a)} \\ cu & \text{for } v = uc, \text{ where } u \text{ starts and ends with } a \text{ or } b \text{ (1b)} \\ ua^2 & \text{for } v = aua, \text{ where } u \text{ starts and ends with } c \text{ (2b)} \\ ub^2 & \text{for } v = bub, \text{ where } u \text{ starts and ends with } c \text{ (3b)}. \end{cases}$$

The grasshopper can read $(v')^3$ from $h(w)$. Note that v' is always well-defined. Indeed, if (1) v ends with c and (1a) $v = cuc$ or $v = uc^2$ for some u , then u can neither start nor end with c (since, by the definition of h , grasshopper can't read c^2 from $h(w)$). If we have (1) and (1a) does not apply, then (1b) $v = uc$, where u starts and ends with a or b . If (1) does not apply, then (2) v ends with a or (3) v ends with b . Note that (2a) $v = bua$ is not possible, because the grasshopper can't read ab from $h(w)$, so we have (2b) $v = aua$ and u must start and end with c , since every double occurrence of a read by the grasshopper from $h(w)$ must be preceded and succeeded by at least one c . The case (3) is symmetric to (2).

Now, we can divide v into subwords from the set

$$\{ca, caa, cca, ccaa, cb, cbb, ccb, cccb\}.$$

Recall that, by definition of h , $h(w)$ does not contain single occurrences of any letter, so we can replace in v every block from $\{ca, caa, cca, ccaa\}$ by $ccaa$ (respectively by 0) and every block from $\{cb, cbb, ccb, cccb\}$ by $ccbb$, obtaining a word x such that x^3 is a subword of $h(w)$. Since $h^{-1}(ccaa) = 0$ and $h^{-1}(ccbb) = 1$, $(h^{-1}(x))^3$ is a subword of w , which contradicts the properties of Thue-Morse sequence. \square

The following lemma was proved by Grytczuk, Kosiński and Zmarz [14, Lemma 7]. We give the proof for completeness.

Lemma 12 (Gryczuk, Kosiński, Zmarz [14]). *The pattern x^2 is grasshopper avoidable on 6 symbols.*

Proof. Let w be the Thue sequence over the alphabet $\{0, 1, 2\}$ and h be a morphism defined as:

$$\begin{aligned} h(0) &= aa' \\ h(1) &= bb' \\ h(2) &= cc'. \end{aligned}$$

Denote $h(w) = z_0 z_1 z_2 z_3 \dots$. Let us suppose that x^2 occurs with jumps in $h(w)$ and let $(i_1, i_2, \dots, i_{2l})$ be a sequence satisfying $i_{j+1} \in \{i_j + 1, i_j + 2\}$ for $j = 1, 2, \dots, 2l - 1$ such that $z_{i_1} z_{i_2} \dots z_{i_{2l}} = v^2$, where v is a word from $\{a, a', b, b', c, c'\}^*$.

Note that without loss of generality, we can assume that $z_{i_l} = z_{i_{2l}} \in \{a', b', c'\}$ and $z_{i_1} = z_{i_{l+1}} \in \{a, b, c\}$. Otherwise, we can consider another occurrence of x^2 with jumps:

$$\begin{cases} z_{i_1}^* z_{i_1} \dots z_{i_l} | z_{i_{l+1}}^* z_{i_{l+1}} \dots z_{i_{2l}} & \text{if } z_{i_l} \in \{a', b', c'\} \text{ and } z_{i_{l+1}} \in \{a', b', c'\} \\ z_{i_1} \dots z_{i_l} z_{i_l}^* | z_{i_{l+1}} \dots z_{i_{2l}} z_{i_{2l}}^* & \text{if } z_{i_l} \in \{a, b, c\} \text{ and } z_{i_{l+1}} \in \{a, b, c\} \\ z_{i_1}^* z_{i_1} \dots z_{i_{l-1}} | z_{i_l} z_{i_{l+1}} \dots z_{i_{2l-1}} & \text{if } z_{i_l} \in \{a, b, c\} \text{ and } z_{i_{l+1}} \in \{a', b', c'\} \end{cases}$$

Where $|$ divides two occurrences of x and the $*$ operation changes a letter from $\{a, b, c\}$ to the corresponding letter from $\{a', b', c'\}$ and vice versa.

Then, $z_{i_1} z_{i_{l+1}} \dots z_{i_{2l}}$ is an occurrence of x^2 (without jumps!) in $h(w)$. It follows that $h^{-1}(z_{i_1} z_{i_{l+1}} \dots z_{i_{2l}})$ is a repetition in w , which contradicts the property of Thue sequence. Therefore, the proof is complete. \square

Lemma 13. *The pattern x^2 is grasshopper unavoidable on 5 symbols.*

Proof. The Lemma follows from an exhaustive computer search (the longest words on 5 symbols without an occurrence of x^2 with jumps has length 22). \square

As mentioned, Theorem 8 leaves one case still open: we do not know whether the grasshopper avoidability index of x^5 equals 2 or 3. However, we feel that the answer should be 2. The longest binary sequence avoiding x^5 with jumps that we found is of length 104. It was generated by an exhaustive computer search (for bigger lengths, the program does not finish in reasonable time).

Problem 14. Determine $\mu'(x^5)$.

3 Avoiding patterns via entropy compression

Fix a pattern p on k variables $\alpha_1, \dots, \alpha_k$ and an alphabet A . Consider the following procedure, parametrized by a number m and a sequence S of length m over A . It uses S to generate some word w over A avoiding p with jumps and produces a log L that will allow us to recreate S from w .

1. Let w be an empty word

2. Initialize an empty log L
3. for $n = 1, \dots, m$
 - (a) Append the n -th symbol from S to w
 - (b) If some suffix r of w is an occurrence with jumps of p
 - i. Encode r and n in L
 - ii. Erase r from w
4. return the pair (w, L)

We will show that the output (w, L) uniquely determines the input sequence S , so there should be $|A|^m$ possible outputs (w, L) . Moreover, the returned word w avoids p with jumps. Our plan to show that a pattern p is grasshopper avoidable is the following. First, we suppose that p is not grasshopper avoidable on some sufficiently large alphabet A , so that the length of any word over A that avoids p with jumps is bounded by some constant N . We show that, for sufficiently large m , there are at most $(|A| - \epsilon)^m$ different logs L that can be returned by our procedure - thus the procedure can have at most $|A|^N (|A| - \epsilon)^m$ different outputs, which is a contradiction. Therefore, p must be grasshopper avoidable (and the procedure will return a word w of length $> N$ that avoids p with jumps).

The log L contains the following information:

- The number x (number of erasures)
- The strictly increasing sequence $e = (e_1, e_2, \dots, e_x)$ (steps n in which the erasure was performed)
- The number v (the total length of all words assigned to all variables in all erased occurrences of p)
- The sequence $\ell = (l_1, l_2, \dots, l_{kx})$ (lengths of words assigned to each of the k variables of p)
- The sequence V of length v over A (concatenation of all words assigned to all variables in all erased occurrences with jumps of p)
- The number y (the total number of indices skipped in erased occurrences with jumps of p)
- The sequence Y of length y over A (symbols skipped in erased occurrences with jumps of p)
- Binary sequence B of length $b \leq m - y$ (that for every index not skipped in occurrences with jumps of p encodes if the next index was skipped)

Assume $w_{i_1} \dots w_{i_d}$ (with $i_d = n$) is the created occurrence with jumps of p . Encoding of the pair (r, n) is the following: we increase x by 1 and v by the sum of length of words assigned to all variables of p , we append i to e , we append to ℓ the sequence $(|f(\alpha_1)|, |f(\alpha_2)|, \dots, |f(\alpha_k)|)$, we append to V the concatenation $f(\alpha_1)f(\alpha_2) \dots f(\alpha_k)$, we increase y by the number of indices skipped in r , we append to Y all symbols skipped in r , for every $j \in \{1 \dots d-1\}$ we append 0 to B if $i_{j+1} = i_j + 1$ (meaning that there is no skipped symbol after index i_j), otherwise we append 1.

Note that given the state (w_n, L_n) of our procedure after i -th step, we can determine the n -th symbol of S and the state (w_{n-1}, L_{n-1}) after $(n-1)$ -th step. Indeed, if $e_x = n$, we reverse the above encoding and obtain word w' and log L' , and otherwise we set $w' = w$ and $L' = L$. Now, the n -th element of S is the last symbol of w' , w_{n-1} is the first $|w| - 1$ symbols from w' and $L_{n-1} = L'$. By induction we get the following claim.

Claim 15. *Given the output (w, L) , the input sequence S can be uniquely determined.*

Claim 16. *Let $c \geq 2$. If $x \leq \frac{m}{c}$ and m is sufficiently large, then there are less than γ_c^m possible sequences e , where*

$$(i) \quad \gamma_c \leq 2,$$

$$(ii) \quad \gamma_c \text{ goes to } 1 \text{ when } c \text{ goes to } \infty.$$

Note that there is a bijection from the set of possible sequences e to the family of subsets of the set $[m]$ of order at most $\frac{m}{c}$, so we immediately have (i). To get (ii), we observe that the estimated number is at most $\frac{m}{c} \binom{m}{\frac{m}{c}}$ (since there are more subsets of $[m]$ of size $\frac{m}{c}$ than of any fixed size less than $\frac{m}{c}$) and we bound it using Stirling's formula. We have

$$\binom{m}{\frac{m}{c}} \approx \frac{\left(\frac{m}{e}\right)^m \sqrt{2\pi m}}{\left(\frac{\frac{m}{c}}{e}\right)^{\frac{m}{c}} \sqrt{2\pi \frac{m}{c}} \left((1 - \frac{1}{c}) \frac{m}{e}\right)^{(1-\frac{1}{c})m} \sqrt{2\pi (1 - \frac{1}{c}) m}} = O(\sqrt{m}) \left(\frac{1}{\frac{1}{\sqrt{c}} (1 - \frac{1}{c})^{1-\frac{1}{c}}} \right)^m.$$

Clearly $\sqrt[c]{c}$ and $(1 - \frac{1}{c})$ tend to 1 with c tending to infinity, so we get (ii).

Claim 17. *There are at most 2^v possible sequences ℓ .*

Note that ℓ is a sequence of positive integers that adds up to v . If we represent the number v as v intervals separated by $v-1$ delimiters, then each possible sequence ℓ corresponds to a subset of those delimiters, so the claim follows.

Theorem 18. *If p is a doubled pattern, then p is grasshopper avoidable on $2^7 + 1$ symbols.*

Proof. Let A be an alphabet of size $2^7 + 1$ and suppose for the contrary that p is grasshopper unavoidable over A and let n be the length of the longest word over A avoiding p . Let m be larger than some m_0 (that implicitly follows from the proofs) and consider the number of possible outputs of our procedure.

There are $|A|^N$ choices for w and at most m^3 possible assignment of values to x, v and y . By Claim 16 there are less than 2^m choices for e and by Claim 17, at most 2^v choices for ℓ . Moreover, we have $y \leq \frac{m}{2}$ (as every skipped letter is preceded by a non-skipped letter) and $v \leq \frac{m-y}{2}$ (as p is doubled, i.e., every variable occurs at least twice). Now, there are $|A|^v \leq |A|^{\frac{m-y}{2}}$ choices for V , $|A|^y$ choices for Y and at most 2^{m-y} choices for B . It follows that the number of possible outputs (w, L) of the procedure is at most

$$|A|^N \cdot m^3 \cdot 2^m \cdot 2^v \cdot |A|^{\frac{m-y}{2}} \cdot |A|^y \cdot 2^{m-y} = |A|^N m^3 |A|^{\frac{m+y}{2}} 2^{\frac{5}{2}m - \frac{3}{2}y} \leq |A|^N m^3 |A|^{\frac{3}{4}m} 2^{\frac{7}{4}m},$$

where the inequality follows by the fact that, since $|A| > 2^3$, the left side attains maximum value for $y = \frac{m}{2}$. The resulting value is less than $|A|^m$, so we get a contradiction by Claim 15, which completes the proof. \square

Lemma 19. *If p is a doubled pattern on k variables of length at least 2^k , then p is grasshopper avoidable on 37 symbols.*

Proof. For $k = 1, 2$ the result directly follows by Theorem 8, because in these cases the pattern must contain a square. For $k \geq 3$, the length of p is at least 8. The proof goes exactly the same as the proof of Theorem 18, but the bound on the number of possible sequences e obtained from Claim 16 is improved from 2^m to γ_8^m , where a straightforward calculation shows that $\gamma_8 \leq 1.46$ (note that $(1.46)^4 < 4.6$). Therefore the number of possible outputs (w, L) of the procedure is at most

$$|A|^N \cdot m^3 \cdot \gamma_c^m \cdot 2^v \cdot |A|^{\frac{m-y}{2}} \cdot |A|^y \cdot 2^{m-y} = |A|^N m^3 |A|^{\frac{m+y}{2}} 2^{\frac{3(m-y)}{2}} \gamma_8^m.$$

This expression for $y = \frac{m}{2}$ attains maximum value

$$|A|^N m^3 |A|^{\frac{3}{4}m} 2^{\frac{3}{4}m} \gamma_8^m = |A|^N m^3 (|A|^3 2^3 \gamma_8^4)^{\frac{m}{4}} < |A|^N m^3 36.95^m.$$

The resulting value is (for sufficiently large m) smaller than $37^m = |A|^m$, so we get a contradiction which completes the proof. \square

We use the following observation of Ochem and Pinlou.

Claim 20 ([22, Claim 1]). *If p is a pattern on k variables of length at least 2^k , then p contains as a factor a doubled pattern p' on k' variables and of length at least $2^{k'}$.*

As a direct corollary of the above Claim and Lemma 19 we get the following.

Theorem 21. *If p is a pattern on k variables of length at least 2^k then p is grasshopper avoidable on 37 symbols.*

Theorem 22. *There exists a constant c such that every pattern on 2 variables of length at least c is grasshopper avoidable on 3 symbols.*

Proof. Let A be an alphabet of size 3, let c be sufficiently large constant (so that we have $\gamma_c 2^{\frac{4}{c}} |A|^{\frac{4}{c}} < 1.5^m$) and take a pattern p on 2 variables of length at least c . Note that if one of variables of p appears less than $\frac{c}{4}$ times in p , then p must contain the pattern x^3 as a factor, so in this case the result follows by Theorem 8. Therefore, we may assume that each variable of p appears at least $\frac{c}{4}$ times.

Now, we repeat the argument from the proof of Theorem 18 with the following modifications:

- Since each variable of p appears at least $\frac{c}{4}$ times, we have $v \leq \frac{4m}{c}$
- The bound 2^m , obtained from Claim 17, improves to $2^{\frac{4m}{c}}$
- The bound 2^m on the number of possible sequences e is replaced by γ_c^m .

Now, the number of possible outputs becomes at most

$$|A|^N \cdot m^3 \cdot \gamma_c^m \cdot 2^v \cdot |A|^{\frac{4m}{c}} \cdot |A|^y \cdot 2^{m-y} = |A|^N m^3 (2^{\frac{4}{c}} |A|^{\frac{4}{c}} \gamma_c)^m |A|^m \left(\frac{2}{|A|} \right)^{m-y}.$$

By our conditions on c , the above expression is less than $|A|^N m^3 \left(\frac{2 \cdot 1.5}{3} \right)^m |A|^m = |A|^N m^3 |A|^m$, which is less than $|A|^m$ for sufficiently large m . It is the desired contradiction and completes the proof. \square

It remains an interesting question whether for sufficiently long binary patterns we can achieve the minimal nontrivial number of symbols, that is, use an alphabet with only 2 elements.

Problem 23. Is it possible to replace ‘3 symbols’ with ‘2 symbols’ in Theorem 22?

All results from this section remain valid in the list version of the problem (that is, when we are given a family of alphabets $(A_i)_i$, called lists, and demand that i -th element of the pattern-avoiding word is from A_i for all i ; a pattern p is list grasshopper k -avoidable if a word avoiding p with jumps can be constructed for every family of k -element lists). The possibility of such a generalization is inherent to the entropy compression method – for examples of results in classic pattern avoidance actually proven also for the corresponding list version see work of Ochem and Pinlou on ‘long’ patterns [22] and the work of Zydroń [32].

4 Longer jumps

We used grasshopper avoidability of a pattern p to solve the problem of line p -free coloring of the plane. Our considerations can be generalised to higher dimensions if we allow the grasshopper to make longer jumps, up to some fixed length. Formally, we say that a pattern $p \in E^+$ occurs with j -jumps in a word $w = a_1 a_2 \dots a_k \in A^+$, if there exist a

non-erasing morphism f from E^* to A^* and a sequence (i_1, i_2, \dots, i_l) satisfying $i_{n+1} \in \{i_n + 1, i_n + 2, \dots, i_n + j\}$ for $n = 1, 2, \dots, l - 1$, such that

$$f(p) = a_{i_1} a_{i_2} \dots a_{i_l}.$$

For example there is an occurrence with 3-jumps of a pattern xx in $abcdaeab$. Pattern p is said to be j -grasshopper k -avoidable, if for a k -elemental alphabet A the set A^* contains infinitely many words in which p does not occur with j -jumps. Pattern p is j -grasshopper avoidable, if it is j -grasshopper k -avoidable for some finite k . For every such p we define j -grasshopper avoidability index

$$\mu_j(p) = \min\{k : p \text{ is } j\text{-grasshopper } k\text{-avoidable}\}.$$

Note that 1-grasshopper avoidability is exactly the classic avoidability.

Proposition 5 generalizes to higher dimensions (the proof is the same, except that there are more dimensions and $\sqrt{2}$ is replaced by \sqrt{d}).

Proposition 24. *For every $\lceil \sqrt{d} \rceil$ -grasshopper k -avoidable pattern p , there exists a line p -free k^2 -coloring of \mathbb{R}^d .*

The j -grasshopper avoidability index of x^2 is at most $3j$. To see this, take the Thue sequence over the alphabet $\{a, b, c\}$ and replace each letter with a sequence of j distinct symbols (for example, replace a with $a_1 a_2 \dots a_j$). It is easy to see that the resulting sequence j -grasshopper avoids x^2 (see [30, Remark 9]). The bound $\mu_j(x^2) \leq 3j$ is tight for $j \in \{1, 2\}$ and it is interesting whether it is also the case for larger j .

One can ask for a generalization of Conjecture 6 for grasshopper j -avoidance. We believe that allowing longer jumps should not change the class of avoidable patterns.

Conjecture 25. A pattern is avoidable if and only if it is j -grasshopper avoidable for every $j \in \mathbb{Z}_+$.

Observe that using a similar argument as in Section 1.4 we can show that Conjecture 7 would imply Conjecture 25. Therefore, Conjecture 25 is true for doubled patterns and patterns of length at least 2^k , where k stands for the number of variables.

Our approach from Section 3 (in particular, Theorem 22) can be generalised to j -grasshopper avoidability. However, the required number of symbols would rapidly grow with j . It would be interesting to find a way to reduce the required size of the alphabet. In particular, the following generalization of Theorem 22 seems plausible.

Conjecture 26. For every j there exists a constant $c = c(j)$ such that every pattern on 2 variables of length at least c is j -grasshopper avoidable on 3 symbols.

Acknowledgements

Research supported by the National Science Center of Poland, grant 2015/17/B/ST1/02660.

References

- [1] D. R. Bean, A. Ehrenfeucht, and G. McNulty. Avoidable patterns in strings of symbols. *Pacific Journal of Mathematics*, 85:261–294, 1979.
- [2] J. P. Bell and T. L. Goh. Exponential lower bounds for the number of words of uniform length avoiding a pattern. *Information and Computation*, 205(9):1295–1306, 2007.
- [3] J. Berstel and D. Perrin. The origins of combinatorics on words. *European Journal of Combinatorics*, 28(3):996–1022, 2007.
- [4] J. Cassaigne. Unavoidable binary patterns. *Acta Informatica*, 30:385–395, 1993.
- [5] J. Cassaigne. *Motifs évitables et régularités dans les mots*. PhD thesis, Université Paris VI, 7 1994.
- [6] J. Cassaigne. Unavoidable patterns. In Lothaire, editor, *Algebraic Combinatorics on Words*. Cambridge University Press, Cambridge, 2002.
- [7] R. J. Clark. The existence of a pattern which is 5-avoidable but 4-unavoidable. *International Journal of Algebra and Computation*, 16(2):351–367, 2006.
- [8] J. D. Currie. Pattern avoidance: themes and variations. *Theoretical Computer Science*, 339(1):7–18, 2005.
- [9] J. D. Currie and J. Simpson. Non-repetitive tilings. *The Electronic Journal of Combinatorics*, 9:2–8, 2002, #R28.
- [10] M. Dębski, J. Grytczuk, U. Pastwa, B. Pilat, J. Sokół, M. Tuczyński, P. Wenus, and K. Węsek. On avoiding r -repetitions in \mathbb{R}^2 . Unpublished results.
- [11] J. Grytczuk. Thue-like sequences and rainbow arithmetic progressions. *The Electronic Journal of Combinatorics*, 9, 2002, #R44.
- [12] J. Grytczuk. Pattern avoidance on graphs. *Discrete Mathematics*, 307(11–12):1341–1346, 2007.
- [13] J. Grytczuk. Thue type problems for graphs, points, and numbers. *Discrete Mathematics*, 308(19):4419–4429, 2008.
- [14] J. Grytczuk, K. Kosiński, and M. Zmarz. Nonrepetitive colorings of line arrangements. *European Journal of Combinatorics*, 51:275–279, 2016.
- [15] J. Kranjc, B. Luzar, M. Mockovciaková, and R. Soták. On a generalization of Thue sequences. *The Electronic Journal of Combinatorics*, 22(2), 2015, #P2.33 .
- [16] M. Lothaire. *Algebraic combinatorics on words*. Cambridge University Press, Cambridge, UK, 2002.
- [17] M. Lothaire. *Applied combinatorics on words*. Cambridge University Press, Cambridge, UK, 2005.
- [18] M. Morse. A one-to-one representation of geodesics on a surface of negative curvature. *American Journal of Mathematics*, 43:35–51, 1921.

- [19] R. A. Moser. A constructive proof of the Lovász local lemma. In *Proceedings of the Forty-first Annual ACM Symposium on Theory of Computing*, STOC '09, pages 343–350, 2009.
- [20] P. Ochem. A generator of morphisms for infinite words. *RAIRO - Theoretical Informatics and Applications*, 40(3):427–441, 10 2006.
- [21] P. Ochem. Doubled Patterns are 3-Avoidable, 23 *The Electronic Journal of Combinatorics*, 23(1), 2016.
- [22] P. Ochem and A. Pinlou. Application of entropy compression in pattern avoidance. *The Electronic Journal of Combinatorics*, 21(2), 2014, #P2.7.
- [23] E. Prouhet. Memoire sur quelques relations entre les puissances des nombres. *C. R. Acad. Sci. Paris*, 33:31, 1851.
- [24] P. Roth. Every binary pattern of length six is avoidable on the two-letter alphabet. *Acta Informatica*, 29(1):95–107, 1992.
- [25] U. Schmidt. Motifs inévitables dans les mots. Technical report, Rapport LITP, Paris VI, 1986. Rapport LITP.
- [26] U. Schmidt. Avoidable patterns on two letters. *Theoretical Computer Science*, 63(1):1–17, 1989.
- [27] A. Soifer. Chromatic number of the plane & its relatives, history, problems and results: An essay in 11 parts. In A. Soifer, editor, *Ramsey Theory*, volume 285 of *Progress in Mathematics*, pages 121–161. Birkhäuser Boston, 2011.
- [28] A. Thue. Über unendliche zeichenreichen. *Norske Vid. Selsk. Skr., I Mat. Nat. Kl., Christiania*, 7:1–22, 1906. (Reprinted in *Selected Mathematical Papers of Axel Thue*, T. Nagell, editor, Universitetsforlaget, Oslo, Norway (1977), pp. 139-158).
- [29] A. Thue. Über die gegenseitige lage gleicher teile gewisser zeichenreihen. *Norske Vid. Selsk. Skr., I Mat. Nat. Kl., Christiania*, 1:1–67, 1912. (Reprinted in *Selected Mathematical Papers of Axel Thue*, T. Nagell, editor, Universitetsforlaget, Oslo, Norway (1977), pp. 413-478).
- [30] P. Wenus and K. Węsek. Nonrepetitive and pattern-free colorings of the plane. *European Journal of Combinatorics*, 54:21–34, 2016.
- [31] A. I. Zimin. Blocking sets of terms. *Mathematics of the USSR-Sbornik*, 47(2):353, 1984.
- [32] A. Zydrón. Unikalność bezjednostkowych wzorców o dużej liczbie zmiennych. Master’s thesis, Jagiellonian University, Poland, 2013.