For Dagstuhl Seminar

Combinatorial and Algorithmic Aspects of Sequence Processing 11081 Organizers: Maxime Crochemore, Lila Kari, Mehryar Mohri and Dirk Nowotka

Date: 20.02.2011-25.02.2011

Observations and Problems on k-abelian avoidability *

Mari Huova and Juhani Karhumäki

Department of Mathematics and TUCS University of Turku 20014 Turku, FINLAND email: {mari.huova, karhumak}@utu.fi

Theory of avoidability is among the oldest and most studied topic in Combinatorics on Words. The first result in this area, or in fact in the whole field, were obtained by Norwegian Axel Thue as early as at the beginning of 20th century [Th1, Th2]. He showed, among other things, the existence of an infinite binary word, which does not contain three consecutive factors of a word, that is a cube. Similarly, he showed that squares can be avoided in infinite ternary words

Since late 1960's commutative variants of the above problems were studied. Apparently, first nontrivial results were obtained by Evdokimov [Ev] who showed that commutative squares could be avoided in infinite words over a 25-letter alphabet. The size of the alphabet was reduced to five by Pleasant [Pl], until the optimal value four was found by Keränen [Ke], solving one celebrated problems of the topic. The optimal value for the size of the alphabet where abelian cubes were avoidable was proved earlier by Dekking [De], the value being three.

Interesting in all these results is that the required words are obtained by iterating a morphism.

We introduce in this note new variants of the problems by defining repetitions via new equivalence relations which lie properly in between equality and commutative equality, that is abelian equality.

Let $k \geq 1$ be a natural number. We say that words u and v in Σ^+ are k-abelian equivalent, in symbols $u \equiv_{a,k} v$, if

- 1. $\operatorname{pref}_{k-1}\left(u\right)=\operatorname{pref}_{k-1}\left(v\right)$ and $\operatorname{suf}_{k-1}\left(u\right)=\operatorname{suf}_{k-1}\left(v\right),$ and
- 2. for all $w \in \Sigma^k$, the number of occurrences of w in u and v coincide, i.e. #(w,u) = #(w,v).

^{*}Supported by the Academy of Finland under the grant 121419 and by the Väisälä Foundation.

Here $\operatorname{pref}_{k-1}$ (resp. \sup_{k-1}) is used to denote the prefixies (resp. suffixies) of length k-1 of words.

It is straightforward to see that $\equiv_{a,k}$ is an equivalence relation and, moreover,

$$u = v \Rightarrow u \equiv_{a,k} v \Rightarrow u \equiv_a v$$
,

where \equiv_a denotes the abelian equivalence, and that

$$u = v \Leftrightarrow u \equiv_{a,k} v \ \forall \ k \ge 1.$$

Now, notions like k-abelian repetitions are naturally defined. For instance, w = uv is a k-abelian square if and only if $u \equiv_{a,k} v$.

Natural variants of the above Thue's problems ask what are the smallest alphabets where k-abelian square and cubes can be avoided. A goal of this note is to point out that these problems are not trivial even in the case k = 2. Before going into that we make a few preliminary simple observations.

First, in the binary case 2- and 3-abelian words are fairly easy to characterize.

Example 1. In a binary alphabet $\Sigma = \{a, b\}$ the characterization of equivalence classes of 2-abelian words, via their representatives, can be given in the form:

$$aa^kb^l(ab)^ma^n$$
 or $bb^ka^l(ba)^mb^n$,

where $k, l, m \ge 0$ and $n \in \{0, 1\}$. And in the same alphabet the characterization of equivalence classes of 3-abelian words can be given in the following form containing eight possible combinations:

$$\left. \begin{array}{l} aaa^kb^l(aabb)^m \\ bbb^ka^l(aabb)^m \\ abb^ka^l(aabb)^m \\ baa^kb^l(aabb)^m \end{array} \right\} * \quad \text{connected with} \quad * \left\{ \begin{array}{l} (aab)^g(ab)^hb^ia^j \\ (abb)^g(ab)^hb^ia^j, \end{array} \right. \text{or}$$

where $k, l, m, g, h \ge 0$, $i \in \{0, 1\}$ and $j \in \{0, \dots, 2-i\}$. Here the characterizations are not unique in few cases.

The above allows to estimate the sizes of the corresponding equivalence classes. They are of order $\Theta(n^2)$ and $\Theta(n^4)$, see [HKSS]. Recently, A. Saarela [Sa] showed that in general the number of k-abelian equivalence classes of words of length n is polynomial in n but the degree of the polynomial grows exponentially in k (in a fixed but arbitrary alphabet).

Our next example shows that the ordinary method of iterating a morphism might not give answers to our problems.

Example 2. In all of the following cases where repetition free infinite word is obtained by iterating a morphism, a 2-abelian cube is found fairly early from the beginning. For overlap- and cube-free words see [AJ].

- Infinite overlap-free Thue-Morse word (morphism: $0 \rightarrow 01, 1 \rightarrow 10$): $01\overline{101001}\overline{100101}\overline{101001}011...$
- Cube-free infinite word (morphism: 0 \rightarrow 001, 1 \rightarrow 011): 001001 011001001011 001011 011...

- Morphism $0 \to 001011$, $1 \to 001101$, $2 \to 011001$ maps ternary cube-free words to binary cube-free words, see [Br], but $001011 \equiv_{a,2} 001101 \equiv_{a,2} 011001$, thus images of all words mapped with this morphism contains 2-abelian cubes.
- A binary overlap-free word w can also be gained in form $w = c_0 c_1 c_2 \dots$, where c_n means the number of zeros (mod 2) in the binary expansion of n. Again, a 2-abelian cube of length 6 begins as early as from the fifth letter: $w = 0010 \ \overline{011010} \ \overline{010110} \ \overline{011010} \ 011...$

In order to go into our problems we recall the following Table 1 which summarizes the results we mentioned at the beginning and at the same time tells the limits of our problems:

Avoidability of squares				Avoidability of cubes			
	type of rep.				type of rep.		
size of the alph.	=	$\equiv_{a,2}$	\equiv_a	size of the alph.	=	$\equiv_{a,2}$	\equiv_a
2	_	_	_	2	+	?	_
3	+	?	_	3	+	+	+
4	+	+	+				

Table 1: Avoidability of different types of repetitions in infinite words.

We were able to settle the first one of the above question marks by computer checking.

Example 3. The longest ternary word which is 2-abelian square-free has length 537, which shows that there does not exist an infinite 2-abelian square-free word over any ternary alphabet. This longest word, given below, is unique up to the permutations of the alphabet, $\Sigma = \{a, b, c\}$.

Example 4. We can also construct the whole tree containing each ternary 2-abelian square-free word and analyze the sizes of the sets containing words from length 1 to 537, respectively. There exist 404 286 words of length 105 and for other lengths the number of ternary 2-abelian square-free words is less. The number of words grows monotonically from 3 up to 403 344 when considering lengths from 1 to 103. After the length 105 there appears more oscillation between the numbers of words of different lengths. The sizes of the sets containing ternary 2-abelian square-free words are shown in Figure 1 with respect to the lengths of words.

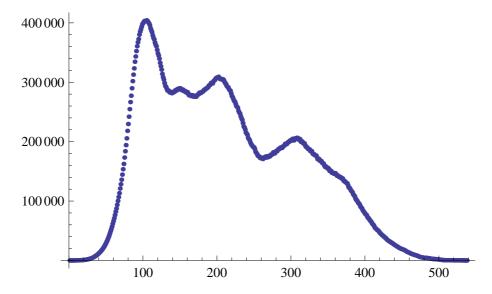


Figure 1: The number of 2-abelian square-free words with respect to their lengths.

To solve the other question mark we also did some computer checking - and obtained evidence that the answer is likely to be different compared to the first one.

Example 5. With a computer we were able to construct a binary word of more than 100 000 letters that still avoids 2-abelian cubes. This shows that there exist, at least, very long binary 2-abelian cube-free words.

Example 6. Similarly, we can examine the number of binary 2-abelian cubefree words of given length as in the previous case concerning ternary 2-abelian square-free words. The numbers of the words with lengths from 1 to 60 grow approximately with a factor 1,3 at each increment of the length, see Figure 2. So that the number of binary 2-abelian cube-free words of length 60 is already 478 456 030. And already, with length 12 there exist more binary 2-abelian cube-free words (254) than ternary 2-abelian square-free words (240).

We also chose some binary 2-abelian cube-free prefixies and counted the numbers of binary 2-abelian cube-free words having these fixed preixies. In this way we can check how many suitable extensions the chosen 2-abelian cube-free word has. As a result we found examples of binary 2-abelian cube-free words with a property that the number of their extensions grows again approximately with a factor 1,3 when increasing the length of extensions by one. In Figure 3 this is done fore a fixed prefix of length 2000.

These examples support the conjecture that there would exist an infinite binary word that avoids 2-abelian cubes. As a conclusion, our two considered problems would behave differently: one like words and the other like abelian words.

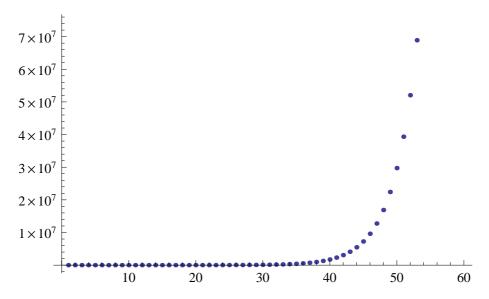


Figure 2: The number of 2-abelian cube-free words with respect to their lengths for small values of length.

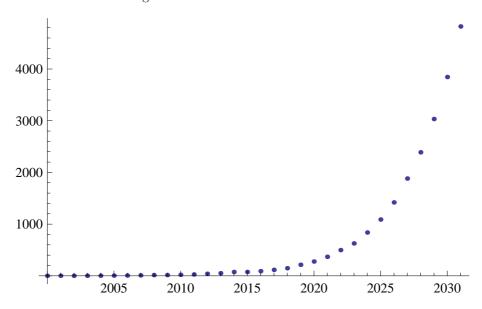


Figure 3: The number of 2-abelian cube-free words with respect to their lengths from 2000 to 2031 and with a common prefix of length 2000.

References

- [AJ] J.-P. Allouche, J. Shallit: Automatic Sequences: Theory, Applications, Generalizations. Cambridge University Press, Cambridge, 2003.
- [Br] F.-J. Brandenburg: Uniformly growing k-th power-free homomorphisms. Theoret. Comput. Sci. 23, 69-82 (1983).

- [De] F. M. Dekking: Strongly non-repetitive sequences and progression-free sets. J. Combin. Theory Ser. A 27(2), 181-185 (1979).
- [Ev] A. A. Evdokimov: Strongly asymmetric sequences generated by a finite number of symbols. Dokl. Akad. Nauk SSSR 179, 1268-1271 (1968); English translation in Soviet Math. Dokl. 9, 536-539 (1968).
- [HKSS] M. Huova, J. Karhumäki, A. Saarela, K. Saari: Local squares, periodicity and finite automata. LNCS Festschrift for Hermann Maurer, Springer, (to appear).
- [Ke] V. Keränen: Abelian squares are avoidable on 4 letters. In: W. Kuich (ed.) ICALP 1992. LNCS, vol. 623, 41-52. Springer, Heidelberg, 1992.
- [Pl] P. A. B. Pleasant: Non-repetitive sequences. Proc. Cambridge Philos. Soc. 68, 267-274 (1970).
- [Sa] A. Saarela: Private communication.
- [Th1] A. Thue: Über unendliche Zeichenreihen. Norske vid. Selsk. Skr. Mat. Nat. Kl. 7, 1-22 (1906).
- [Th2] A. Thue: Über die gegenseitige Lage gleicher Teile gewisser Zeichenreihen. Norske vid. Selsk. Skr. Mat. Nat. Kl. 1, 1-67 (1912).