

Square-free words

David Stanovský

Abstract. Construction of an infinite square-free word via the word of Thue-Morse is presented. Some additional properties of these words are mentioned.

Introduction

Combinatorics on words is engaged in looking for regularities in words. For instance van der Waerden's Theorem shows, that every sufficiently long word possesses prescribed arithmetic progression of one letter. We will introduce another situation – not every (long) word contains a square, although it seems on the contrary (try to find any!).

No knowledge of combinatorics on words is supposed. All definitions are included, almost all techniques are elementary. More advanced parts of the text are typed by small font and it is not necessary to read and understand them.

Let \mathbb{N} denote the set of all natural numbers $\{0, 1, \dots\}$. By an *alphabet* we mean a finite nonempty set, its elements are called *letters*. A *word over alphabet* A is a finite sequence of letters from A . Empty word (sequence of length 0) is denoted ε . The set of all words over an alphabet A we denote A^* and $A^+ = A^* \setminus \{\varepsilon\}$.

Concatenation of words u and v is denoted by uv . A *morphism* between A^* and B^* is a map $f : A^* \rightarrow B^*$ such that $f(\varepsilon) = \varepsilon$ and $f(uv) = f(u)f(v)$ for every $u, v \in A^*$. A word u is called a *factor* of v , if there exist words x, y such that $v = xuy$. The word u is called a *left factor*, if $x = \varepsilon$. If u is a word, than $|u|$ means length of the word and u^R is a word read in an opposite direction. *Palindrome* is such word, that $u = u^R$.

$(A^*, \text{concatenation}, \varepsilon)$ forms a free monoid. A morphism is a monoid homomorphism.

A *square* is a word of the form uu , where u is some nonempty word. Word contains a square, if one of its factors is square. Otherwise we call the word *square-free*. E.g. $abcacbacbc$ contains the square $acbacb$, but $abcacbabcb$ is square-free (as will be shown later).

We will construct an infinite square-free word over an alphabet with three letters. Clearly, then there exist infinitely many finite square-free words. There exists no square-free word over two-letter alphabet of the length more than 3 (the only ones are a, b, ab, ba, aba, bab). The infinite square-free word will be derived from the so-called word of Thue-Morse, which contains no factor of the form $avava$, where a is a letter and v is a word.

Axel Thue, Norwegian mathematician, was first who was interested in this topic. He constructed the same words as we will in his papers written in 1906 and 1912. This was independently described and improved by M. Morse in 1921. Then many other papers were written on related topics.

The fact that there exist infinite square-free words is equivalent to the fact that A^0/\equiv is infinite, where A^0 denotes the monoid A^* with adjoined zero (i.e. $a0 = 0a = 00 = 0$ for every $a \in A^*$ and \equiv is a congruence generated by $uu \equiv 0$ for every $u \in A^+$).

Preliminaries

Let A denote an alphabet. Let $u, v \in A^+$, u occurs at least twice in v . Then there exist $x, y, x', y' \in A^*$ such that $|x| < |x'|$, $|y| > |y'|$ and $v = xuy = x'u'y'$. Occurrences of u in v are called

- (1) *disjoint*, if $|x'| > |xu|$, i.e. $v = xuzuy'$ for some z .
- (2) *adjacent*, if $|x'| = |xu|$, i.e. $v = xuyu'$.
- (3) *overlapping*, if $|x'| < |xu|$.

A good description of the third possibility is provided by the following lemma. By an *overlapping factor* we mean a factor of the form $avava$, where $a \in A, v \in A^*$.

Lemma. *A word $w \in A^*$ contains two overlapping occurrences of some nonempty word, iff it contains some overlapping factor.*

Proof.

I. Let $w = xuy = x'uy'$ such that $0 \leq |x| < |x'| < |xu| < |x'u| \leq |w|$. Then $x' = xs, xu = x'z, x'u = xut$ for some nonempty words s, z, t . Then $(*) u = sz = zt$. Denote a the first letter of s , i.e. of z too (by $(*)$). So $s = as', z = az'$. Thus $u = sz = as'az'$, so $w = x'uy' = xsuy' = xas'as'az'y'$ and clearly $as'as'a$ is an overlapping factor in w .

II. If $w = xavavay$, then ava has an overlapping occurrence in w . □

By *word* we mean always finite word. Now we will define an *infinite word*. It is an infinite sequence of letters, i.e. a function $\mathbf{a} : \mathbb{N} \rightarrow A$, denoted by $\mathbf{a} = a_0a_1a_2\dots$, where $a_i = \mathbf{a}(i)$ for each $i \in \mathbb{N}$. Let us define $\mathbf{a}^{[k]} = a_0\dots a_{k-1}$ and call it a *left factor* of \mathbf{a} of the length k . If $u = \mathbf{a}^{[k]}$, then we shall write $\mathbf{a} = u\mathbf{b}$, where \mathbf{b} is such that $b_i = a_{i+|u|}$, $i \in \mathbb{N}$. A word u we call a *factor* of \mathbf{a} , if $\mathbf{a} = xub$ for some x and \mathbf{b} .

Infinite words are useful for description of properties of finite words which are *stable for factors*. It means that if some word possesses this property, then so do all its factors. Clearly square-freeness is stable for factors.

We say, that an infinite word \mathbf{a} has a property P , if all its factors do so. This clarifies the sense of the term "infinite square-free word".

Let us denote L_P the set of all words with the property P . Thus, if P is stable for factors and $w \in L_P$, then all factors of w are in L_P .

Lemma. *Let P be a property of words over A stable for factors. Then L_P is infinite, iff there exist an infinite word over A having the property P .*

Proof.

I. Suppose L_P infinite and A finite. There must exist some $a_0 \in A$ such that infinitely many words from L_P start with a_0 . Let us denote $L_0 = \{b \in L_P : b = a_0y \text{ for some } y \in A^*\}$. The same argument allows us to construct by induction sets L_1, L_2, \dots of words starting by $a_0a_1, a_0a_1a_2, \dots$. Letters a_0, a_1, \dots form an infinite word $\mathbf{a} = a_0a_1\dots$ with the property P .

II. Converse direction is quite clear. If \mathbf{a} is an infinite word with the property P , then for every i natural $\mathbf{a}^{[k]} \in L_P$, so L_P is infinite. □

The proof shows us an algorithm for derivation of the infinite word with P from infinitely many finite words possessing P .

There is one more interesting fact about properties of words.

Proposition. *Let P be a property of words over A such that L_P is a two-sided ideal in A^* (i.e. if $v \in L_P$, then every word over A containing v as a factor is in L_P). Then every infinite word over A has a factor in L_P , iff $A^* \setminus L_P$ is finite.*

Proof. If L_P is two-sided ideal in A^* , then $\neg P$ is stable for factors. By the previous lemma, $L_{\neg P}$ is infinite, iff there exist an infinite word with $\neg P$. Note that $L_{\neg P} = A^* \setminus L_P$ and the existence of an infinite word without P means, that $(\exists \mathbf{a})(\forall u \text{ factor of } \mathbf{a}) u \text{ possesses } \neg P$. So $A^* \setminus L_P$ is finite, iff for every infinite word \mathbf{a} there exist a factor u of \mathbf{a} such that u has P . So $u \in L_P$ and the Proposition holds. □

Let us consider a sequence w_0, w_1, \dots of words over A such that w_n is a left factor of w_{n+1} for all n natural. Denote \mathbf{a} an infinite word satisfying $\mathbf{a}^{[k]} = w_n$ for all $k = |w_n|, n \in \mathbb{N}$. We write $\mathbf{a} = \lim w_n$ and call it a *limit* of sequence $(w_n)_{n=0}^\infty$.

Imagine this special case. Let $\alpha : A^* \rightarrow A^*$ be a morphism satisfying $\alpha(a) \neq \varepsilon$ for all $a \in A$ and $\exists a_0 \in A$ such that $\alpha(a_0) = a_0u$ for some $u \in A^+$ (we say that α satisfies (\heartsuit) for a_0). Thus for every n natural $\alpha^{n+1}(a_0) = \alpha^n(\alpha(a_0)) = \alpha^n(a_0u) = \alpha^n(a_0)\alpha^n(u)$, so $\alpha^n(a_0)$ is a left factor of $\alpha^{n+1}(a_0)$. The limit of this sequence is called *limit of iterating α on a_0* and it is denoted $\alpha^\omega(a_0)$.

There is natural extension of a morphism $\alpha : A^* \rightarrow A^*$ to infinite words over A . For $\mathbf{b} = b_0b_1\dots$ is $\alpha(\mathbf{b}) = \alpha(b_0)\alpha(b_1)\dots$ — it is an infinite word because of the first condition in (\heartsuit) .

Lemma. *Let α satisfies (\heartsuit) for a_0 and $\mathbf{a} = \alpha^\omega(a_0)$. Then $\alpha(\mathbf{a}) = \mathbf{a}$.*

Proof. If u is a left factor of $\alpha(\mathbf{a})$, then so is $\alpha(u)$. Thus every $\alpha^n(a_0)$ is a left factor of $\alpha(\mathbf{a})$. But $\alpha(\mathbf{a})$ starts with a_0 (the second condition in (\heartsuit)), so $\alpha(\mathbf{a}) = \lim \alpha^n(a_0) = \alpha^\omega(a_0) = \mathbf{a}$. \square

Words od Thue-Morse

Let $A = \{a, b\}$ in the rest of the paper. For every $w \in A^*$ we denote \bar{w} the word obtained from w by replacing a to b and vice versa.

Let $\mu : A^* \rightarrow A^*$ is a morphism defined by $\mu(a) = ab$ and $\mu(b) = ba$. Clearly μ satisfies (\heartsuit) for a and b . Denote

$$\mathbf{t} = \mu^\omega(a) = abbabaabbaababbabaababbaabbabaab\dots$$

$$\bar{\mathbf{t}} = \mu^\omega(b) = baababbaababbabaababbabaabbaababba\dots$$

Lemma. *Let $u_0 = a, v_0 = b, u_{n+1} = u_nv_n, v_{n+1} = v_nu_n$ for all natural n . Then for all $n \in \mathbb{N}$ hold*

- (1) $u_n = \mu^n(a), v_n = \mu^n(b)$.
- (2) $v_n = \overline{u_n}, u_n = \overline{v_n}$.
- (3) u_{2n}, v_{2n} are palindromes, $u_{2n+1}^R = v_{2n+1}$

Proof. By induction on n . The case $n = 0$ is clear.

- (1) $u_{n+1} = u_nv_n = \mu^n(a)\mu^n(b) = \mu^n(ab) = \mu^n(\mu(a)) = \mu^{n+1}(a)$. The rest is similar.
- (2) $v_{n+1} = v_nu_n = \overline{u_nv_n} = \overline{u_nv_n} = \overline{u_{n+1}}$. The rest is similar.
- (3) $u_{2n+2} = u_{2n+1}v_{2n+1} = u_{2n+1}u_{2n+1}^R$ which is a palindrome. For v_n similar. $u_{2n+1}^R = (u_{2n}v_{2n})^R = v_{2n}^R u_{2n}^R = v_{2n}u_{2n} = v_{2n+1}$. \square

There is a lot of equivalent definitions of $\mathbf{t} = t_0t_1\dots$. For instance this one. Let $d(n)$ denote the number of ones in the binary expansion of n

Proposition. *For all n natural is t_n either a (if $d(n)$ is even) or b (if $d(n)$ is odd).*

Proof. By the last lemma in the previous section $\mathbf{t} = \mu(\mathbf{t}) = \mu(t_0)\mu(t_1)\dots$. Hence $\mu(t_n) = t_{2n}t_{2n+1}$, because $|\mu(a)| = |\mu(b)| = 2$. So (look at the definition of μ) $t_{2n} = t_n$ and $t_{2n+1} = \overline{t_n}$ holds.

The proposition we prove by induction. For $n = 0$ it holds. Now for all $k < n$ proposition holds. At first let $n = 2m$. Then $d(n) = d(m)$ and $t_n = t_{2m} = t_m$ and proposition holds for n too. Now let $n = 2m + 1$. Then $d(n) = d(m) + 1$ and $t_n = t_{2m+1} = \overline{t_m}$, so it holds too. \square

Now we will prove, that \mathbf{t} contains no overlapping factor. Then \mathbf{t} is also cube-free, because if uuu is a factor of \mathbf{t} , $u = au'$ for $a \in A$, then $au'au'a$ is an overlapping factor of t providing contradiction.

We will need two lemmas.

Lemma 1. *If $X = \{ab, ba\}$, $x \in X^*$, then $axa \notin X^*$, $bx b \notin X^*$.*

Proof. Let $x \in X^*$. We use induction on $|x|$. For $|x| = 0$ is $aa, bb \notin X^*$. Now let x satisfies $axa \in X^*$ and for all shorter words proposition holds. Let us write $axa = u_0\dots u_k$, $u_i \in X$. Thus must be $u_0 = ab$ and $u_k = ba$. So $u = u_1\dots u_{k-1} \in X^*$, u is shorter then x and $bub = x \in X^*$. That is contradiction with an induction assumption. For $bx b$ similarly. \square

Lemma 2. *If $w \in A^+$ contains no overlapping factor, then neither does $\mu(w)$.*

Proof. Suppose $\mu(w)$ contains an overlapping factor. Then $\mu(w) = xcvcvcy$ for some $x, v, y \in A^*$, $c \in A$. Note, that $\mu(w) \in X^*$ for $X = \{ab, ba\}$ and thus $|\mu(w)|$ is even. But $|cvcvc|$ is odd, so either

- (1) $|x|$ is odd, $|y|$ is even and thus $xc, vcvc, y \in X^*$, or
- (2) $|x|$ is even, $|y|$ is odd and thus $x, cvcv, cy \in X^*$.

In both cases is $|v|$ odd (if it is even, then $cvc \in X^*$, $v \in X^*$ contradicting Lemma 1). So either $vc \in X^*$ or $cv \in X^*$.

- (1) We can write $w = rsst$ so that $\mu(r) = xc$, $\mu(s) = vc$, $\mu(t) = y$. Words r, s finish by the same letter \bar{c} , so $r = r'\bar{c}$, $s = s'\bar{c}$. Thus $w = r'\bar{c}s'\bar{c}s't$ contains an overlapping factor, contradiction.
- (2) We can write $w = rsst$ so that $\mu(r) = x$, $\mu(s) = cv$, $\mu(t) = cy$. Words s, t start by the same letter c , so $s = cs'$, $t = ct'$. Thus $w = rcs'cs'ct'$ contains an overlapping factor, contradiction. \square

Theorem. *An infinite word \mathbf{t} contains no overlapping factor.*

Proof. Let x is an overlapping factor in \mathbf{t} . There must exist (sufficiently large) k such that $\mu^k(a)$ has an overlapping factor x . But a does'nt have any overlapping factor, so by lemma 2 also $\mu(a)$ does'nt, so also $\mu^2(a)$, etc., also $\mu^k(a)$ has no overlapping factor. That is contradiction. \square

Two interesting properties of \mathbf{t} are presented.

Proposition. *If $h : A^* \rightarrow A^*$ is a morphism such that $\mathbf{t} = h(\mathbf{t})$, then $h = \mu^n$ for some $n \in \mathbb{N}$.*

Proof. One can prove the following technical lemma (proof is not included).

Lemma. *If wuu is a left factor of \mathbf{t} , then there exist $n \in \mathbb{N}$ such that $|u|$ is 2^n or $3 \cdot 2^n$ and $|w|$ is an integer multiple of 2^n .*

The word abb is a left factor of \mathbf{t} . Hence $h(abb) = h(a)h(b)h(b)$ is a left factor of $h(\mathbf{t}) = \mathbf{t}$, so by the lemma $|h(a)| = k2^r$ and $|h(b)| = l2^r$ for some k, r natural, $l \in \{1, 3\}$.

The word $abbabaa$ is a left factor of \mathbf{t} . Hence $h(abbabaa) = h(abbab)h(a)h(a)$ is a left factor of \mathbf{t} , so by the lemma $|h(abbab)| = 3|h(b)| + 2|h(a)| = i2^s$ and $|h(a)| = j2^s$ for some i, s natural, $j \in \{1, 3\}$.

Let us make a simple computation. $|h(a)| = k2^r = j|2^s|$, so $k2^{r-s} = j$. The left side is even for $r > s$, the right side is odd, thus $r \leq s$. $3|h(b)| + 2|h(a)| = i2^s = j2^s + l2^r$, so $(i-j)2^{s-r} = l$. The left side is even for $r < s$, the right side is odd, thus $s \leq r$. Hence $r = s$.

Now the word $w = abbabaabbaababbab$ is also a left factor of \mathbf{t} . Hence $h(w) = h(abbabaabbaa)h(bab)h(bab)$ is a left factor of \mathbf{t} , so by the lemma $|h(bab)| = m2^t$ for some t natural, $m \in \{1, 3\}$. But $m2^t = |h(bab)| = |h(a)| + 2|h(b)| = (j+2l)2^r$, so $2l+j = m2^{t-r}$. The left side is odd. The right side is for $t > r$ even, for $t < r$ not integer (m is not multiple of 2), so the only possibility is $t = r$. Then $2l+j = m$ and $l, j, m \in \{1, 3\}$. The only solution is $l = j = 1, m = 3$ and thus $|h(a)| = |h(b)| = 2^r$. Because $h(a)h(b)$ is a left factor in \mathbf{t} of the length 2^{r+1} , the identity $h = \mu^r$ must hold. \square

The formal power series over \mathbb{Z}_2 form a ring $\mathbb{Z}_2[[x]]$. Let $y \in \mathbb{Z}_2[[x]]$ be the power series $\sum_{i=0}^{\infty} a_i x^i$, where $a_i = 0$ if $t_i = a$ and $a_i = 1$ if $t_i = b$. Let $\bar{y} \in \mathbb{Z}_2[[x]]$ be the power series $\sum_{i=0}^{\infty} b_i x^i$, where $b_i = 0$ if $t_i = b$ and $b_i = 1$ if $t_i = a$.

Proposition. *Formal power series y and \bar{y} are solutions of the equation $(1+x)^3 z^2 + (a+x)^2 z + x = 0$ in $\mathbb{Z}_2[[x]]$.*

Proof. Let $z = \sum_{i=0}^{\infty} z_i x^i$ is a solution of the given equation in $\mathbb{Z}_2[[x]]$. Denote $s_i = \sum_{j=0}^i a_j a_{i-j}$. Then $z^2 = \sum_{i=0}^{\infty} s_i x^i$. Hence z is a solution, iff

$$\sum_{i=0}^{\infty} s_i x^i + \sum_{i=1}^{\infty} s_{i-1} x^i + \sum_{i=2}^{\infty} s_{i-2} x^i + \sum_{i=3}^{\infty} s_{i-3} x^i + \sum_{i=0}^{\infty} z_i x^i + \sum_{i=2}^{\infty} z_{i-2} x^i + x = 0,$$

iff

$$s_0 + a_0 = 0, \quad s_0 + s_1 + z_1 + 1 = 0, \quad s_0 + s_1 + s_2 + z_2 + z_0 = 0$$

$$s_i + s_{i-1} + s_{i-2} + s_{i-3} + a_i + a_{i-2} = 0 \quad (i \geq 3).$$

One can see that for every k natural $s_{2k} = z_k$ and $s_{2k+1} = 0$ (because $yy = y$ and $y + y = 0$ in \mathbb{Z}_2). Thus first and third equation hold for every $z \in \mathbb{Z}_2[[x]]$.

The proposition about an equivalent definition of \mathbf{t} gives the following observation. If a_n are coefficients of y (i.e. $-a_n$ are coefficients of \bar{y}), then $a_n = d(n) \bmod 2$. Because $d(2n) = d(n)$, (i) $a_{2n} = a_n$ holds for every n natural. Similary, $d(2n+1) = d(n) + 1$, so (ii) $a_{2n+1} = -a_n$.

If $z_i = \pm a_i$, second equation holds ($z_0 + z_1 = 0$, because z_0, z_1 are different). For $i \geq 3$ odd we have $i = 2j + 1$, and thus $a_j + a_{j-1} + a_{2j-1} + a_{2j+1} = 0$ using (ii). For $i \geq 3$ even we have $i = 2j$, and thus $a_j + a_{j-1} + a_{2j} + a_{2j-2} = 0$ using (i). (For $z_i = -a_i$ similarly.) \square

Square-free words

We know, that no square-free words longer then 3 occur over two-letter alphabet. So let $B = \{a, b, c\}$ and

$$\delta : B^* \rightarrow A^*, \quad \delta(a) = abb, \quad \delta(b) = ab, \quad \delta(c) = a.$$

If \mathbf{a} is an infinite word without overlapping factors starting with letter a , then there is a unique factorization $\mathbf{a} = y_0 y_1 \dots$, where $y_n \in \{a, ab, abb\} = \delta(B)$ for all $n \in \mathbb{N}$. It is true, because every a in \mathbf{a} is followed by at most two letters b and then again by a (\mathbf{a} is cube-free). So with every a starts some $w_n \in \delta(B)$ of the length (number of b)+1. Thus it is clear, that there exist a unique infinite word \mathbf{b} over B such that $\delta(\mathbf{b}) = \mathbf{a}$.

Theorem. *If \mathbf{a} is an infinite word over A starting with a without overlapping factors and \mathbf{b} is such that $\delta(\mathbf{b}) = \mathbf{a}$, then \mathbf{b} is square-free.*

Proof. Let \mathbf{b} contains a square uu , denote d next letter after one of its occurrences (i.e. $\mathbf{b} = xuudc$ for some x, c). Hence $\delta(uud)$ is a factor of \mathbf{a} . Denote v, w words satisfying $\delta(u) = av$, $\delta(d) = aw$. Then $\delta(uud) = \delta(u)\delta(u)\delta(d) = avavaw$ is a factor of \mathbf{a} , so \mathbf{a} contains an overlapping factor. Contradiction. \square

Let us denote \mathbf{m} the square-free word obtained from \mathbf{t} (i.e. $\delta(\mathbf{m}) = \mathbf{t}$). One can check, that

$$\mathbf{m} = abcacbabcbacabcacbacabcabcbacabcbacabcbac \dots$$

The theorem does not hold conversely – there exists a square-free infinite word \mathbf{b} over B , such that $\delta(\mathbf{b})$ has an overlapping factor.

There exists a lot of another constructions of square-free words. E.g. $\mathbf{m} = \varphi^\omega(a)$ for a morphism defined by

$$\varphi : B^* \rightarrow B^*, \quad a \mapsto abc, \quad b \mapsto ac, \quad c \mapsto b$$

We can also systematically generate finite square-free words by so called *square-free morphisms*. That are such morphisms $\alpha : A^* \rightarrow B^*$ that satisfy $\alpha(A) \neq \{\varepsilon\}$ and for every square-free word w is $\alpha(w)$ square-free. E.g. a morphism

$$\varphi : B^* \rightarrow B^*, \quad a \mapsto abcab, \quad b \mapsto acacb, \quad c \mapsto acbcab$$

is square-free. An important theorem due to Bean, Ehrenfeucht and McNulty (1979) describes square-free morphisms.

Theorem. *Let $\alpha : A^* \rightarrow B^*$ be a morphism satisfying*

- (1) $\alpha(A) \neq \{\varepsilon\}$,
- (2) for every square-free word w of the length at most 3 is $\alpha(w)$ square-free,
- (3) for all $a, b \in A$ no $\alpha(a)$ is a proper factor of $\alpha(b)$.

Then α is a square-free morphism.

Literature

The paper was written using the book M. Lothaire, *Combinatorics on Words*, Cambridge University Press, 1983, 1997, chapter 2.