

Crochemore factorization of infinite words

Jean Berstel, Alessandra Savelli

Basic definition

Crochemore factorization: the *c-factorization* $c(x)$ of a word x is

$$c(x) = (x_1, x_2, \dots, x_m, x_{m+1}, \dots)$$

where x_m is the longest prefix of $x_mx_{m+1} \dots$ occurring twice in $x_1x_2 \dots x_m$, or x_m is a letter a if a does not occur in $x_1 \dots x_{m-1}$.

The *c-factorization* of $x = ababaab$ is (a, b, aba, ab) , since *aba* occurs twice in *ababa*.

Fibonacci word

The **Fibonacci word** is defined as the limit of the sequence $f_{-1} = b$, $f_0 = a$, and $f_{n+2} = f_{n+1}f_n$:

$$f_0 = a$$

$$f_1 = ab$$

$$f_2 = aba$$

$$f_3 = abaab$$

...

$$\mathbf{f} = abaababaaba \cdots = \tilde{f}_0 \tilde{f}_1 \tilde{f}_2 \tilde{f}_3 \cdots$$

The c -factorization of \mathbf{f} is exactly:

$$c(\mathbf{f}) = (a, b, a, aba, baaba, \dots) = (a, b, a, \tilde{f}_2, \tilde{f}_3, \dots).$$

Fibonacci word

The c -factorization is closely related to two other factorizations of the Fibonacci word.

$$h(\mathbf{f}) = (a, b, a, ab, aba, abaab, \dots)$$

$$w(\mathbf{f}) = (a, b, aa, bab, aabaa, \dots)$$

The three factorizations can be visualized through the following scheme:

$h :$	a	b	a	a	b	a	b	a	a	b	\dots	
$w :$	a	b	a	a	b	a	b	a	a	b	a	\dots
$c :$	a	b	a	a	b	a	b	a	a	b	a	\dots

Sturmian words

A **standard Sturmian word** is defined as the limit of

$$s_{-1} = b, s_0 = a, \text{ and } s_n = s_{n-1}^{d_n} s_{n-2},$$

where d_i is a positive integer for all $i > 0$.

Similarly to the Fibonacci word, the Sturmian words have a decomposition in reverse finite words s_n :

$$s = \tilde{s}_0^{d_1} \tilde{s}_1^{d_2} \tilde{s}_2^{d_3} \dots$$

The c -factorization of standard Sturmian words is closely related to that decomposition:

$$c(s) = (a, a^{d_1-1}, b, a^{d_1} \tilde{s}_1^{d_2-1}, \tilde{s}_2^{d_3}, \tilde{s}_3^{d_4}, \dots, \tilde{s}_n^{d_{n+1}}, \dots).$$

Thue-Morse word

Let τ be the Thue-Morse morphism on a two-letter alphabet defined by $\tau(a) = ab$ and $\tau(b) = ba$, the **Thue-Morse infinite word** $\mathbf{t} = abbabaabbaababba \dots$ is defined as the limit of the sequence

$$t_0 = a, \quad t_n = \tau(t_{n-1}).$$

Each factor in the c -factorization of \mathbf{t} can be obtained from the previous ones by applying the morphism τ :

$$c(\mathbf{t}) = (c_1, c_2, \dots), \quad c_{n+2} = \tau(c_n) \text{ for every } n > 6.$$

$$c(\mathbf{t}) = (a, b, b, ab, a, abba, \textcolor{red}{aba}, \textcolor{green}{bbabaab}, \textcolor{red}{abbaab}, \textcolor{green}{babaabbaababba}, \dots)$$

Thue-Morse generalized words

The **Thue-Morse generalized word** on a m -letter alphabet $A = \{a_1, a_2, \dots, a_m\}$ obtained as the limit of the sequence

$$t_0^{(m)} = a_1, \quad t_{n+1}^{(m)} = \tau_m(t_n^{(m)}),$$

where τ_m is the morphism defined by

$$\tau_m(a_i) = a_i a_{i+1} \cdots a_m a_1 \cdots a_{i-1} \quad (i = 1, \dots, m).$$

$$c_{n+2(m-1)}^{(m)} = \tau_m(c_n) \text{ for every } n > m \text{ and } m \geq 3.$$

Example for $m = 3$:

$$c(\mathbf{t}^{(3)}) = (a, b, c, bc, a, \textcolor{red}{ca}, b, bcacab, abc, \textcolor{red}{cababc}, bca, \dots)$$

Period doubling sequence

Let δ be the morphism on a two-letter alphabet defined by

$$\delta(a) = ab, \quad \delta(b) = aa.$$

The **period doubling sequence** is the limit of the sequence

$$q_0 = a \text{ and } q_{n+1} = \delta(q_n).$$

Period doubling sequence

Similarly to the case of standard Sturmian words, the period doubling sequence admits the decomposition:

$$\mathbf{q} = \tilde{q}_0 \tilde{q}_1 \tilde{q}_2 \cdots ,$$

and the c -factorization of \mathbf{q} is

$$c(\mathbf{q}) = (a, b, a, aa, ba, baba, aaba, \dots) = (\tilde{q}_0, \tilde{q}_0', \tilde{q}_0, \tilde{q}_1', \tilde{q}_1, \tilde{q}_2', \tilde{q}_2, \dots),$$

where q'_n is q_n with just the last letter changed to its opposite.

Note that $q_{n+1} = q_n q'_n$.

Crochemore vs. Ziv-Lempel

Ziv-Lempel factorization: the *z-factorization* $z(x)$ of a word x is

$$z(x) = (y_1, y_2, \dots, y_m, y_{m+1}, \dots)$$

where y_m is the shortest prefix of $y_m y_{m+1} \dots$ which occurs only once in $y_1 y_2 \dots y_m$.

The *c-factorization* of $x = ababaab$ is $(a, b, abaa, b)$.

Crochemore vs. Ziv-Lempel

Facts:

- A Crochemore factor cannot properly include a Ziv-Lempel factor.
- If a Ziv-Lempel factor includes a Crochemore factor, then it ends at most a letter after.

For example, let x be the word $x = aabaaccbaabaabaa$. The c -factorization and the z -factorization of x are:

$$c(x) = (a, a, b, aa, c, c, baa, baabaa)$$

$$z(x) = (a, ab, aac, cb, aabaab, aa).$$

Crochemore vs. Ziv-Lempel

Consider for example the period doubling sequence

$$q = abaaabababaaabaa \dots$$

Each Ziv-Lempel factor of q properly includes a Crochemore factor by ending just a letter before:

$z :$	a	b	a	a	a	b	a	b	a	b	a	a	a	b	a	a	\dots
$c :$	a	b	a	a	a	b	a	b	a	b	a	a	a	b	a	a	\dots