

## SUBSTITUTIONS ON AN INFINITE ALPHABET: FIRST RESULTS

SÉBASTIEN FERENCZI

ABSTRACT. We give a few examples of substitutions on an infinite alphabet, and the beginning of a general theory of the associated dynamical systems.

### 1. SUBSTITUTIONS

Let  $A$  be a finite or countable set, called the **alphabet**, and its elements will be called **letters**.

**Definition 1.1.** A **word** is a finite string  $w_1\dots w_k$  of elements of  $A$ ; the concatenation of two words  $w$  and  $w'$  is denoted multiplicatively, by  $ww'$ . A word  $w_1\dots w_k$  is said to **occur** at place  $i$  in the infinite sequence or finite word  $u$  if  $u_i = w_1, \dots, u_{i+k-1} = w_k$ ; when  $u$  is finite, we denote by  $N(w, u)$  the number of these occurrences.

A **substitution** is an application from an alphabet  $A$  into the set  $A^*$  of finite words on  $A$ ; it extends to a morphism of  $A^*$  for the concatenation by  $\sigma(ww') = \sigma w \sigma w'$ .

It is called **primitive** if there exists  $k$  such that  $a$  occurs in  $\sigma^k b$  for any  $a \in A$ ,  $b \in A$ .

It is called **of constant length**  $q$  if  $\sigma a$  is of length  $q$  for any  $a \in A$ .

A **fixed point** of  $\sigma$  is an infinite sequence  $u$  with  $\sigma u = u$ .

For any sequence  $u = (u_n, n \in \mathbb{N})$  on a finite alphabet  $A$ , we can define the (topological) **symbolic dynamical system** associated to  $u$ : we first take  $\Omega = A^{\mathbb{N}}$ , equipped with the product topology (each copy of  $A$  being equipped with the discrete topology) and  $T$  the one-sided shift

$$T(x_0x_1x_2\dots) = x_1x_2x_3\dots$$

then  $X_u$  is the closure of the orbit of  $u$  under  $T$ . The dynamical system associated to a primitive substitution is the symbolic system  $(X_u, T)$  associated to any of its fixed points.

In the usual case when  $A$  is finite, the theory is well-established, see for example [QUE], [PYT]: under the (relatively mild) assumption of primitivity, the symbolic system  $X_u, T$  is **minimal**: the closed orbit of any point under  $T$  is the whole  $X_u$ , or, equivalently, for all  $m$  there exists  $n$  such that every word of length  $m$  occurring in  $u$  occurs in every word of length  $n$  occurring in  $u$ . Under the same assumption,

the system is **uniquely ergodic**: it admits a unique invariant probability measure  $\mu$ . The measure-theoretic dynamical systems built from primitive substitutions give many interesting examples in ergodic theory, such as

**Example 1.2** (The Morse substitution).

$$\begin{aligned} a &\rightarrow ab \\ b &\rightarrow ba \end{aligned}$$

For any sequence  $u$ , the **language**  $L(u)$  is the set of all words occurring in  $u$ ; the **complexity** of  $u$  is the function  $p(n)$  which associates to each  $n \in \mathbb{N}$  the number of words of length  $n$  in  $L(u)$ . For fixed points of primitive substitutions, the complexity is always bounded by  $Cn$ , and this implies the system has topological (and hence measure-theoretic) **entropy** zero.

## 2. A FUNDAMENTAL EXAMPLE

The following broad question was asked by C. Mauduit: what can be said of the following substitution on  $A = \mathbb{Z}$ ?

**Example 2.1** (The drunken man substitution).

$$n \rightarrow (n-1)(n+1)$$

for all  $n \in \mathbb{Z}$

The first obstacle is that, if we look at the  $k$ -th image of 0, it is made only of even (resp. odd) numbers if  $k$  is even (resp. odd); this reflects the fact that the matrix has period two (see section 3 below). Hence the right substitution to consider is

**Example 2.2** (The squared drunken man substitution).

$$n \rightarrow (n-2)nn(n+2)$$

for all  $n \in A = 2\mathbb{Z}$

This substitution, which we denote by  $\sigma$ , has no fixed point; but we can define a subset  $X$  of  $A^{\mathbb{N}}$  to be the set of all sequences  $x = x_0x_1\dots$  such that every word occurring in  $x$  occurs also in  $\sigma^n 0$  for at least one  $n > 0$ .  $X$  is then a closed subset of the (noncompact) set  $\mathbb{Z}^{\mathbb{N}}$  equipped with the product topology (each copy of  $\mathbb{Z}$  being equipped with the discrete topology), and is invariant by the shift  $T$ . We say that  $(X, T)$  is the (non-compact) symbolic system associated to the substitution  $\sigma$ .

It is trivially false that, in any given sequence  $x$  of  $X$ , for all  $m$  there exists  $n$  such that every word of length  $m$  occurring in  $x$  occurs in every word of length  $n$  occurring in  $x$ ; but on an infinite alphabet the minimality of the system  $(X, T)$  would be equivalent to a weaker property, namely that any word occurring in one element of  $X$  occurs in every element of  $X$ . But in fact this property is not satisfied here, as there exist infinite sequences in  $X$  without any occurrence of the letter 0: take for example the sequence beginning by  $\sigma^n(2n)$  for all  $n$ .

Hence

**Proposition 2.3.**  $(X, T)$  is not minimal.

Though individual sequences may have strange properties, we are looking at good statistical properties for “typical” sequences of  $X$ . This involves looking for invariant measures; but here the situation is also different from the finite case, as

**Proposition 2.4.** There is no finite measure on  $X$  invariant under  $T$ .

**Definition 2.5.** For any words  $v$  and  $w$ , we say that  $v$  is an **ancestor** (under  $\sigma$ ) of  $w$  with multiplicity  $m$  if  $w$  occurs in  $\sigma v$  at  $m$  different places. If  $w = w_0 \dots w_s$ , the **cylinder**  $[w]$  is the set  $\{x \in X; x_0 = w_0, \dots, x_s = w_s\}$ .

We define the **natural measure**  $\mu$  on  $(X, T)$  by assigning to each cylinder  $[n]$ ,  $n \in A$ , or  $T^k[n]$ , the measure 1, and to a cylinder  $[w]$ , or  $T^k[w]$ , the measure  $\frac{1}{4} \sum \mu[v]m(v)$ , the sum being taken on all its ancestors  $v$  and  $m(v)$  denoting their multiplicities.

**Proposition 2.6.**  $\mu$  is an infinite measure on  $X$  invariant under  $T$ .

**Lemma 2.7.** The system  $(X, T, \mu)$  is generated by a countable family of **Rokhlin stacks**: namely, for every  $n \in \mathbb{N}$ ,  $X$  is, up to sets of  $\mu$ -measure zero, the disjoint union of the  $T^k[\sigma^n j]$ ,  $j \in 2\mathbb{Z}$ ,  $0 \leq k \leq 4^n - 1$ .

**Proposition 2.8.** The system  $(X, T, \mu)$  is **recurrent**: namely, for every set  $E$  with  $0 < \mu(E)$ ,  $\mu\{x \in E; T^n x \notin E \text{ for every } n > 0\} = 0$ .

**Proposition 2.9.** The system  $(X, T, \mu)$  is **ergodic**: namely, for every set  $E$  with  $0 < \mu(E)$  and  $\mu(E\Delta TE) = 0$ , either  $\mu(E) = 0$  or  $\mu(X/E) = 0$ .

In fact, to prove ergodicity, we prove that though we cannot define frequencies for words, we may define ratios of frequencies: namely, for almost all  $x \in X$ , and words  $w$  and  $w'$ ,  $\frac{1}{n}N(w, x_0 \dots x_{n-1})$  has limit zero when  $n \rightarrow +\infty$ , but  $\frac{N(w, x_0 \dots x_{n-1})}{N(w', x_0 \dots x_{n-1})}$  does converge to  $\frac{\mu[w]}{\mu[w']}$ .

Note that there are many  $T$ -invariant measures concentrated on the same set as  $\mu$ ; in particular, if in the definition of  $\mu$  we replace the natural constant 4 by any  $C \geq 4$ , we get another infinite measure on  $X$  invariant under  $T$  - necessarily nonergodic.

Because of the recurrence, it makes sense to study the **induced**, or first return, map of  $(X, T, \mu)$  on the cylinder  $[0]$ . Let  $(Y, S, \nu)$  be this system.

**Proposition 2.10.** The system  $(Y, S, \nu)$  is measure-theoretically isomorphic to the symbolic system associated to the substitution  $\tau$  on  $A = \mathbb{N} \times \mathbb{Z}$ , equipped with its natural measure, which is an invariant probability measure.

where  $\tau$  is the

**Example 2.11** (The induced drunken man substitution).

$$(m, n) \rightarrow \prod_{j=0}^{n-1+m^+} (j, 1) \quad (m, n+1) \quad \prod_{i=-n+1+m^-}^{-1} (i, 1)$$

for all  $m \in \mathbf{Z}$  and  $n \geq 1$ .

and its natural measure is defined by  $\nu[m, n] = \nu(S^k[m, n]) = 2^{-|m|-2n}$  and a cylinder  $[w]$ , or  $S^k[w]$ , has measure  $\frac{1}{4} \sum \mu[v]m(v)$ , the sum being taken on all its ancestors (under  $\tau$ )  $v$  and  $m(v)$  denoting their multiplicities.

**Proposition 2.12.** *The system  $(Y, S, \nu)$  is not minimal and not uniquely ergodic.*

The system  $(Y, S, \nu)$  being a finite measure-preserving system, we can compute its measure-theoretic **entropy**  $h(S, \nu)$ . Note that  $\tau$  has a fixed point  $u$ , which is the infinite sequence beginning by  $\tau^n(0, 1)$  for every  $n$ .

**Lemma 2.13.** *If, for given  $M$ , the sequence  $v(M)$  is deduced from  $u$  by replacing each  $(m, n)$  with the symbol  $\omega$  when  $|m| > M$  or  $n > M$ , then its complexity is bounded by  $C(M)n^2$ .*

**Corollary 2.14.**  $h(S, \nu) = 0$ .

### 3. GENERAL THEORY

**Definition 3.1.** *The **matrix** of a substitution  $\sigma$  is defined by  $M = ((m_{ij}))$  where  $m_{ij}$  is the number of occurrences of the letter  $j$  in the word  $\sigma i$ .*

We define the substitution  $\sigma_l$  on the alphabet  $A^l$  by associated to the  $l$ -letter  $v_1 \dots v_l$  the  $l$ -word made by enumerating all the words of length  $l$  occurring in  $\sigma(v_1 \dots v_l)$ , starting from the first position. We denote by  $M_l$  the matrix of this substitution.

**Definition 3.2.** *Let  $M$  be a matrix on a countable alphabet. We denote by  $m_{ij}(n)$  the coefficients of  $M^n$ ;  $M$  is **irreducible** if for every  $(i, j)$  there exists  $l$  such that  $M_{ij}(l) > 0$ . An irreducible  $M$  has **period**  $d$  if for every  $i$   $d = \text{GCD}\{l; m_{ii}(l) > 0\}$ , and is **aperiodic** if  $d = 1$ .*

An irreducible aperiodic matrix admits a **Perron-Frobenius** eigenvalue  $\lambda$  defined as  $\lim_{n \rightarrow +\infty} m_{ij}(n)^{\frac{1}{n}}$ .  $M$  is **transient** if

$$\sum_n m_{ij}(n) \lambda^{-n} < +\infty,$$

**recurrent** otherwise. For a recurrent  $M$ , we define  $l_{ij}(1) = m_{ij}$ ,  $l_{ij}(n+1) = \sum_{r \neq i} l_{ir}(n) m_{rj}$ ;  $M$  is **null recurrent** if

$$\sum_n n l_{ii}(n) \lambda^{-n} < +\infty,$$

and **positive recurrent** otherwise.

The reference for all the definitions and results on infinite matrices above is [KIT]. The vocabulary comes from the theory of random walks: a matrix is positive recurrent if it is the matrix of a random walk which returns to each point with probability one and the expectation of the waiting time is finite, it is null recurrent

if it is the matrix of a random walk which returns to each point with probability one and the expectation of the waiting time is infinite, and it is transient if it is the matrix of a random walk which does not return to each point with probability one. And of course, the matrix of the drunken man substitution is the matrix of the famous random walk of the same name, though the dynamical systems we can associate to these two objects are completely different.

Now, for a given substitution  $\sigma$  on a countable alphabet  $A$ , we define the dynamical system associated to  $\sigma$  in the same way as in the previous section.

**Proposition 3.3.** *If  $\sigma$  has a positive recurrent matrix, the associated system  $(X, T)$  admits a natural invariant measure which is a probability, and ergodic if  $\sigma$  is of constant length.*

The natural measure is defined by taking  $(\mu[n], n \in A)$  to be the normalized left eigenvector of  $M$  for its Perron-Frobenius eigenvalue  $\lambda$ ,  $(\mu[w], w \in A^l)$  to be the normalized left eigenvector of  $M_l$  for its (same) Perron-Frobenius eigenvalue  $\lambda$ , and  $\mu(T^k[w]) = \mu[w]$  for all cylinders. When  $\sigma$  is of constant length,  $\lambda$  is the common length of the  $\sigma n$ ,  $n \in A$ .

**Proposition 3.4.** *If  $\sigma$  has a null recurrent matrix, the associated system  $(X, T)$  admits a natural infinite invariant measure.*

The natural measure is defined as in the previous case by Perron-Frobenius eigenvectors; as for the transient case, such a measure may exist or not.

#### 4. FURTHER EXAMPLES

**Example 4.1** (The one-sided drunken man substitution).

$$n \rightarrow (n - 1)(n + 1)$$

for all  $n \geq 1$ , and

$$0 \rightarrow 1.$$

As for its two-sided counterpart, this substitution has a matrix of period 2, hence we study its square.

**Example 4.2** (The squared one-sided drunken man substitution).

$$n \rightarrow (n - 2)nn(n + 2)$$

for all even  $n \geq 2$ , and

$$0 \rightarrow 02.$$

Here the matrix is transient, but there exists an infinite invariant measure, whose value on letters is given by  $\mu[2n] = 2n + 1$ . When we induce it on the cylinder  $[0]$  we get:

**Example 4.3** (The induced one-sided drunken man substitution).

$$n \rightarrow 123 \dots (n + 1)$$

for all  $n \geq 1$ .

Turning to positive recurrent examples, we have:

**Example 4.4** (The one step forward, two step backwards, substitution).

$$n \rightarrow (n-1)(n-1)(n+1)$$

for all  $n \geq 1$ , and

$$0 \rightarrow 111.$$

As its matrix is of period two,

**Example 4.5** (The squared one step forward, two step backwards, substitution).

$$n \rightarrow (n-2)(n-2)n(n-2)(n-2)nnn(n+2)$$

for all even  $n \geq 2$ , and

$$0 \rightarrow 002002002.$$

The system has a natural invariant ergodic probability measure, which gives measure  $\frac{1}{3}$  to  $[0]$  and  $2^{-2n+1}$  to  $[2n]$ ,  $n \geq 1$ . But it is still not minimal, and not uniquely ergodic.

**Example 4.6** (The golden ratio substitution).

$$n \rightarrow (n-2)(n+1)$$

for all  $n \geq 2$ ,

$$0 \rightarrow 01,$$

$$1 \rightarrow 02.$$

It has an aperiodic matrix, positive recurrent, and the natural invariant ergodic probability gives to  $[n]$  the measure  $\frac{2^n(3-\sqrt{5})}{2(1+\sqrt{5})^n}$ .

**Example 4.7** (The infini-Bonacci substitution).

$$n \rightarrow 1(n+1)$$

for all  $n \geq 1$ .

This is a very special case, as the symbolic system is minimal and uniquely ergodic. Measure-theoretically, the system is isomorphic to the dyadic odometer, with an explicit coding to and from the system generated by the *period-doubling substitution* on two letters,  $1 \rightarrow 12$ ,  $2 \rightarrow 11$ . From the combinatorial point of view, the infini-Bonacci fixed point was used by Cassaigne to build many interesting new sequences and thus earned the unofficial nickname of *the universal counter-example*.

## REFERENCES

- [KIT] B. KITCHENS: Symbolic dynamics. One-sided, two-sided and countable state Markov shifts, Universitext. (1998), Springer-Verlag.
- [PYT] N. PYTHEAS FOGG: Substitutions in dynamics, arithmetics and combinatorics, Lecture Notes in Math. vol. 1794 (2002), Springer-Verlag.
- [QUE] M. QUEFFELEC: Substitution dynamical systems - Spectral analysis, Lecture Notes in Math. vol. 1294 (1987), Springer-Verlag.

INSTITUT DE MATHÉMATIQUES DE LUMINY CNRS, UPR 9016163 AV. DE LUMINY F13288  
MARSEILLE CEDEX 9(FRANCE) AND FÉDÉRATION DE RECHERCHE DES UNITÉS DE MATHÉMATIQUES  
DE MARSEILLE CNRS - FR 2291

*E-mail address:* `ferenczi@iml.univ-mrs.fr`