



**HAL**  
open science

# On the asymptotic behaviour of the correlation measure of sum-of-digits function in base 2

Jordan Emme, Alexander Prikhodko

► **To cite this version:**

Jordan Emme, Alexander Prikhodko. On the asymptotic behaviour of the correlation measure of sum-of-digits function in base 2. 2015. hal-01138865v1

**HAL Id: hal-01138865**

**<https://hal.science/hal-01138865v1>**

Preprint submitted on 7 Apr 2015 (v1), last revised 8 Dec 2017 (v2)

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

# On the asymptotic behaviour of the correlation measure of sum-of-digits function in base 2

Jordan Emme\*, and Alexander Prikhod'ko†

## Abstract

Let  $s_2(x)$  denote the number of digits “1” in a binary expansion of any  $x \in \mathbb{N}$ . We study the mean distribution  $\mu_a$  of the quantity  $s_2(x+a) - s_2(x)$  for a fixed positive integer  $a$ . It is shown that solutions of the equation

$$s_2(x+a) - s_2(x) = d$$

are uniquely identified by a finite set of prefixes in  $\{0, 1\}^*$ , and that the probability distribution of differences  $d$  is given by an infinite product of matrices whose coefficients are operators of  $l^1(\mathbb{Z})$ .

Then, denoting by  $l(a)$  the number of patterns “01” in the binary expansion of  $a$ , we give the asymptotic behaviour of this probability distribution as  $l(a)$  goes to infinity as well as estimates of the variance of the probability measure  $\mu_a$

## Introduction

In this article we are interested in the statistic behaviour of the difference of the number of digits 1 in the binary expansion of an integer  $x$  before and after its summation with  $a$ . This kind of question can be linked with carry propagation problems developed in [5] and [9] and has to do with computer arithmetics as in [8] or [3] but our approach is different. Moreover, the function  $s_2$  modulo 2 was extensively studied for its links with the Thue-Morse sequence (as we can find in [4]) for example or for arithmetic reason as in [2] and [6]. But in this paper our motivation is not of the same nature and we will not look at  $s_2 \bmod 2$  but just at the function  $s_2$  which is the sum of the digits in base 2.

Because of the simplicity of base 2 and of the summation process, we can give detailed behaviour of the correlation between the number of digits 1 in the binary expansion of  $x$  and  $x+a$  for any integer  $a$ . More precisely, we are able to compute the distribution of the difference  $s_2(x+a) - s_2(x)$  where  $s_2(y)$  denotes the number of digits 1 in the binary expansion of any integer  $y$ . We essentially give two methods in Section 1. The first one, very elementary, just helps understanding the basics of our construction. The second that follows is a little bit more involved but allows us to present our result in a much more usable form for what interests us and what will be developed in the next section.

---

\*Aix-Marseille Université, CNRS, Centrale Marseille, I2M, UMR 7373, 13453 Marseille, France. E-mail : jordan.emme@univ-amu.fr

†Moscow Institute of Physics and Technology, Moscow, Russia. E-mail : sasha.prikhodko@gmail.com

Section 2 is a technical part that allows us to understand the asymptotical behaviour of such a correlation as  $a$  gets bigger in a certain, non trivial sense. This section develops the idea that our correlation gets smaller as the number of subwords “01” in the binary expansion of  $a$  increases. To go further and develop this idea of increasing the number of patterns “01” in the binary expansion of  $a$ , we demonstrate a limit theorem for the random variable consisting of the difference of 1 in  $x$  and  $x + a$  for different  $a$  going to infinity in the previous sense.

Let us now give some definitions and notations that will be used throughout this article.

**Definition 0.0.1.** For any integer  $x \in \mathbb{N}$  whose binary expansion is given by :

$$x = \sum_{k=0}^n x_k 2^k, \quad \forall k \in \{0, \dots, n\} \quad x_k \in \{0, 1\},$$

we define the quantity

$$s_2(x) = \sum_{k=0}^n x_k$$

which is the number of digits 1 in the binary expansion of  $x$ .

We will also define the following :

**Definition 0.0.2.** For any integer  $x$  whose binary expansion is given by :

$$x = \sum_{k=0}^n x_k 2^k,$$

we denote by  $\underline{x}$  the word  $x_0 \dots x_n$  in  $\{0, 1\}^*$ .

**Remark 0.0.3.** Since we are working in the free monoïd  $\{0, 1\}^*$  of binary words, let us remark that we will denote the cylinder set of a word  $w = w_0 \dots w_n$  by the standard notation  $[w] := \{v \in \{0, 1\}^{\mathbb{N}} \mid v_0 \dots v_n = w_0 \dots w_n\}$ . These sets form a topological basis of clopen of  $\{0, 1\}^{\mathbb{N}}$  for the product topology. Moreover, we endow the set of binary configurations  $\{0, 1\}^{\mathbb{N}}$  with the natural probability measure  $\mathbb{P}$  that is the balanced Bernoulli probability measure defined on the Borel sets.

With that in mind, the equation we wish to study is the following, with parameters  $a \in \mathbb{N}$  and  $d \in \mathbb{Z}$  :

$$s_2(x + a) - s_2(x) = d \tag{1}$$

More particularly, we wish to understand the behaviour, for any integer  $a$ , of the correlation between  $s_2(x + a)$  and  $s_2(x)$  given by :

$$\mu_a(d) := \lim_{N \rightarrow \infty} \frac{1}{N} \# \{x \leq N \mid s_2(x + a) - s_2(x) = d\}.$$

**Remark 0.0.4.** Let us remark that for all integer  $a$ , this defines a probability measure on  $\mathbb{Z}$ .

To this end, the first section will essentially be about the combinatorial description of the solutions of (1). We will describe the addition process via the construction of a particular tree. This will allow us to have an explicit formula for the desired distribution of probability only depending on the binary expansion of  $a$  as a product of matrices. The main result is the following :

**Theorem.** *The distribution  $\mu_a$  is calculated via an infinite product of matrices whose coefficients are operators of  $l^1(\mathbb{Z})$  applied to a vector whose coefficients are elements of  $l^1(\mathbb{Z})$*

$$\mu_a = (Id, Id) \cdots A_{a_n} A_{a_{n-1}} \cdots A_{a_1} A_{a_0} \begin{pmatrix} \delta_0 \\ 0 \end{pmatrix},$$

where the sequence  $(a_n)_{n \in \mathbb{N}}$  is the binary expansion of  $a$ ,  $\delta_0$  is the Dirac mass in 0

$$A_0 = \begin{pmatrix} Id & \frac{1}{2}S^{-1} \\ 0 & \frac{1}{2}S \end{pmatrix}, \quad A_1 = \begin{pmatrix} \frac{1}{2}S^{-1} & 0 \\ \frac{1}{2}S & Id \end{pmatrix},$$

and  $S$  is the left shift transformation on  $l^1(\mathbb{Z})$ .

**Remark 0.0.5.** We remind that the set of finite measures on  $\mathbb{Z}$  is in bijection with the elements of  $l^1(\mathbb{Z})$ . We will always identify finite measures on  $\mathbb{Z}$  and elements of  $l^1(\mathbb{Z})$ .

Such a result allows an analytical study of these distributions as the binary expansion of  $a$  is more and more complex. Let us introduce this notion of complexity for  $a$  :

**Definition 0.0.6.** For any  $a \in \mathbb{N}$ , let us denote by  $l(a)$  the number of subwords 01 in the binary expansion of  $a$ .

We are thus interested in what happens as there are more and more patterns 01 in the word  $a$ . We can precisely estimate the asymptotic behaviour of the  $l^2(\mathbb{Z})$  norm of this distribution as  $a$  tends to infinity by increasing the number of subwords 01.

Namely, the theorem is as follows :

**Theorem.** *There exists a real constant  $C_0$  such that for any integer  $a$  we have the following :*

$$\|\mu_a\|_2 \leq C_0 \cdot l(a)^{-1/4}.$$

Finally, in the last section, we wish to obtain a much more precise result regarding the behaviour of such a distribution than just estimates of the  $l^2$  norm. So we study in which way the variance of the random variable of probability law  $\mu_a$  is linked to the number of subwords 01 in the binary expansion of  $a$ . We have bounds on this variance as shown in this result :

**Theorem.** *For any integer  $a$  such that  $l(a)$  is large enough, the variance  $-2V(a)$  of  $\mu_a$  has bounds :*

$$l(a) - 1 \leq -2V(a) \leq 2(2l(a) + 1).$$

# 1 Statistics of binary sequences

In this section we wish to understand the following quantity, for any given positive integer  $a$  and any integer  $d$  :

$$\mu_a(d) := \lim_{N \rightarrow \infty} \frac{1}{N} \# \{x \leq N \mid s_2(x+a) - s_2(x) = d\}.$$

We know from [1] that such a limit exists and has been studied in [7] for example, but we will give a proof using the structure of solutions of the following equation :

$$s_2(x+a) - s_2(x) = d.$$

for any  $a$  and  $d$ .

Let us investigate such solutions as well as their construction in order to understand the distribution of probability  $\mu_a$  of differences  $d = s_2(x+a) - s_2(x)$ . We prove that this distribution is given by an infinite product of matrices whose sequence is given by the binary expansion of  $a$ .

## 1.1 Combinatorial description of summation tree

First we prove the following lemma :

**Lemma 1.1.1.** *For all  $a \in \mathbb{N}$  and  $d \in \mathbb{Z}$ , there exists a finite set of words  $\mathcal{P}_{a,d} = \{p_1, \dots, p_k\} \subset \{0, 1\}^*$  such that  $x$  is solution of (1) if and only if :*

$$\underline{x} \in \bigcup_{j=1}^k [p_j]$$

where  $[w]$  denotes the cylinder set of configurations whose prefix is  $w$  for any word  $w \in \{0, 1\}^*$ .

**Remark 1.1.2.** Before we prove Lemma 1.1.1, let us remark that for any even integer  $n$  the following holds :

$$s_2(n) = s_2\left(\frac{n}{2}\right).$$

*Proof.* Let us prove this lemma by induction on  $a$ . It is obviously true that for  $a = 0$  and  $a = 1$ , there exists a set of words for any  $d \in \mathbb{Z}$  (possibly empty) describing the solutions of Equation 1. Let us assume it is true for every integer not greater than a given  $a \in \mathbb{N}$ .

If  $a$  is even, then, for any  $d$ , let the set of words  $\mathcal{P}_{a+1,d}$  be  $\{0w, w \in \mathcal{P}_{\frac{a}{2}, d-1}\} \cup \{1w, w \in \mathcal{P}_{\frac{a}{2}+1, d+1}\}$ .

Indeed,  $x \in \mathbb{N}$  being an even solution of  $s_2(x+a+1) - s_2(x) = d$  is equivalent to having  $s_2(x+a) + 1 - s_2(\frac{x}{2}) = d$  which can be written  $s_2(\frac{x}{2} + \frac{a}{2}) - s_2(\frac{x}{2}) = d - 1$ . From the induction hypothesis follows that  $\underline{x}$  starts with a 0 followed by a word beginning by a word in  $\mathcal{P}_{\frac{a}{2}, d-1}$ .

Moreover,  $x \in \mathbb{N}$  being an odd solution of  $s_2(x+a+1) - s_2(x) = d$  is equivalent to having  $s_2(x-1+2+a) - (s_2(\frac{x-1}{2}) + 1) = d$  which can be written  $s_2(\frac{x-1}{2} + \frac{a}{2} + 1) - s_2(\frac{x-1}{2}) = d + 1$ . Then we can state that  $\underline{x}$  must start with a 1 followed by a word in  $\mathcal{P}_{\frac{a}{2}+1, d+1}$  by induction hypothesis.

If  $a$  is odd, we have  $\mathcal{P}_{a+1,d} = \{0w, w \in \mathcal{P}_{\frac{a+1}{2},d}\} \cup \{1w, w \in \mathcal{P}_{\frac{a+1}{2},d}\}$  for any  $d$ .

In fact, let  $x \in \mathbb{N}$  be an odd solution of  $s_2(x+a+1) - s_2(x) = d$ . This is of course equivalent to having the equality  $s_2(x+a) + 1 - (s_2(\frac{x-1}{2}) + 1) = d$  which can be written  $s_2(\frac{x-1}{2} + \frac{a+1}{2}) - s_2(\frac{x-1}{2}) = d$ . Moreover,  $x$  being an even solution of  $s_2(x+a+1) - s_2(x) = d$  is equivalent to having the equality  $s_2(\frac{x}{2} + \frac{a+1}{2}) - s_2(\frac{x}{2}) = d$ . These last equalities, with the induction hypothesis can be summarized as follows :  $\underline{x}$  begins either by a 0 or a 1 and is followed by a word in  $\mathcal{P}_{\frac{a+1}{2},d}$ . □

From this lemma follows immediatly that for all positive integer  $a$  and integer  $d$ , the sequence  $(\#\{x \leq N \mid s_2(x+a) - s_2(x) = d\})_{N \in \mathbb{N}}$  is a sum of sequences which contain an arithmetic progression and thus the following limit exists :

$$\mu_a(d) := \lim_{N \rightarrow \infty} \frac{1}{N} \#\{x \leq N \mid s_2(x+a) - s_2(x) = d\}.$$

Moreover, a quick computation yields

$$\mu_a(d) = \sum_{p \in \mathcal{P}_{a,d}} \mathbb{P}(p)$$

where  $\mathbb{P}$  is the balanced Bernouilli probability measure on  $\{0,1\}^{\mathbb{N}}$ .

The proof of Lemma 1.1.1 naturally gives the idea of an inductive way to compute the prefixes. A comfortable way to do that is to span, for every  $a$ , a tree  $\tau_a$  whose paths will represent the binary decomposition of some  $x$  and which will allow us to keep track of the quantity  $s_2(x+a) - s_2(x)$  as the summation goes. Let us state the following lemma :

**Lemma 1.1.3.** *For each  $a \in \mathbb{N}$ , there exists a tree  $\tau_a$  with vertices labeled in  $\{0, \dots, a\} \times \mathbb{Z}$  and edges labelled in  $\{0,1\}$  such that for any  $d \in \mathbb{Z}$ , words in  $\mathcal{P}_{a,d}$  are exactly paths from the vertex  $(a,0)$  to a vertex  $(0,d)$ .*

*Proof.* Let us choose a particular  $a$  in  $\mathbb{N}$  and construct the associated tree  $\tau_a$ . The tree  $\tau_a$  will have vertices labelled in  $\{0, \dots, a\} \times \mathbb{Z}$  and edges labelled in the binary alphabet  $\{0,1\}$ . This tree will also have its vertices on different levels ordered in  $\mathbb{N}$ . We will construct our tree inductively on each level.

If the tree is defined up to level  $n \in \mathbb{N}$ . From every vertex we construct the vertices at level  $n+1$  in the following way.

Starting from a vertex labelled by  $(k,c)$  :

- If  $k$  is even we add a vertex at level  $n+1$  labelled by  $(\frac{k}{2}, c)$  and add one edge of each type between these vertices.
- If however  $k$  is odd, we add two vertices labelled by  $(\frac{k-1}{2}, c+1)$  and  $(\frac{k+1}{2}, c-1)$  and we construct an edge of type 0 between  $(k,c)$  and  $(\frac{k-1}{2}, c+1)$  and an edge of type 1 between  $(k,c)$  and  $(\frac{k+1}{2}, c-1)$ .

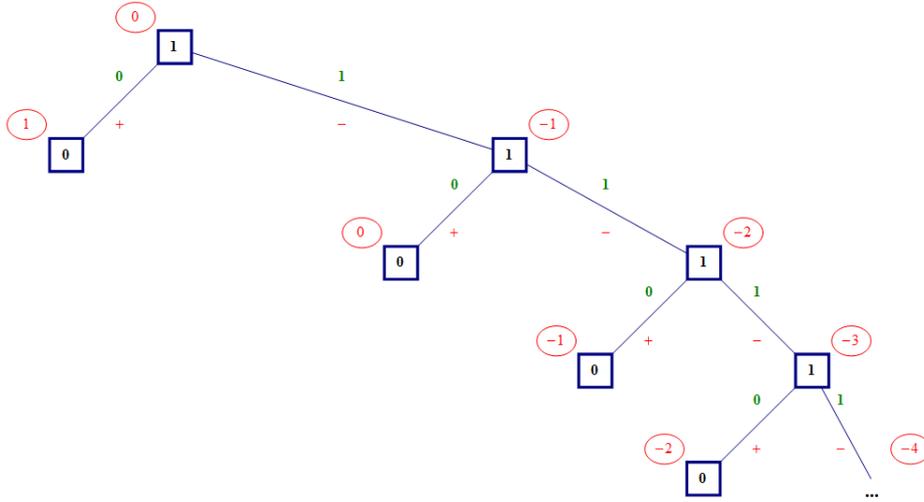


Figure 1: Infinite subtree spanned by every vertex labelled by 1

The construction starts at level 0 with only one vertex labelled by  $(a, 0)$ .

Of course such a construction actually spans an infinite tree for we eventually build a vertex labelled by 1. Then, notice that a vertex 1 is parent of two vertices 0 and 1. Iterating the procedure builds an infinite subtree rooted in 1 as given in the example of figure 1.

The goal of this construction is then to read prefixes which are solutions of the equation as paths on this graph. Let us start at the vertex labelled  $(a, 0)$ . The second coordinate of labels of vertices is a counter that keeps track of the difference  $s_2(x + a) - s_2(x)$  as the summation is processed digit by digit, and, in our case, that is represented by which vertex we go through as we follow  $\underline{x}$  as a path in  $\tau_a$  starting from  $(a, 0)$ .

Starting from an odd vertex, if we follow the edge labelled by 0 then we add one to our counter and if we follow the edge labelled by 1 then we subtract 1 to our counter.

If however we start from an even vertex, the counter remains unchanged, regardless of which edge we follow.

This allows to keep track of the difference  $s_2(x + a) - s_2(x)$ . Once we reach a vertex  $(0, d)$ , the summation is completed, so in order to find a prefix we only need to find a path from the vertex  $(a, 0)$  to a vertex  $(0, d)$  such that the counter reaches the desired  $d$ .

For any  $a$ , the associated tree allows us to compute the set of prefixes  $\mathcal{P}_{a,d}$  for any  $d$ . Let us choose a positive integer  $x$ . We will follow the path associated to  $\underline{x}$  starting from the vertex  $a$ . Then, on every level, following the path and keeping doing the indicated operations on our counter is just doing the summation digit by digit while keeping track of the 1 we have lost or added during this summation. If we are not careful, we might not end on a 0 (if the binary expansion of  $x$  is too short and/or a carry is still being propagated). However, adding enough 0 to  $\underline{x}$  ensures that we do not encounter this problem.

This tree is then a convenient way to compute  $d = s_2(x + a) - s_2(x)$  and as such, to find the prefixes we are interested in as the different paths starting from  $(a, 0)$  and ending on a vertex  $(0, d)$ .

□

We can see an example of such a construction on figure 2 for  $a = 5$ . Vertices' labels are boxed, edges labels are above them, the counter numbers are circled and, as a reminder, the increment or decrement we need to apply after each step are indicated on the edges.

**Example 1.1.4.** Let  $a = 5$ . Then we start our construction on the first level of the tree by creating the root labeled by  $(5,0)$ . The number 5 being odd, we add the vertices  $(2,1)$  and  $(3,-1)$  on the second level and add an edge of type 0 between vertices  $(2,1)$  and  $(5,0)$  and an edge of type 1 between vertices  $(3,-1)$  and  $(5,0)$ .

We continue the construction on third level by :

- adding a child to vertex  $(2,1)$  labelled by  $(1,1)$  since 2 is even.
- adding two children to vertex  $(3,-1)$  labelled by vertices  $(1,0)$  and  $(2,-2)$  since 3 is odd. The edge between vertices  $(1,0)$  and  $(3,-1)$  is of type 0 and the edge between vertices  $(2,-2)$  and  $(3,-1)$  is of type 1.

And we continue this process on each vertex on each level to obtain the tree of figure 2.

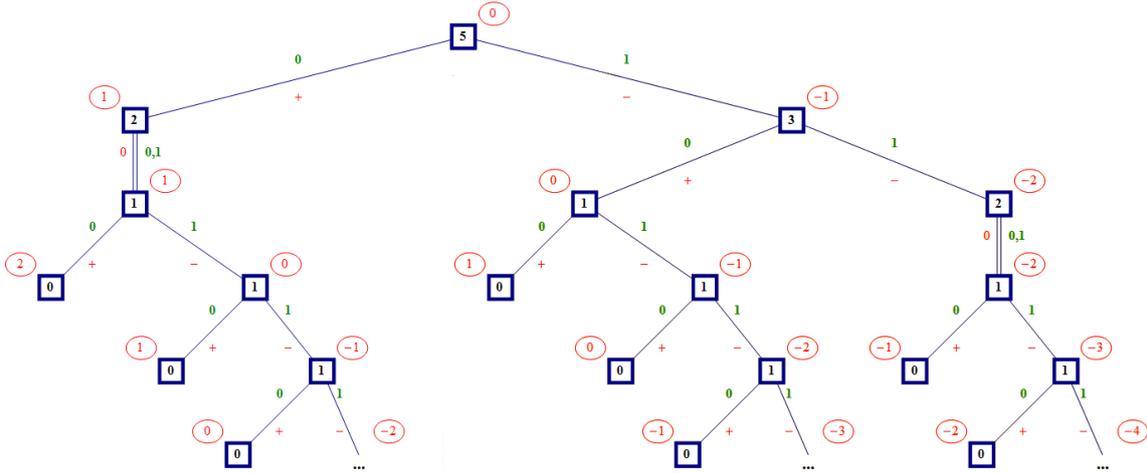


Figure 2: Part of the tree  $\tau_5$ .

Let us compute the words in  $\mathcal{P}_{5,0}$ . It can be read on the tree  $\tau_5$  that this set is given by the words 00110, 01110 and 1010.

We wish to remark that because of the way  $\underline{x}$  is defined, the word 00110 for example is not the binary expansion of 6 but of 12 as everything is mirrored.

Let us now remark that having such a family of trees proves the following proposition :

**Proposition 1.1.5.** For every  $a$  in  $\mathbb{N}$  and every  $d$  in  $\mathbb{Z}$ , we have the following identities :

$$\mu_{2a}(d) = \mu_a(d)$$

and

$$\mu_{2a+1}(d) = \frac{1}{2}\mu_a(d-1) + \frac{1}{2}\mu_{a+1}(d+1).$$

*Proof.* For any positive integer  $a$  and any integer  $d$ , we can read from the tree  $\tau_{2a}$  that words in  $\mathcal{P}_{2a,d}$  are exactly words beginning by either 0 or 1 and followed by a word in  $\mathcal{P}_{a,d}$ , thus we have :

$$\mu_{2a}(d) = \mu_a(d)$$

Notice that words in  $\mathcal{P}_{2a+1,d}$  are exactly words beginning by a 0 and followed by a word in  $\mathcal{P}_{a,d-1}$  and the ones beginning by a 1 followed by a word in  $\mathcal{P}_{a+1,d+1}$ . It follows that :

$$\mu_{2a+1}(d) = \frac{1}{2}\mu_a(d-1) + \frac{1}{2}\mu_{a+1}(d+1).$$

□

**Remark 1.1.6.** This proposition allows to compute explicitly any distribution  $\mu_a$  since it is trivial to give explicit closed formulas for  $\mu_0$  and  $\mu_1$ . It also gives an understanding of the link between  $\mu_a$  and the binary expansion of  $a$ . However, we prefer another presentation of such a result. Namely, one that assembles in a close formula these identities as we follows the binary decomposition of  $a$ . Such a formula is given in Theorem 1.2.1.

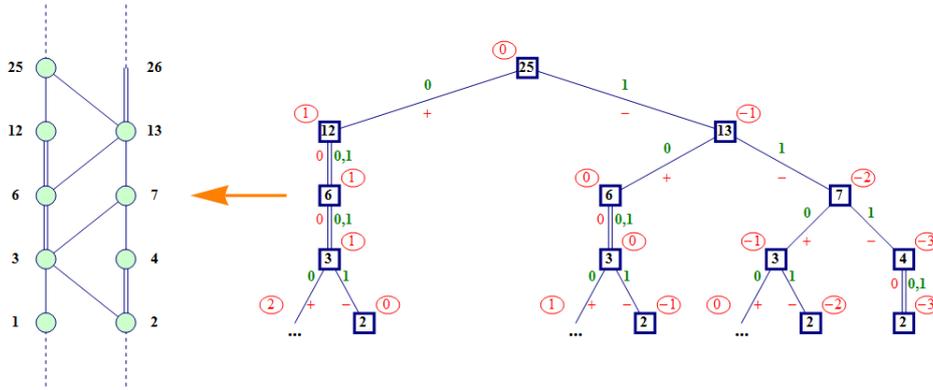


Figure 3: Collapsing  $\tau_{25}$ .

One way to do such a thing is to remark that if we collapse the tree  $\tau_a$  by identifying the vertices whose label's first coordinate are the same, we obtain a graph whose structure is studied in the next section.

Looking at this graph after collapsing is quite interesting for it gives us an understanding of the link between the summation process and the binary expansion of  $a$ . It also allows the statistical study of the behaviour of  $s_2(x+a) - s_2(x)$ . An example of the tree  $\tau_{25}$  collapsing is given in Figure 3.

## 1.2 Distribution of $s_2(x+a) - s_2(x)$

The goal of this section is to prove the following theorem :

**Theorem 1.2.1.** *The distribution  $\mu_a$  is calculated via an infinite product of matrices whose coefficients are operators of  $l^1(\mathbb{Z})$*

$$\mu_a = (Id, Id) \cdots A_{a_n} A_{a_{n-1}} \cdots A_{a_1} A_{a_0} \begin{pmatrix} \delta_0 \\ 0 \end{pmatrix},$$

where the sequence  $(a_n)_{n \in \mathbb{N}}$  is the binary expansion of  $a$ ,  $\delta_0$  is the Dirac mass in 0

$$A_0 = \begin{pmatrix} Id & \frac{1}{2}S^{-1} \\ 0 & \frac{1}{2}S \end{pmatrix}, \quad A_1 = \begin{pmatrix} \frac{1}{2}S^{-1} & 0 \\ \frac{1}{2}S & Id \end{pmatrix},$$

and  $S$  is the left shift transformation on  $l^1(\mathbb{Z})$ .

One can notice that only two different patterns can appear between two levels of the collapsed graph. Those patterns are given by figure 4.

**Remark 1.2.2.** Let us remark that there are always only two vertices on each level (except at level 0 where we can add the vertex  $a+1$  for coherence) labelled by two consecutive integers and that the order in which each pattern appears and how vertices are labelled are given by the binary expansion of  $a$ .

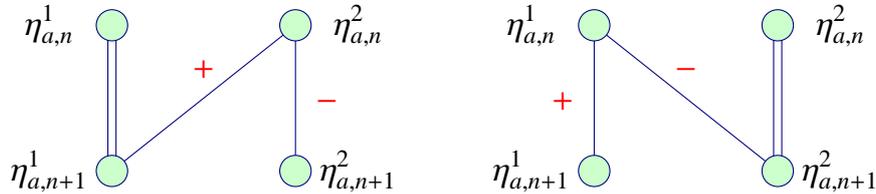


Figure 4: The two different patterns encountered in the collapsed graph

First, we begin by a definition :

**Definition 1.2.3.** Let us define for each positive integer  $a$  and each level  $n$  of the collapsed graph the probability measures  $\eta_{a,n}^1$  and  $\eta_{a,n}^2$  on  $\mathbb{Z}$  such that  $\eta_{a,n}^i(d)$  is the probability that a path of length  $n$  starting from  $a$  in the associated collapsed graph ends on the  $i^{\text{th}}$  vertex at level  $n$  with the counter value being  $d$ . Here we still endow  $\{0, 1\}^*$  with the balanced Bernoulli probability measure.

This seems to only define a sequence in  $l^1(\mathbb{Z})$ . Let us remind that we can identify finite measures on  $\mathbb{Z}$  with sequences in  $l^1(\mathbb{Z})$ . It is a quick check to verify that we did define probability measures.

**Lemma 1.2.4.** *The  $\eta_{a,n}^i$  are given by induction :*

$$\forall n \in \mathbb{N}, \quad \begin{pmatrix} \eta_{a,n+1}^1 \\ \eta_{a,n+1}^2 \end{pmatrix} = A_{a_n} \begin{pmatrix} \eta_{a,n}^1 \\ \eta_{a,n}^2 \end{pmatrix}$$

where

$$a = \sum_{k=0}^{+\infty} a_k 2^k.$$

*Proof.* Let  $n \in \mathbb{N}$ ,  $d \in \mathbb{Z}$  and  $a \in \mathbb{Z}_+$ . We wish to compute  $\eta_{a,n+1}^1$  and  $\eta_{a,n+1}^2$ . It is obvious that the pattern we see between levels  $n$  and  $n+1$  are given by  $a_n$ .

If  $a_n = 0$ , then we encounter the left pattern of figure 4. In this case we have :

$$\eta_{a,n+1}^1(d) = \eta_{a,n}^1(d) + \frac{1}{2}\eta_{a,n}^2(d-1).$$

This being true for all  $d$ , we can write :

$$\eta_{a,n+1}^1 = \eta_{a,n}^1 + \frac{1}{2}S^{-1}(\eta_{a,n}^2),$$

where  $S$  is the left shift transformation on the space  $l^1(\mathbb{Z})$ . For the same pattern we also have :

$$\eta_{a,n+1}^2(d) = \frac{1}{2}\eta_{a,n}^2(d+1)$$

which can be rewritten

$$\eta_{a,n+1}^2 = \frac{1}{2}S(\eta_{a,n}^2).$$

So in this particular case we can write the relations in the following way :

$$\begin{pmatrix} \eta_{a,n+1}^1 \\ \eta_{a,n+1}^2 \end{pmatrix} = \begin{pmatrix} Id & \frac{1}{2}S^{-1} \\ 0 & \frac{1}{2}S \end{pmatrix} \begin{pmatrix} \eta_{a,n}^1 \\ \eta_{a,n}^2 \end{pmatrix}.$$

The same arguments for the right pattern of figure 4 (which corresponds to the case where  $a_n = 1$ ) yields :

$$\begin{pmatrix} \eta_{a,n+1}^1 \\ \eta_{a,n+1}^2 \end{pmatrix} = \begin{pmatrix} \frac{1}{2}S^{-1} & 0 \\ \frac{1}{2}S & Id \end{pmatrix} \begin{pmatrix} \eta_{a,n}^1 \\ \eta_{a,n}^2 \end{pmatrix}$$

which ends the proof of the lemma. □

This lemma allows us to prove Theorem 1.2.1 :

*Proof.* Let  $a \in \mathbb{N}$  and  $\mu_a \in l^1(\mathbb{Z})$  be the distribution of probability of the difference  $s_2(x+a) - s_2(x)$ . From Lemma 1.1.1 we know that the set of solution for  $s_2(x+a) - s_2(x) = d$  for any  $d \in \mathbb{Z}$  is given by a finite set of prefixes. We can consider them to be all of size  $m_d$  by adding prefixes if needed (which means partitionning our cylinder sets from Lemma 1.1.1 into smaller

cylinders of constant size  $(\frac{1}{2})^{m_d}$ . We can also assume that  $m_d$  is bigger than the highest  $n$  such that  $a_n = 1$  with the same argument of adding prefixes (i.e. partitionning again).

Then the probability  $\mu_a(d)$  is given by  $\eta_{a,m_d}^1(d)$  since we need to end the summation on a vertex labelled by 0. Let us also notice that  $\eta_{a,m_d}^2(d) = 0$  for otherwise we would have a solution prefix whose length would be greater than  $m_d$ . With the previous lemma we thus have the following :

$$\mu_a(d) = \left( (Id, Id) A_{a_{m_d}} A_{a_{m_d-1}} \cdots A_{a_1} A_{a_0} \begin{pmatrix} \eta_{a,0}^1 \\ \eta_{a,0}^2 \end{pmatrix} \right) (d).$$

It is also worth noticing that after rank  $m_d$  keeping multiplying on the left by  $A_0$  will not have any influence on the final result so this allows for a more convenient writing of our result :

$$\mu_a(d) = \left( (Id, Id) \cdots A_{a_n} A_{a_{n-1}} \cdots A_{a_1} A_{a_0} \begin{pmatrix} \eta_{a,0}^1 \\ \eta_{a,0}^2 \end{pmatrix} \right) (d)$$

which proves the theorem. □

## 2 Asymptotic properties of distributions $\mu_a$

In this section we study the asymptotic behaviour of the family of distributions  $\mu_a$  as  $l(a)$  increases. We prove that the  $l^2(\mathbb{Z})$  norm of  $\mu_a$  tends to 0 as  $l(a)$  increases hence the densities of  $\mu_a$  tend to zero.

Let us remind that  $l(a)$  denotes the number of distinct patterns 01 in the binary expansion of  $a$ .

**Remark 2.0.5.** For any integer  $a$  there exists  $k$  integers  $p_1, \dots, p_k$  such that  $\underline{a} = 0^{p_1} 1^{p_2} \dots 1^{p_k}$  or  $\underline{a} = 1^{p_1} 0^{p_2} \dots 1^{p_k}$  depending on  $a$  being even or odd. In the first case,  $k = 2l(a)$ , and in the second,  $k = 2l(a) + 1$ .

In other words, the number of distinct continuous maximal blocks of digits in  $\underline{a}$  is either  $2l(a) + 1$  or  $2l(a)$ .

### 2.1 Convergence to zero

We start by proving the following theorem.

**Theorem 2.1.1.** *There exists a constant  $C_0$  such that for any integer  $a$  we have the following :*

$$\|\mu_a\|_2 \leq C_0 \cdot l(a)^{-1/4}.$$

The proof is established in a series of lemmas. Let us recall the notation

$$A_0 = \begin{pmatrix} Id & \frac{1}{2}S^{-1} \\ 0 & \frac{1}{2}S \end{pmatrix}, \quad A_1 = \begin{pmatrix} \frac{1}{2}S^{-1} & 0 \\ \frac{1}{2}S & Id \end{pmatrix},$$

where  $S$  is the left shift on  $l^1(\mathbb{Z})$ . Let us define a norm on the space of  $n \times n$  matrices  $Y = (y_{i,j})_{i,j \in \{1, \dots, n\}^2}$  by the following :

$$\|Y\| = \max_{1 \leq j \leq n} \sum_{i=1}^n |y_{i,j}|.$$

**Remark 2.1.2.** Let us remark right away that this defines a submultiplicative norm.

We consider the Fourier transform  $\hat{\mu}_a$  defined on the circle  $\mathbb{T}$  by :

$$\forall \theta \in [0, 2\pi), \quad \hat{\mu}_a(\theta) = \sum_{d \in \mathbb{Z}} e^{-id\theta} \mu_a(d).$$

Notice that applying the operator  $S$  to an element of  $l^1(\mathbb{Z})$  transforms to the multiplication by  $e^{i\theta}$  of the fourier transform hence the linear operators  $A_0, A_1$  become ordinary linear transformations in  $\mathbb{C}^2$  for any given  $\theta$  defined in the following way

$$\hat{A}_0(\theta) := \begin{pmatrix} 1 & \frac{1}{2}e^{-i\theta} \\ 0 & \frac{1}{2}e^{i\theta} \end{pmatrix}, \quad \hat{A}_1(\theta) := \begin{pmatrix} \frac{1}{2}e^{-i\theta} & 0 \\ \frac{1}{2}e^{i\theta} & 1 \end{pmatrix},$$

**Lemma 2.1.3.** *We have the following elementary inequalities for any  $\theta$ :*

$$\|\hat{A}_0(\theta)\|_1 = \|\hat{A}_1(\theta)\|_1 = \|\hat{A}_0(\theta)\hat{A}_1(\theta)\|_1 = \|\hat{A}_1(\theta)\hat{A}_0(\theta)\|_1 = 1 \text{ and}$$

$$\|\hat{A}_0(\theta)\hat{A}_1(\theta)\hat{A}_0(\theta)\|_1 = \|\hat{A}_1(\theta)\hat{A}_0(\theta)\hat{A}_1(\theta)\|_1 := \phi(\theta) = \frac{1 + \sqrt{5 + 4 \cos \theta}}{4}.$$

Notice that  $\phi$  is strictly less than 1 except when  $\theta = 0$ .

*Proof.* For any  $\theta \in [0, 2\pi)$ , we have :

$$\hat{A}_0(\theta)\hat{A}_1(\theta)\hat{A}_0(\theta) = \begin{pmatrix} \frac{1}{2}e^{-i\theta} + \frac{1}{4} & \frac{1}{4}e^{-2i\theta} + \frac{1}{8}e^{-i\theta} + \frac{1}{4} \\ \frac{1}{4}e^{2i\theta} & \frac{1}{8}e^{i\theta} + \frac{1}{4}e^{2i\theta} \end{pmatrix}.$$

and

$$\hat{A}_1(\theta)\hat{A}_0(\theta)\hat{A}_1(\theta) = \begin{pmatrix} \frac{1}{8}e^{-i\theta} + \frac{1}{4}e^{-2i\theta} & \frac{1}{4}e^{-2i\theta} \\ \frac{1}{4}e^{2i\theta} + \frac{1}{8}e^{i\theta} + \frac{1}{4} & \frac{1}{2}e^{i\theta} + \frac{1}{4} \end{pmatrix}.$$

so he have

$$\|\hat{A}_0(\theta)\hat{A}_1(\theta)\hat{A}_0(\theta)\|_1 = \|\hat{A}_1(\theta)\hat{A}_0(\theta)\hat{A}_1(\theta)\|_1.$$

Moreover, the triangular inequality yields

$$\left| \frac{1}{4}e^{-2i\theta} + \frac{1}{8}e^{-i\theta} + \frac{1}{4} \right| + \left| \frac{1}{4}e^{2i\theta} + \frac{1}{8}e^{i\theta} \right| \leq \left| \frac{1}{4}e^{-2i\theta} \right| + \left| \frac{1}{8}e^{-i\theta} + \frac{1}{4} \right| + \left| \frac{1}{4}e^{2i\theta} + \frac{1}{8}e^{i\theta} \right|$$

so we have

$$\left| \frac{1}{4}e^{-2i\theta} + \frac{1}{8}e^{-i\theta} + \frac{1}{4} \right| + \left| \frac{1}{4}e^{2i\theta} + \frac{1}{8}e^{i\theta} \right| \leq \left| \frac{1}{4}e^{2i\theta} \right| + \left| \frac{1}{2}e^{-i\theta} + \frac{1}{4} \right|$$

hence the sum of the modulus of the terms in the first column is less than the sum of the modulus of the terms of the second column, which implies that

$$\|\hat{A}_0(\theta)\hat{A}_1(\theta)\hat{A}_0(\theta)\|_1 = \left| \frac{1}{4}e^{-2i\theta} \right| + \left| \frac{1}{2}e^{i\theta} + \frac{1}{4} \right|$$

so finally

$$\|\hat{A}_0(\theta)\hat{A}_1(\theta)\hat{A}_0(\theta)\|_1 = \frac{1 + \sqrt{5 + 4 \cos \theta}}{4}.$$

□

**Lemma 2.1.4.** *For any  $k \in \mathbb{N}^*$  and any  $\theta \in [0, 2\pi)$*

$$\|\hat{A}_0(\theta) (\hat{A}_1(\theta))^k \hat{A}_0(\theta)\|_1 \leq \|\hat{A}_0(\theta) \hat{A}_1(\theta) \hat{A}_0(\theta)\|_1.$$

*Proof.* First let us state that, for all positive integer  $k$  :

$$\forall \theta \in [0, 2\pi), \quad (\hat{A}_1(\theta))^k = \begin{pmatrix} \frac{e^{-ik\theta}}{2^k} & 0 \\ s_k(\theta) & 1 \end{pmatrix}$$

where

$$s_k(\theta) = e^{2i\theta} \sum_{j=1}^k \frac{e^{-ij\theta}}{2^j}.$$

We can compute that :

$$\forall k \in \mathbb{N}, \quad \hat{A}_0(\theta) (\hat{A}_1(\theta))^k \hat{A}_0(\theta) = \begin{pmatrix} \frac{e^{-ik\theta}}{2^k} + \frac{e^{-i\theta}}{2} s_k(\theta) & \frac{e^{-i(k+1)\theta}}{2^{k+1}} + \frac{e^{-2i\theta}}{4} s_k(\theta) + \frac{1}{4} \\ \frac{e^{i\theta}}{2} s_k(\theta) & \frac{1}{4} s_k(\theta) + \frac{e^{2i\theta}}{4} \end{pmatrix}.$$

Moreover, noticing that for every positive integer  $k$  and every  $\theta$  in  $[0, 2\pi)$ ,

$$s_{k+1}(\theta) = \frac{e^{-i\theta}}{2} s_k(\theta) + \frac{e^{i\theta}}{2}$$

implies that

$$\left| \frac{e^{-i(k+1)\theta}}{2^{k+1}} + \frac{e^{-i\theta}}{2} s_{k+1}(\theta) \right| + \left| \frac{e^{i\theta}}{2} s_{k+1}(\theta) \right| = \left| \frac{e^{-i(k+1)\theta}}{2^{k+1}} + \frac{e^{-2i\theta}}{4} s_k(\theta) + \frac{1}{4} \right| + \left| \frac{1}{4} s_k(\theta) + \frac{e^{2i\theta}}{4} \right|.$$

Which means that the sum of modulus of coefficient on the first column of matrix  $\hat{A}_0(\theta) (\hat{A}_1(\theta))^{k+1} \hat{A}_0(\theta)$  is equal to the sum of the modulus of the coefficients of the second column of the matrix  $\hat{A}_0(\theta) (\hat{A}_1(\theta))^k \hat{A}_0(\theta)$ . Hence to prove that

$$\forall k \in \mathbb{N}, \quad \forall \theta \in [0, 2\pi), \quad \|\hat{A}_0(\theta) (\hat{A}_1(\theta))^k \hat{A}_0(\theta)\|_1 \leq \|\hat{A}_0(\theta) \hat{A}_1(\theta) \hat{A}_0(\theta)\|_1$$

suffices to actually show the following :

$$\forall k \geq 2, \quad \forall \theta \in [0, 2\pi), \quad \left| \frac{e^{-ik\theta}}{2^k} + \frac{e^{-i\theta}}{2} s_k(\theta) \right| + \left| \frac{e^{i\theta}}{2} s_k(\theta) \right| \leq \|\hat{A}_0(\theta) \hat{A}_1(\theta) \hat{A}_0(\theta)\|_1.$$

Using triangular inequality, we have

$$\forall k \geq 2, \quad \forall \theta \in [0, 2\pi), \quad \left| \frac{e^{-ik\theta}}{2^k} + \frac{e^{-i\theta}}{2} s_k(\theta) \right| + \left| \frac{e^{i\theta}}{2} s_k(\theta) \right| \leq \frac{1}{2^k} + |s_k(\theta)|$$

which, still using triangular inequality, yields

$$\forall k \geq 2, \quad \forall \theta \in [0, 2\pi), \quad \left| \frac{e^{-ik\theta}}{2^k} + \frac{e^{-i\theta}}{2} s_k(\theta) \right| + \left| \frac{e^{i\theta}}{2} s_k(\theta) \right| \leq \frac{1}{2^k} + \left| \frac{e^{i\theta}}{2} + \frac{1}{4} \right| + \sum_{j=3}^k \frac{1}{2^j}$$

and

$$\forall \theta \in [0, 2\pi), \quad \left| \frac{e^{-i2\theta}}{2^2} + \frac{e^{-i\theta}}{2} s_2(\theta) \right| + \left| \frac{e^{i\theta}}{2} s_2(\theta) \right| \leq \frac{1}{2^2} + \left| \frac{e^{i\theta}}{2} + \frac{1}{4} \right|$$

So, in any case, we get :

$$\forall k \geq 2, \quad \forall \theta \in [0, 2\pi), \quad \left| \frac{e^{-ik\theta}}{2^k} + \frac{e^{-i\theta}}{2} s_k(\theta) \right| + \left| \frac{e^{i\theta}}{2} s_k(\theta) \right| \leq \left| \frac{e^{i\theta}}{2} + \frac{1}{4} \right| + \frac{1}{4}$$

which completes the proof since

$$\left| \frac{e^{i\theta}}{2} + \frac{1}{4} \right| + \frac{1}{4} = \|\hat{A}_0(\theta) \hat{A}_1(\theta) \hat{A}_0(\theta)\|_1$$

□

**Lemma 2.1.5.** *If  $-\pi \leq \theta \leq \pi$  then*

$$\phi(\theta) \leq e^{-\theta^2/15}.$$

*Hence for any positive integer  $N$  :*

$$\|\phi^N\|_2 \leq \|e^{-N\theta^2/15}\|_2 = \left( \frac{15\pi}{2N} \right)^{1/4}$$

*with  $\phi$  as defined in 2.1.3.*

The proof is left to the reader.

*Proof of theorem 2.1.1.* Recall that  $\mu_a = \mu_{a,1} + \mu_{a,2}$  is the sum of the components of the vector  $\bar{\mu}_a$ , which is calculated as an infinite product

$$\bar{\mu}_a = \begin{pmatrix} \mu_{a,1} \\ \mu_{a,2} \end{pmatrix} = \dots A^{a_n} A^{a_{n-1}} \dots A^{a_1} A^{a_0} \begin{pmatrix} \delta_0 \\ 0 \end{pmatrix}.$$

Let us represent the binary expansion of  $a$  as a sequence of  $l(a)$  groups  $11\dots 1$  separated by zeros:

$$\underline{a} = 0^{m_1} \underline{11\dots 1} 0^{m_2} \underline{11\dots 1} 0 \dots 0^{m_{l(a)}} \underline{11\dots 1}.$$

We apply Lemma 2.1.4 to half of the patterns  $011\dots 10$  and use the bound of Lemma 2.1.3 :

$$\|\hat{A}_0(\theta) (\hat{A}_1(\theta))^k \hat{A}_0(\theta)\|_1 \leq \phi(\theta).$$

We can do this only on half of patterns  $011\dots 10$  to avoid problematic patterns like  $\dots 0111011110\dots$ , since we need at least two “0” between two blocks on “1” to be able to apply Lemma 2.1.4 twice but we have only one “0” in between. Thus, we get

$$\|\dots \hat{A}^{a_n}(\theta) \dots \hat{A}^{a_0}(\theta)\|_1 \leq \phi(\theta)^N, \quad N = \frac{l(a) - 1}{2},$$

and hence, for each  $\theta \in [0, 2\pi)$  and  $j \in \{1, 2\}$

$$|\hat{\mu}_{a,j}(\theta)| \leq \phi(\theta)^N \cdot \left\| \begin{pmatrix} 1 \\ 0 \end{pmatrix} \right\|_1 = e^{-N\theta^2/15},$$

and

$$\|\mu_{a,j}\|_2 = \frac{1}{\sqrt{2\pi}} \|\hat{\mu}_{a,j}\|_2 \leq \frac{1}{\sqrt{2\pi}} \|\phi^N\|_2 = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} e^{-Nt^2/15} dt = \left( \frac{15}{8\pi N} \right)^{1/4}.$$

Finally, we get using lemma 2.1.5

$$\|\mu_a\|_2 \leq \sqrt{2} \left( \frac{15}{\pi(l(a) - 1)} \right)^{1/4} = O(l(a)^{-1/4}).$$

□

**Remark 2.1.6.** It follows directly from Theorem 2.1.1 that the density  $\mu_a(j) \rightarrow 0$  as  $l(a) \rightarrow \infty$ .

Let us also remark that this theorem actually gives an important information as to the fact that the probability measure  $\mu_a$  is linked with the complexity of the binary expansion of  $a$ , complexity which is measured by  $l(a)$ . In fact, let us remark that for any integer  $n$ ,  $\mu_{2^n} = \mu_1$  so it is possible to find arbitrarily large  $a$  such that  $\|\mu_a\|_2$  does not tend to zero (actually whenever  $l(a)$  does not tend to infinity).

## 2.2 Asymptotic mean and variance of $\mu_a$

Let us first recall the following :

$$\forall \theta \in [0, 2\pi), \quad \hat{A}_0(\theta) := \begin{pmatrix} 1 & \frac{1}{2}e^{-i\theta} \\ 0 & \frac{1}{2}e^{i\theta} \end{pmatrix}, \quad \hat{A}_1(\theta) := \begin{pmatrix} \frac{1}{2}e^{-i\theta} & 0 \\ \frac{1}{2}e^{i\theta} & 1 \end{pmatrix},$$

We begin by using Taylor’s expansion on the matrices  $\hat{A}_0(\theta)$  and  $\hat{A}_1(\theta)$  we have :

$$\hat{A}_k(\theta) = I_k + \theta \alpha_k + \theta^2 \beta_k + O(\theta^3), \quad k \in \{0, 1\},$$

where

$$\begin{aligned} I_0 &= \begin{pmatrix} 1 & \frac{1}{2} \\ 0 & \frac{1}{2} \end{pmatrix}, & \alpha_0 &= \frac{i}{2} \begin{pmatrix} 0 & -1 \\ 0 & 1 \end{pmatrix}, & \beta_0 &= -\frac{1}{4} \begin{pmatrix} 0 & 1 \\ 0 & 1 \end{pmatrix}, \\ I_1 &= \begin{pmatrix} \frac{1}{2} & 0 \\ \frac{1}{2} & 1 \end{pmatrix}, & \alpha_1 &= \frac{i}{2} \begin{pmatrix} -1 & 0 \\ 1 & 0 \end{pmatrix}, & \beta_1 &= -\frac{1}{4} \begin{pmatrix} 1 & 0 \\ 1 & 0 \end{pmatrix}, \end{aligned}$$

and we observe that the following relations hold

$$(1, 1)I_k = (1, 1), \quad (2)$$

$$(1, 1)\alpha_k = (0, 0), \quad (1, 1)\beta_0 = -\frac{1}{2}(0, 1), \quad (1, 1)\beta_1 = -\frac{1}{2}(1, 0). \quad (3)$$

**Lemma 2.2.1.** *For all  $\theta \in [0, 2\pi)$ , the product of matrices  $(\hat{A}_0(\theta))^k$  converges as  $k$  goes to  $+\infty$ . We denote the limit  $\hat{A}_0^\infty(\theta)$*

$$\hat{A}_0^\infty(\theta) = \begin{pmatrix} 1 & \frac{e^{-i\theta}}{2-e^{i\theta}} \\ 0 & 0 \end{pmatrix}.$$

*Proof.* Remarking the following :

$$\forall \theta \in [0, 2\pi), \quad (\hat{A}_0(\theta))^k = \begin{pmatrix} 1 & \overline{s_k(\theta)} \\ 0 & \frac{e^{ik\theta}}{2^k} \end{pmatrix}$$

where

$$\overline{s_k(\theta)} = e^{-2i\theta} \sum_{j=1}^k \frac{e^{ij\theta}}{2^j}.$$

yields the lemma. □

Let us now remark that the asymptotic expansion of  $\hat{A}_0^\infty$  is :

$$\hat{A}_0^\infty(\theta) = I_\infty + \theta^2 \beta_\infty + O(\theta^3),$$

where

$$I_\infty = \begin{pmatrix} 1 & 1 \\ 0 & 0 \end{pmatrix}, \quad \beta_\infty = -\begin{pmatrix} 0 & 1 \\ 0 & 0 \end{pmatrix}.$$

**Definition 2.2.2.** Given  $a \in \mathbb{N}$  of binary length  $N$ , let us define  $\Pi_a(\theta) = \hat{A}_0^\infty \hat{A}_{a_N}(\theta) \cdots \hat{A}_{a_0}(\theta)$ . Let  $D := \{v \in \mathbb{R}^2 : v_1 + v_2 = 1\}$ .

**Lemma 2.2.3.** *For any real vector  $v_0 \in D$  the product  $(1, 1) \cdot \Pi_a(\theta) \cdot v_0$  has the following asymptotic expansion*

$$(1, 1) \cdot \Pi_a(\theta) \cdot v_0 = 1 + V(a)\theta^2 + O(\theta^3),$$

where

$$V(a) = (1, 1) (\beta_\infty v_{N+1} + \beta_{a_N} v_N + \dots + \beta_{a_0} v_0),$$

and

$$\forall j \in \{1, \dots, N+1\}, \quad v_j = I_{a_{j-1}} \cdots I_{a_0} v_0.$$

*Proof.* Let us first recall that :

$$(1, 1)I_0 = (1, 1)I_1 = (1, 1)I_\infty = (1, 1).$$

Moreover :

$$(1, 1) \cdot \Pi_a(\theta) \cdot v_0 = (1, 1) (I_\infty + \theta^2 \beta_\infty + O(\theta^3)) \prod_{i=1}^N (I_{a_i} + \alpha_{a_i} \theta + \beta_{a_i} \theta^2 + O(\theta^3)) v_0.$$

Hence the constant term in the asymptotic expansion of  $(1, 1) \cdot \Pi_a(\theta) \cdot v_0$  is 1 (since  $v_0$  is in  $D$ ). In addition, we have :

$$(1, 1)\alpha_0 = (1, 1)\alpha_1 = (0, 0)$$

so the term of degree 1 is 0.

With the same argument, we can compute the quadratic term  $V(a)$ . Notice that any coefficient  $\alpha_0$  or  $\alpha_1$  is killed by left multiplication by  $(1, 1)$ . Hence the terms of the form :

$$I_\infty \cdots I_{a_{j+1}} \alpha_{a_j} I_{a_{j-1}} \cdots I_{a_{i+1}} \alpha_{a_i} I_{a_{i-1}} \cdots I_{a_0}$$

disappear after the left multiplication by  $(1, 1)$  and thus do not contribute in any way to the coefficient of the quadratic term.

So the only terms contributing to  $V(a)$  in the product  $\hat{A}_\infty(\theta) \hat{A}_{a_N}(\theta) \dots \hat{A}_{a_0}(\theta)$  are the terms with coefficient  $\beta_0, \beta_1$  or  $\beta_\infty$ . □

**Lemma 2.2.4.** *For any  $j \in \{0, \dots, N\}$ , the value  $(1, 1)\beta_{a_j} v_j$  is estimated as follows*

$$\frac{1}{2^{b_j^1}} \left(1 - \frac{1}{2^{b_j^2}}\right) \leq -(1, 1)\beta_{a_j} v_j \leq \frac{1}{2^{b_j^1}},$$

where  $b_j^1$  is the length of a consecutive and maximal sequence of digits equal to  $a_j$  to the left of  $j$  (including position  $j$ ) in  $\underline{a}$ , and  $b_j^2$  is the length of the next block of identical digits to the left of the block containing the digit  $a_j$ .

*Proof.* Without loss of generality assume that  $a_j = 0$ . Observe that  $I_0$  and  $I_1$  contract the segment  $[(1, 0), (0, 1)]$  to the first and to the second point respectively with the contraction factor 2. Hence the block  $I_1^{b_j^2}$  maps  $[(1, 0), (0, 1)]$  into the segment  $[(2^{-b_j^2}, (1 - 2^{-b_j^2})), (0, 1)]$ . Then the block  $I_0^{b_j^1-1}$  contracts later to

$$\left[ (1 - 2^{-b_j^1+1}(1 - 2^{-b_j^2}), 2^{-b_j^1+1}(1 - 2^{-b_j^2})), (1 - 2^{-b_j^1+1}, 2^{-b_j^1+1}) \right].$$

And multiplying on the left by  $(1, 1)\beta_{a_j} = (1, 1)\beta_0 = -\frac{1}{2}(0, 1)$  means taking the second coordinate with an additional coefficient  $-1/2$ , and the lemma follows. □

**Remark 2.2.5.** The same proof yields :

$$-(1, 1)\beta_\infty v_{N+1} \leq \frac{1}{2^{b_N^1-1}}.$$

Moreover, we have the following lemma :

**Lemma 2.2.6.**

$$l(a) \leq \sum_{j=0}^N \frac{1}{2^{b_j^1}} \leq 2l(a) + 1.$$

*Proof.* This lemma is obtained by noticing that each block of identical digits in the binary expansion of  $a$  spans in  $\sum_{j=0}^N \frac{1}{2^{b_j^1}}$  a partial sum of a geometric sequence of ratio  $\frac{1}{2}$  and first term  $\frac{1}{2}$  so each is smaller than 1 and greater than  $\frac{1}{2}$ . In fact let us assume that there are  $k$  blocks in the binary expansion of  $a$  of lengths  $l_1, \dots, l_k$  (i.e  $a = 0^{l_1}1^{l_2}\dots 0^{l_{k-1}}1^{l_k}$  for example). Then

$$\sum_{j=0}^N \frac{1}{2^{b_j^1}} = \sum_{i=1}^k \sum_{j=1}^{l_i} \left(\frac{1}{2}\right)^j$$

and

$$\sum_{i=1}^k \frac{1}{2} \leq \sum_{i=1}^k \sum_{j=1}^{l_i} \left(\frac{1}{2}\right)^j \leq \sum_{i=1}^k 1.$$

Moreover,  $k$  can be equal to either  $2l(a) + 1$  or  $2l(a)$  which completes the proof. □

**Theorem 2.2.7.** *For any integer  $a$  such that  $l(a)$  is large enough, the variance  $-2V(a)$  of  $\mu_a$  has bounds :*

$$l(a) - 1 \leq -2V(a) \leq 2(2l(a) + 1).$$

*Proof.* The idea of this theorem is to see that each block of same digits in the binary expansion of  $a$  has an impact on  $V(a)$  estimated by a constant.

Let us estimate the value of  $V(a)$  :

$$-V(a) = -\sum_{j=0}^N (1, 1)\beta_{a_j} v_j - (1, 1)\beta_\infty v_{N+1} \leq \sum_{j=0}^N \frac{1}{2^{b_j^1}} + \frac{1}{2^{b_N^1-1}},$$

by Lemma 2.2.4. Moreover, the upper bound of Lemma 2.2.6 yields:

$$-2V(a) \leq 2(2l(a) + 2).$$

Now we will prove the lower bound :

$$-V(a) \geq \frac{1}{2}(l(a)).$$

Let us remark that

$$-V(a) = -\sum_{j=0}^N (1, 1)\beta_{a_j} v_j - (1, 1)\beta_{\infty} v_{N+1} \geq -\sum_{j=0}^N (1, 1)\beta_{a_j} v_j$$

so, using lemma (2.2.4), we have :

$$-V(a) \geq \sum_{j=0}^N \left( \frac{1}{2^{b_j^1}} \left( 1 - \frac{1}{2^{b_j^2}} \right) \right).$$

Hence,

$$-V(a) \geq \sum_{j=0}^N \frac{1}{2^{b_j^1}} - \sum_{j=0}^N \frac{1}{2^{b_j^1 + b_j^2}}.$$

Let us denote by  $k$  the largest integer such that  $a_0 = \dots = a_k$ . Noticing that  $b_j^2 = 0$  for all  $j \in \{0, \dots, k\}$  we have :

$$-V(a) \geq \sum_{j=k+1}^N \frac{1}{2^{b_j^1}} - \sum_{j=k+1}^N \frac{1}{2^{b_j^1 + b_j^2}}.$$

Moreover, for any integer  $j$  larger than  $k$ ,  $b_j^2 \geq 1$  hence :

$$-V(a) \geq \frac{1}{2} \sum_{j=k+1}^N \frac{1}{2^{b_j^1}}$$

from which we derive :

$$-V(a) \geq \frac{1}{2}(l(a) - 1)$$

by applying the lower bound of Lemma 2.2.6 . This proves Theorem 2.2.7. □

On a final note, it is interesting to remark that the bounds for Theorem 2.2.7 might not be optimal but that for a sequence of integers  $(a(n))_{n \in \mathbb{N}}$  such that  $\lim_{n \rightarrow \infty} l(a(n)) = +\infty$  and such that the limit of  $\frac{-2V(a(n))}{l(a(n))}$  exists, then this limit value is in  $[1, 4]$ . We give some examples of computer simulations in Figure 5.

It would be interesting to understand the asymptotic behaviour of the ratio  $\frac{-2V(a(n))}{l(a(n))}$ . Having a necessary condition for convergence and a precise idea of how it behaves asymptotically would help understanding the variance of the probability measure  $\mu_a$ .

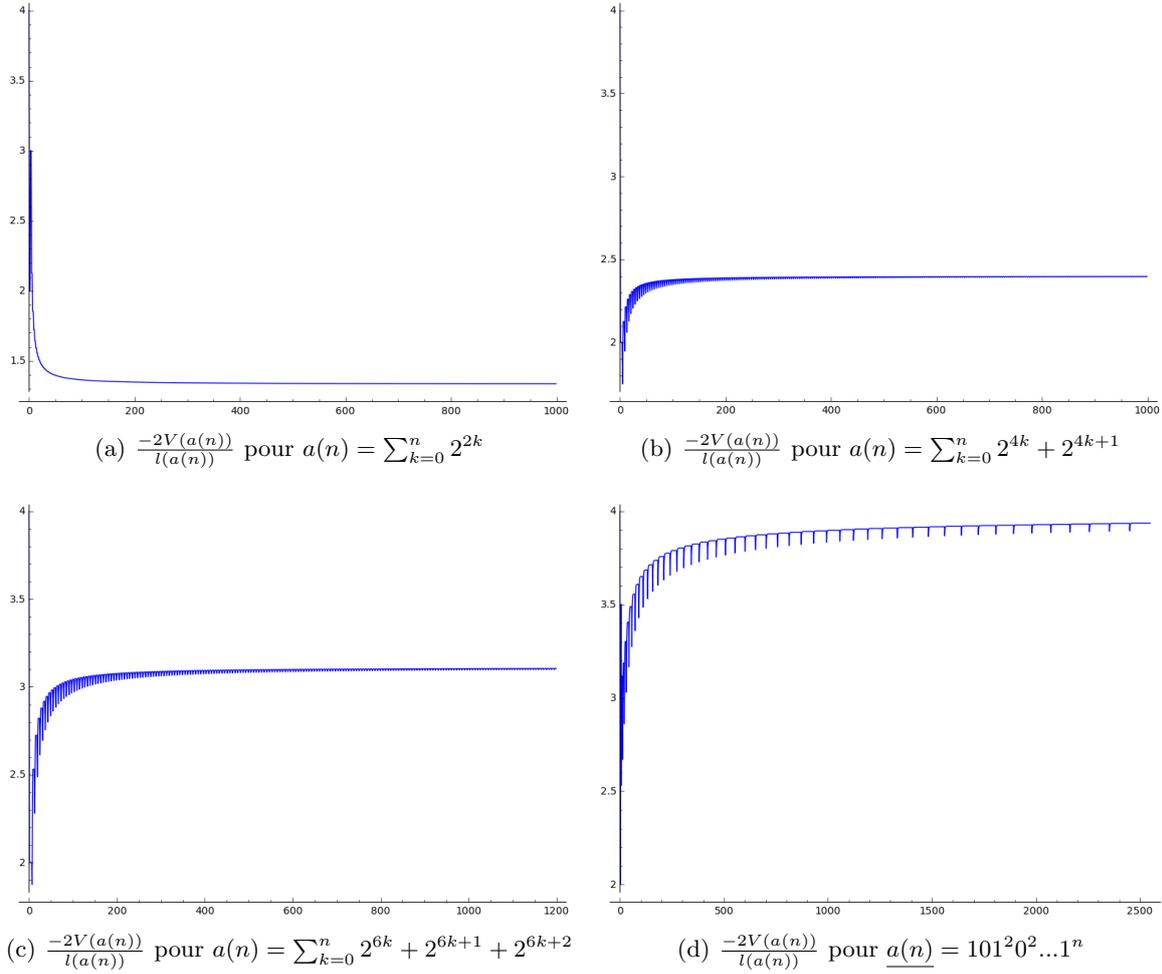


Figure 5: Asymptotics of the ratio  $\frac{-2V(a(n))}{l(a(n))}$  for different sequences  $(a(n))_{n \in \mathbb{N}}$

## References

- [1] Jean Bésineau. Indépendance statistique d'ensembles liés à la fonction "somme des chiffres". In *Séminaire Delange-Pisot-Poitou, 13e année (1971/72), Théorie des nombres, Fasc. 2, Exp. No. 23*, page 8. Secrétariat Mathématique, Paris, 1973.
- [2] Michael Drmota, Christian Mauduit, and Joël Rivat. Primes with an average sum of digits. *Compos. Math.*, 145(2):271–292, 2009.
- [3] Miloš D Ercegovac and Tomas Lang. *Digital arithmetic*. Elsevier, 2003.
- [4] M. Keane. Generalized Morse sequences. *Z. Wahrscheinlichkeitstheorie und Verw. Gebiete*, 10:335–353, 1968.

- [5] Donald E. Knuth. The average time for carry propagation. *Indagationes Mathematicae (Proceedings)*, 81(1):238 – 242, 1978.
- [6] Christian Mauduit and Joël Rivat. Sur un problème de Gelfond: la somme des chiffres des nombres premiers. *Ann. of Math. (2)*, 171(3):1591–1646, 2010.
- [7] Johannes F. Morgenbesser and Lukas Spiegelhofer. A reverse order property of correlation measures of the sum-of-digits function. *Integers*, 12:Paper No. A47, 5, 2012.
- [8] Jean-Michel Muller. *Arithmétique des ordinateurs*. Masson, 1989.
- [9] Nicholas Pippenger. Analysis of carry propagation in addition: An elementary approach. *Journal of Algorithms*, 42(2):317 – 333, 2002.