

ECOLE NORMALE SUPERIEURE DE LYON

COMPUTER SCIENCE DEPARTMENT

MASTER'S INTERNSHIP REPORT

Asymptotics on Lempel-Ziv 78 compression

Intern:

Guillaume DUBOC

Supervisor:

Wojtek SZPANKOWSKI

Center for Science of Information - Purdue University

May 21st - August 8th, 2018



École Normale Supérieure de
Lyon



Purdue University

Asymptotics on Lempel-Ziv 78 compression

Abstract

Diving into universal data compression through the analysis of the Lempel-Ziv78 algorithm - with the help of probabilistic models, computer simulation and complex analysis.

Contents

I	Introduction	3
1	Introduction to the problem	3
1.1	Compression algorithm : Lempel-Ziv 78	3
1.2	Probabilistic models	4
1.3	Entropy	4
1.4	Probabilistic analysis	5
2	Theorems and Goals	5
II	Numerical simulation of LZ78 for Markovian sources	5
1	Conditions of the experiment	5
1.1	Simulation details	6
1.2	Implementation	6
1.3	Raw results	6
1.4	Empirical normalization	7
2	Validating variance candidates	8
2.1	A first expression	8
2.2	Using the Frobenius eigenvalue of $P(s)$	8
3	Conclusion	12
III	Tail symbols analysis	13
1	Markov Independent Model	13
2	Tail symbols	13
2.1	Definition	13
2.2	Recurrence relation	14
3	Total path length	14
3.1	Definition	14
3.2	Recurrence relation	14
4	Poisson tranform differential equation	15
4.1	Covariance recurrence relation	15
4.2	Poisson transform	15
IV	Optimal parsing	16
1	LZ77 and LZ78 comparison	16
1.1	Practical uses of LZ77	16
1.2	16
2	Observations regarding the results and the proof	16
3	Definitions	17
3.1	Definitions	17
3.2	Clues in proving the result	19
3.3	Working with these probabilistic objects	19
4	Proof of an upperbound for the paper	20
5	Research at Purdue	21

Appendices	23
A Asymptotics of the mean $E[M_n]$	23
B Much longer words	23
C Another (more complicated) computation of $\check{\lambda}(-1)$	25

I Introduction

Data compression, source coding, or bit-rate reduction involve encoding information using fewer bits than the original representation.

1 Introduction to the problem

Data compression is achieved by different types of algorithms, but in order to study it formally we use a formal definition of a data compression scheme. First, we describe an algorithm which operates on words - in our case, Lempel-Ziv 78. Then, we introduce a probabilistic model representing the data to be compressed, which allows us to quantify the efficiency of the algorithm.

1.1 Compression algorithm : Lempel-Ziv 78

Definition 1 A compression algorithm is a pair of functions on words $(\mathcal{C}, \mathcal{D})$ - compressor and decompressor.

Definition 2 A lossless compression algorithm decompresses the exact same data that was compressed in the first place.

Remark 1 We will not describe the decompression part of our algorithms, considering that the possibility of decompression is obvious in the case of LZ78.

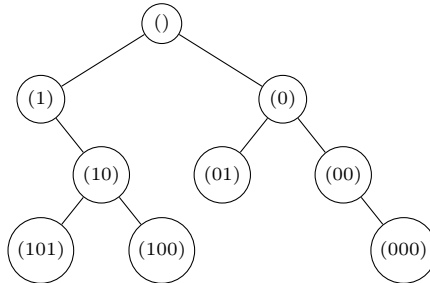
In general, Lempel-Ziv algorithms dictionary-based scheme which exploit previously seen patterns and redundancy to save off coding space. The Lempel-Ziv'78 (LZ'78) that partitions a sequence into phrases or blocks of variable size such that a new phrase is the shortest substring not seen in the past as a phrase. Every such phrase is encoded by the index of its prefix appended by a symbol; thus the LZ'78 code contains the pairs (**pointer**, **symbol**). For example, the string

11001010001000100

of length 17 is parsed as

(1)(10)(0)(101)(00)(01)(000)(100)

which can be encoded in the following digital search tree:



1.2 Probabilistic models

Definition 1 Let \mathcal{A} be an alphabet. An information source is a one-sided infinite sequence of random variables $(X_k)_{k=1}^{+\infty}$ with each X_k having values in \mathcal{A} .

Remark 2 Each realization of an information source is called a sequence or word.

Remark 3 Defining the law of the X_k produces out different models for data generation which can be studied mathematically and simulated.

Definition 2 A memoryless source is an information source for which the X_k are mutually independent, following the uniform law on $\mathcal{A} = \{a_1, \dots, a_V\}$:

$$P(X_k = a_k) = p_k \quad \text{with} \quad \sum_{i=1}^V p_i = 1$$

Remark 4 This is the simplest information source and it has been studied successfully in the past, but it is not a realistic model. We replace it with the following Markov model whenever possible.

Definition 3 A Markov source is an information source with a Markov dependency between successive symbols.

Definition 4 A Markov source of order r is a Markov source for which each symbol apparition depends on the previous r symbols.

Remark 5 We will study Markov sources of order 1 - where each symbol simply depends on the previous one. This is general enough, as Markov sources of superior order can be simulated by expanding the alphabet and using a Markov source of order 1.

1.3 Entropy

In this study, we consider only two information sources: the *memoryless* and the *Markov* source. For the Markov source, we consider that there exists a *stationary distribution*.

Definition 5 The entropy of an information is the average rate at which information is produced by a probabilistic source of information.

Property 1 The entropy of a memoryless source on alphabet $\mathcal{A} = \{a_1, \dots, a_V\}$ is:

$$\sum_{i=1}^V \log p_i$$

Property 2 The entropy of a Markov source with probability $(p_{ij})_{(i,j) \in \{1, \dots, V\}^2}$ and a stationary distribution (π_1, \dots, π_V) is:

$$h = \sum_{i=1}^V \sum_{j=1}^V \pi_i \log p_{ij}$$

The entropy gives a lower bound on the compression ratio of lossless algorithms. Furthermore, in our case, for a memoryless or Markov source, the compression tends asymptotically to the entropy of the source:

Property 3
$$\lim_{n \rightarrow +\infty} \frac{C_n}{n} = h$$

1.4 Probabilistic analysis

Under this context, we can define and conduct a thorough analysis of several random variables with different meanings regarding the effectiveness of compression.

Notation 1 Let n be an integer - the size of the considered words. Defining Ω_n the set of words of size n on alphabet \mathcal{A} . Each word being an event, a natural probability space is given by considering the output of size n of a Markov source.

Remark 6 We will now study random variables defined on this probability space.

Notation 2 W_n denotes words of size n output by a given Markov source.

Notation 3 The number of phrases used to compress words of size n (W_n) with LZ78 is given by $M_n(W_n)$ or simply M_n .

Remark 7 This is one of the most important variables to consider because, as it will appear shortly, it is closely tied to the compression ratio of LZ78.

Definition 6 The codelength is the number of bits required to encode the LZ78-compressed version of a word, denoted by $C(W_n)$ or C_n .

Definition 7 The compression ratio for words of size n is the ratio between the codelength and initial size of a word, $\frac{C(W_n)}{|W_n|} = \frac{C_n}{n}$.

Remark 8 We focus now, as stated in the title of this report, on the asymptotical behavior - for $n \rightarrow +\infty$ - of the compression ratio. Certainly this gives us information about real-world compression, as files on our computer often attain large numbers when measured in bits. However, this assumption can be nuanced. For example, the speed of convergence towards this asymptotic behavior can make a non-negligible difference in real-world application, as we will see later on when we talk about the difference between LZ'77 and LZ'78 and Optimal Parsing.

2 Theorems and Goals

There is a number of interrogations and results that surround the LZ'78 compression scheme, and that have been doing so for some time now. These were proven for the simplistic memoryless source model which will allow us to formulate these as theorems for memoryless sources. But the more interesting case of Markov sources has all these theorems become conjectures. My first job was to simulate the process and output visualizations that would disprove or make these results seem more likely.

▷ For memoryless source

II Numerical simulation of LZ78 for Markovian sources

1 Conditions of the experiment

The first step of my internship was to simulate the probabilistic model of our data representation (*i.e.* Markov sources) and implement the Lempel-Ziv 78 algorithm in order to visualize its behavior on the model. In particular, to see if our formulation of the Central Limit Theorem would hold for Markov sources, and begin to identify the asymptotics of the mean and variance of the number of phrases in particular.

1.1 Simulation details

The process for doing this simulation is pretty straightforward, as follows:

- Generating a random Markov chain of size 2 of matrix

$$\begin{pmatrix} p_{00} & p_{01} \\ p_{10} & p_{11} \end{pmatrix}$$

- Generating $n_{\text{exp}} \sim 10^3$ words of length n (or n_{word}), with $n \sim 10^6$ or 10^7
- Applying LZ78 on each of these words to estimate, for each n , the number of phrases M_n . A simple histogram of these values can be seen in figure 1.
- From this sampling of the random variable M_n and other parameters such as the entropy of the Markov chain, computing
 - the empirical mean (μ) and the empirical variance (σ^2)
 - different theoretical expressions for the variance
- Using these expressions to standardize M_n in different ways, plotting
 - the probability distribution of M_n (standardized)
 - the cumulative distribution function of M_n (standardized)
- Finally, comparing the different theoretical expressions for the variance by plotting their differences for a large range of values of n , and a constant number of experiments n_{exp} .

1.2 Implementation

I took great care to implement this process, in *Python*, in a way that would be easy to use and easy to modify during my internship as well as for future research - for example, I could not find a Python library that would implement arbitrary Markov sources, and had to implement my own. The whole code is available on a *Github* repository: <https://github.com/gliboc/lz-compression>.

1.3 Raw results

After implementation, I could start visualizing the probability distribution of M_n for different values of n . A probability distribution visualization can be obtained by normalizing raw histograms such as the one in Figure 1, which represents the counts of the different values taken by M_n for $n = 10^6$.

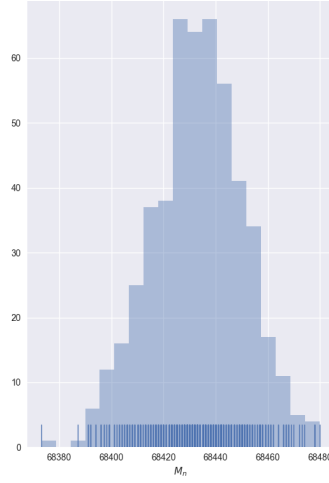
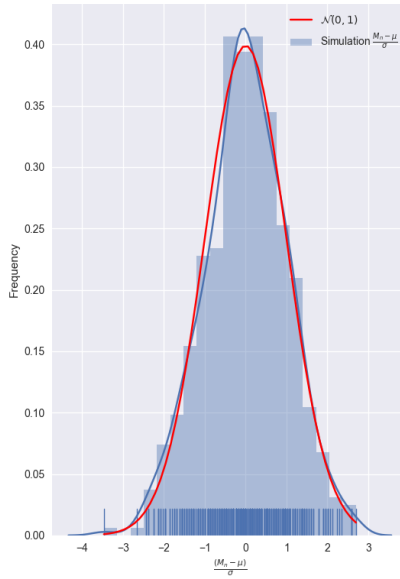


Figure 1: Each tick on the x-axis is a data point.

$$n_{\text{word}} = 500, n_{\text{exp}} = 10^6$$

1.4 Empirical normalization

Using the empirical mean (μ) and standard deviation (σ) of the dataset to normalize M_n , this is a plot of the normalized distribution, compared to the normal distribution in red :

Figure 2: Standardized M_n distribution *v.s.* normal distribution

$$n_{\text{word}} = 500, n_{\text{exp}} = 10^6$$

and its cumulative distribution function in green, compared to the normal one in red:

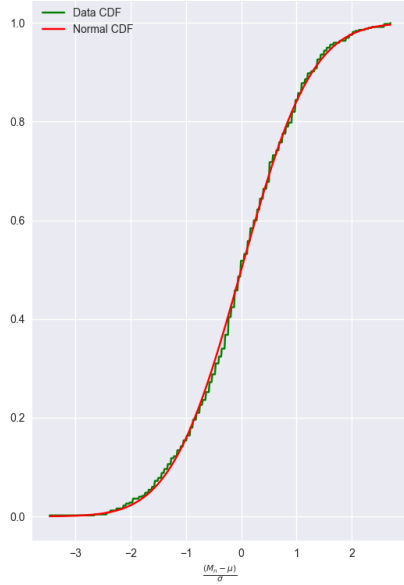


Figure 3: CDF of standardized M_n vs. normal law
 $n_{\text{word}} = 500$, $n_{\text{exp}} = 10^6$

These simulations and figures strongly indicate that the general distribution of M_n respects the central limit theorem. We now experiment with candidates for the variance of M_n : V_n

2 Validating variance candidates

2.1 A first expression

First, I was asked to review the draft of a paper (see [7]), in particular to review their own candidate for the variance of M_n . For a Markov source of order 1, with probabilities p_{00}, p_{01}, p_{10} and p_{11} :

$$V_n^{(1)} = \frac{H^3 \sigma^2 n}{\log_2^2(n)}$$

where $H = \pi_0 H_0 + \pi_1 H_1$ and $\sigma^2 = \sigma_0^2 + \sigma_1^2$ with, for $i \in \{0, 1\}$,

$$\sigma_i^2 = \frac{\pi_i p_{i0} p_{i1}}{H^3} \left(\log_2 \left(\frac{p_{i0}}{p_{i1}} \right) + \frac{H_1 - H_0}{p_{01} + p_{10}} \right)^2$$

and $H_i = -p_{i0} \log_2(p_{i0}) - p_{i1} \log_2(p_{i1})$ $H = \pi_0 H_0 + \pi_1 H_1$

I could first prove that this expression rejoins the expression of the variance of M_n for memoryless sources - which was encouraging - and I could also see, from simulations, that this variance fitted the asymptotical behavior of the empirical variance. To be exact, their difference seemed to be a sublinear function of n growing to $+\infty$ but which would be negligible compared to the first order of the variance. See Figure 4 for comparison.

2.2 Using the Frobenius eigenvalue of $P(s)$

In walking towards the end-goal of proving the Central Limit Theorem for M_n and Markov sources, a major step is finding an expression that fits V_n well

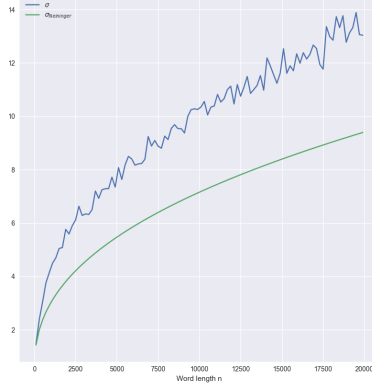


Figure 4: Variance $V_n^{(1)}$ compared to empirical variance
 $n_{\text{exp}} = 500$

enough. The candidate we saw previously which came from [6] was only defined for Markov sources of order 1 and was not proved to be a fitting candidate. Therefore, I was tasked to investigate another expression, which appeared in the analysis of P. Jacquet and my supervisor's old paper [4]. What I achieved here is to find several expressions of this variance candidate and run them against my simulations - eventually finding out what seems to be the right candidate for the variance.

First, some definitions in order to introduce this candidate, and compute the exact expressions needed for simulation:

Definition 8 We define the matrix $P(s)$ from the probability matrix P of a Markov source as:

$$\forall s \in \mathbb{C} \quad P(s) = (p_{ij}^{-s})_{(i,j) \in \{1, \dots, V\}^2}$$

Definition 9 For $s \in \mathbb{C}$, the Frobenius eigenvalue of $P(s)$ is the highest eigenvalue of $P(s)$ in absolute value. It is denoted $\lambda(s)$.

Ignoring the difficult steps it took to obtain it in [4], this is the variance expression:

$$V_n = \left(\ddot{\lambda}(-1) - \dot{\lambda}(-1)^2 \right) \frac{n}{\ln^2 n}$$

Let's compute $\ddot{\lambda}(-1)$ with a Markov chain of order 1. In [4],

$$\ddot{\lambda}(-1) = \pi \ddot{P}(-1) \psi + 2\dot{\pi}(-1) \dot{P}(-1) \psi - 2\dot{\lambda}(-1) \dot{\pi}(-1) \psi$$

However, the relations defining $\pi(s)$:

$$\begin{cases} \pi(s)P(s) &= \lambda(s)\pi(s) \\ P(s)\psi(s) &= \lambda(s)\psi(s) \\ \pi(s)\psi(s) &= \lambda(s) \end{cases}$$

did not seem to allow me to directly compute $\dot{\pi}(s)$ (it seemed like I need one more). Therefore, I computed $\lambda(s)$ as the greatest eigenvalue of $P(s)$. Let χ the characteristic polynomial of $P(s)$, and Δ its discriminant

$$\chi = (X - p_{00}^{-s})(X - p_{11}^{-s}) - (p_{01} p_{10})^{-s}$$

$$\begin{aligned} \text{and} \quad \Delta &= (p_{00}^{-s} + p_{11}^{-s})^2 - 4[(p_{00} p_{11})^{-s} - (p_{01} p_{10})^{-s}] \\ &= p_{00}^{-2s} + p_{11}^{-s} - 2(p_{00} p_{11})^{-s} + 4(p_{01} p_{10})^{-s} \end{aligned}$$

Informally, we have this expression for $\lambda(s)$ where we need to decide which sign is the correct one :

$$\lambda(s) = \frac{p_{00}^{-s} + p_{11}^{-s} \pm \sqrt{\Delta(s)}}{2}$$

Since

$$\Delta(-1) = (p_{00} + p_{11})^2 - 2p_{00}p_{11} + 4p_{01}p_{10} = (p_{00} + p_{11} - 2)^2$$

then $\sqrt{\Delta(-1)} = 2 - p_{00} - p_{11} = p_{01} + p_{10}$. Thus, picking the + sign in the former expression, we verify that

$$\lambda(-1) = \frac{p_{00} + p_{11} + \sqrt{\Delta(-1)}}{2} = 1$$

Derivating

$$\dot{\lambda}(s) = \frac{1}{2} \left(-\ln p_{00} p_{00}^{-s} - \ln p_{11} p_{11}^{-s} + \frac{\Delta'(s)}{2\sqrt{\Delta(s)}} \right)$$

$$\text{and the expression for } \Delta'(s) \quad \Delta'(s) = -2 \ln p_{00} p_{00}^{-2s} - 2 \ln p_{11} p_{11}^{-2s} + 2 \ln(p_{00}p_{11}) (p_{00} p_{11})^{-s} - 4 \ln(p_{01}p_{10}) (p_{01} p_{10})^{-s}$$

gives

$$\Delta'(-1) = -2 \ln p_{00} p_{00}^2 - 2 \ln p_{11} p_{11}^2 + 2 \ln(p_{00}p_{11}) (p_{00}p_{11}) - 4 \ln(p_{01}p_{10}) (p_{01}p_{10})$$

allowing to compute $\dot{\lambda}(-1)$. Numerically, we verified that $\dot{\lambda}(-1) = h$. Derivating again yields

$$\ddot{\lambda}(s) = \frac{1}{2} \left(\ln^2 p_{00} p_{00}^{-s} + \ln^2 p_{11} p_{11}^{-s} + \frac{\Delta''(s)\sqrt{\Delta(s)} - \Delta'(s) \cdot \Delta'(s)/2\sqrt{\Delta(s)}}{2\Delta(s)} \right)$$

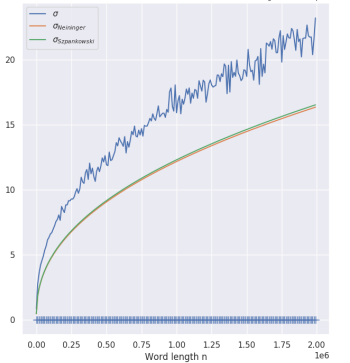
$$\text{with } \Delta''(s) = 4 \ln^2 p_{00} p_{00}^{-2s} + 4 \ln^2 p_{11} p_{11}^{-2s} - 2 \ln^2(p_{00}p_{11}) (p_{00} p_{11})^{-s} + 4 \ln^2(p_{01}p_{10}) (p_{01} p_{10})^{-s}$$

Finally,

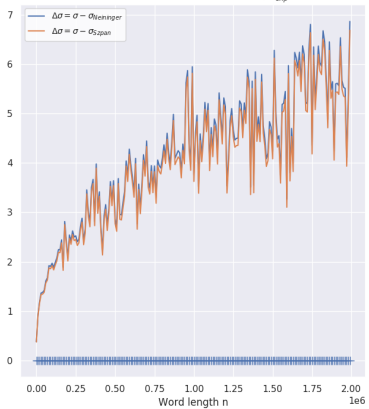
$$\ddot{\lambda}(-1) = \frac{1}{2} \left(\ln^2 p_{00} p_{00} + \ln^2 p_{11} p_{11} + \frac{\Delta''(-1)\sqrt{\Delta(-1)} - \Delta'(-1)^2/2\sqrt{\Delta(-1)}}{2\Delta(-1)} \right)$$

The simulations using this coefficient for the variance are quite good. It also seems that this formula for the variance is equivalent to the one used in the unpublished paper *Probabilistic Analysis of Lempel-Ziv Parsing for Markov Sources* by Leckey, Wormald and Neininger, but our two ways of deriving it differs. Numerical instability might account for the tiny differences found for high n values (10^7), although this hasn't been verified.

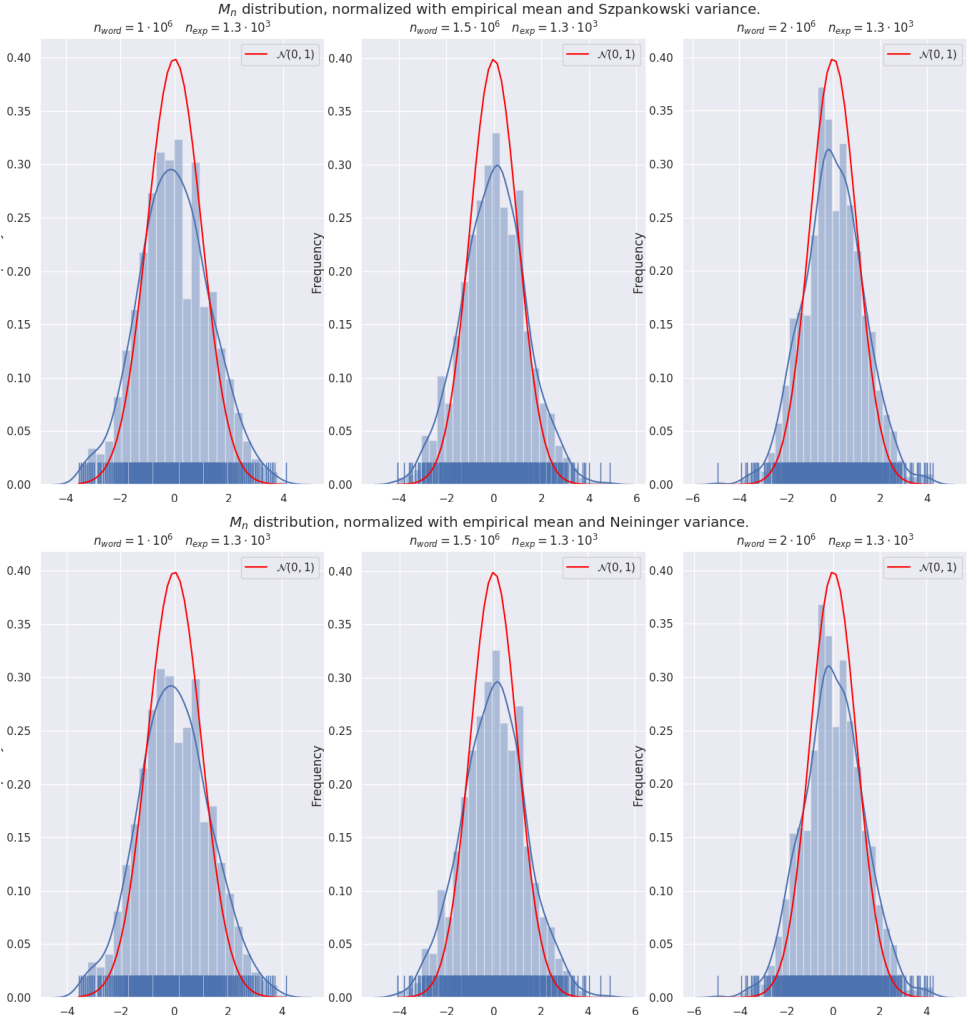
Empirical standard deviation (σ) and theoretical ones ($\sigma_{\text{Lempel-Ziv}}$, σ_{Ziv}), $n_{\text{exp}} = 400$



Difference between standard deviations, $n_{\text{exp}} = 400$



Now, for some distributions of very long words that were normalized using our theoretical standard deviations, and empirical means. The blue plot is a gaussian fit for the simulation results, which also appear as a blue histogram. The two sets of figures are identical but obtained using different expressions.



3 Conclusion

Similar results were obtained for a variety of randomly generated Markov sources, which seem to indicate that this formula for the variance could be proven theoretically correct.

Limits of this work

The figures suffer from imprecision over the computation of the empirical variance : this is due to the difficulties encountered in computing large amounts of long words (of size over 10^6). The figure in appendix B is an example of this limitation : with only 100 experiments, the empirical variance varies a lot. Possible ideas of improvement might come from parallelization, rewriting functions in a computational language such as Julia, or using/devising a datastructure specific to the task of building very long words.

Another problem is that the space of random Markov chains (here:

stochastic matrices of size 2) is not sampled thoroughly. Sampling a small number of Markov chains uniformly according to their entropy might be interesting as a representation of the space, because otherwise it would be hard to sample a large number of stochastic matrices due to the necessity of computing large words for each of them.

Finally, the difference between empirical V_n and our expression seems to be growing very slowly with n . This might be a term that is negligible (*i.e.* of order lower than $n/\log_2^2(n)$), or a small detail in the formula of V_n . For example, the natural base logarithm (in $n/\ln^2(n)$) in the eigenvalue expression works slightly better than the one in base 2 (in $n/\log_2^2(n)$), with the inverse situation happening for the first expression ($\sigma_{\text{Neininger}}$). Anyway, the figures obtained seem to indicate that we are close to the exact solution.

The code for these experiments, with detailed procedures for reproducibility, is hosted on GitHub in one of my private repositories (which I plan to make public eventually). Another (probably bad) way of computing $\lambda(s)$ is in appendix C.

III Tail symbols analysis

Covariance asymptotics

1 Markov Independent Model

Let M a Markov source and n an integer. The Markov Independent Model (MI) considers n infinite words generated by M . The choice of the starting symbol of each sequence is a parameter of the model. For example, all the sequences might start with the same letter of the alphabet $c \in \mathcal{A}$. Or the first symbol could be initialized using the stationary distribution of the Markov chain related to M .

This is an example for $n = 4$. These are the sequences :

$$\begin{aligned} X(1) &= 00000\dots \\ X(2) &= 1010101\dots \\ X(3) &= 1001101\dots \\ X(4) &= 001100111\dots \end{aligned}$$

These sequences are used to build a digital search tree by considering the shortest prefix of each sequence that has not appeared yet in the previously considered sequences.

On our example, it yields the parsed word :

$$()(0)(1)(10)(00)$$

which can be read as the DST :

INSERT TREE

2 Tail symbols

2.1 Definition

Each of the n sequences possess a tail symbol. For each sequence, its tail symbol is the character that immediately follows the prefix phrase inserted into the DST. Therefore the tail symbol is a specific character of this sequence. If

we only have the DST containing the prefix phrases, we cannot recover the tail symbols.

Visually, with the **prefix phrases** in red and the **tail symbol** in green :

$$\begin{aligned} X(1) &= 000000 \dots \\ X(2) &= 1010101 \dots \\ X(3) &= 1001101 \dots \\ X(4) &= 001100111 \dots \end{aligned}$$

Let c be a character from our alphabet $\{a, b\}$. In the case when all the sequences start with a c , we define T_n^c the *number of times a is a tail symbol in the experiment*.

2.2 Recurrence relation

For $n \geq 0$, we have :

$$T_{n+1}^c = \delta_a + \tilde{T}_{N_a}^a + \tilde{T}_{N_b}^b$$

where :

- $\delta_a = \begin{cases} 1 & \text{if } a \text{ is the tail symbol of the first sequence} \\ 0 & \text{else} \end{cases}$
- N_a is the random variable giving *the size of the left subtree which contains phrases whose second letter is a*
- $\tilde{T}_{N_a}^a$ is the number of times a is a tail symbol for the sequences that were used to build the subtree with phrases having a as second symbol.
- T_0^c for all c by convention.

If we take $\{N_a = k\}$, then $\{N_b = n - k\}$, then the count of the tail symbols on the left tree is independent of the one on the right tree : *i.e.* these quantities are conditionnaly independent.

3 Total path length

3.1 Definition

Defining L_n^c as the *total path length of the nodes of the DST that was built with MI model with n sequences starting with letter c* . It is the sum of the lengths of all the prefix phrases.

3.2 Recurrence relation

There is another recursive stochastic relation for this quantity, which is, for all $n \geq 0$:

$$L_{n+1}^c = n + \tilde{L}_{N_a}^a + \tilde{L}_{N_b}^b$$

with the convention that $L_0^c = 0$ for all c .

Same as for the number of tail symbols, this relation is found by considering the DST and its two main subtrees. Except that this time, we count the number of times an edge contributes to the path length. The root with its two nodes contributes as n . The two subtrees contribute respectively for $\tilde{L}_{N_a}^a$ and $\tilde{L}_{N_b}^b$.

It is convenient that these two quantities are conditionnaly independent in the same way as previously seen for the tail symbols.

4 Poisson tranform differential equation

4.1 Covariance recurrence relation

Using the previous recurrence relations, we have for all $n \geq 0$:

$$\text{Cov}(T_{n+1}^c, L_{n+1}^c) = \text{Cov}(\delta_a + \tilde{T}_{N_a}^a + \tilde{T}_{N_b}^b, n + \tilde{L}_{N_a}^a + \tilde{L}_{N_b}^b)$$

Since the covariance is a bilinear function which is equal to zero if its two terms are independent or if one is constant, we can ignore the term n and expand this quantity into six terms :

$$\begin{aligned} \text{Cov}(T_{n+1}^c, L_{n+1}^c) &= \text{Cov}(\delta_a, \tilde{L}_{N_a}^a) + \text{Cov}(\delta_a, \tilde{L}_{N_b}^b) + \text{Cov}(\tilde{T}_{N_a}^a, \tilde{L}_{N_a}^a) \\ &\quad + \text{Cov}(\tilde{T}_{N_a}^a, \tilde{L}_{N_b}^b) + \text{Cov}(\tilde{T}_{N_b}^b, \tilde{L}_{N_a}^a) + \text{Cov}(\tilde{T}_{N_b}^b, \tilde{L}_{N_b}^b) \end{aligned}$$

Since δ_a is given by the tail symbol of the first sequence, which is independent from the rest of the process : $\text{Cov}(\delta_a, \tilde{L}_{N_a}^a) = \text{Cov}(\delta_a, \tilde{L}_{N_b}^b) = 0$

Now, it is not obvious if the pairs $(\tilde{T}_{N_a}^a, \tilde{L}_{N_b}^b)$ and $(\tilde{T}_{N_b}^b, \tilde{L}_{N_a}^a)$ are independent or uncorrelated, because the random variable N_a is not fixed. However they are conditionnaly independent, therefore :

$$\begin{aligned} \text{Cov}(\tilde{T}_{N_a}^a, \tilde{L}_{N_b}^b) &= \sum_{k=0}^n \text{Cov}(\tilde{T}_{N_a}^a, \tilde{L}_{N_b}^b \mid N_a = k) P(N_a = k) \\ &= \sum_{k=0}^n \text{Cov}(\tilde{T}_k^a, \tilde{L}_{n-k}^b) P(N_a = k) \\ &= \sum_{k=0}^n 0 \cdot P(N_a = k) \\ &= 0 \end{aligned}$$

Samely, $\text{Cov}(\tilde{T}_{N_b}^b, \tilde{L}_{N_a}^a) = 0$. Yielding :

$$\boxed{\text{Cov}(T_{n+1}^c, L_{n+1}^c) = \text{Cov}(\tilde{T}_{N_a}^a, \tilde{L}_{N_a}^a) + \text{Cov}(\tilde{T}_{N_b}^b, \tilde{L}_{N_b}^b)}$$

4.2 Poisson transform

Defining

$$\boxed{C_c(z) = \sum_{n \geq 0} \text{Cov}(T_n^c, L_n^c) \frac{z^n}{n!} e^{-z}}$$

Computing, with $p = P(a|c)$ and $q = 1 - p$:

$$\begin{aligned} \sum_{n \geq 0} \text{Cov}(\tilde{T}_{N_a}^a, \tilde{L}_{N_a}^a) \frac{z^n}{n!} e^{-z} &= \sum_{n \geq 0} \sum_{k=0}^n P(N_a = k) \text{Cov}(\tilde{T}_{N_a}^a, \tilde{L}_{N_a}^a \mid N_a = k) \frac{z^n}{n!} e^{-z} \\ &= \sum_{n \geq 0} \sum_{k=0}^n \binom{n}{k} p^k q^{n-k} \text{Cov}(\tilde{T}_k^a, \tilde{L}_k^a) \frac{z^n}{n!} e^{-z} \end{aligned}$$

In this case, \tilde{T}_k^a and T_k^a as well as \tilde{L}_k^a and L_k^a have the same distribution, hence :

$$\begin{aligned}
\sum_{n \geq 0} \text{Cov}(\tilde{T}_{N_a}^a, \tilde{L}_{N_a}^a) \frac{z^n}{n!} e^{-z} &= \sum_{n \geq 0} \sum_{k=0}^n \binom{n}{k} p^k q^{n-k} \text{Cov}(T_k^a, L_k^a) \frac{z^n}{n!} e^{-z} \\
&= \sum_{n \geq 0} \sum_{k=0}^n \left(\frac{(zp)^k}{k!} \text{Cov}(T_k^a, L_k^a) e^{-zp} \right) \left(\frac{(zq)^{n-k}}{(n-k)!} e^{-zq} \right) \\
&= \underbrace{\left(\sum_{n \geq 0} \frac{(zp)^n}{n!} \text{Cov}(T_n^a, L_n^a) e^{-zp} \right)}_{=C_a(zp)} \underbrace{\left(\sum_{n \geq 0} \frac{(zq)^n}{n!} e^{-zq} \right)}_{=1} \\
&= C_a(zp)
\end{aligned}$$

A similar computation gives $\sum_{n \geq 0} \text{Cov}(\tilde{T}_{N_b}^b, \tilde{L}_{N_b}^b) \frac{z^n}{n!} e^{-z} = C_b(zq)$, this time conditioning on $P(N_b = k) = \binom{n}{k} q^k p^{n-k}$.

From what we've seen, when derivating $C_c(z)$ we get :

$$\begin{aligned}
\partial_z C_c(z) &= \sum_{n \geq 0} \text{Cov}(T_n^c, L_n^c) n \frac{z^{n-1}}{n!} e^{-z} - C_c(z) \\
&= \sum_{n \geq 0} \text{Cov}(T_{n+1}^c, L_{n+1}^c) \frac{z^n}{n!} e^{-z} - C_c(z) \\
&= \sum_{n \geq 0} \left[\text{Cov}(\tilde{T}_{N_a}^a, \tilde{L}_{N_a}^a) + \text{Cov}(\tilde{T}_{N_b}^b, \tilde{L}_{N_b}^b) \right] \frac{z^n}{n!} e^{-z} - C_c(z) \\
&= C_a(zp) + C_b(zq) - C_c(z)
\end{aligned}$$

Finally the equation for $C_c(z)$ is :

$$\boxed{\partial_z C_c(z) + C_c(z) = C_a(zp) + C_b(zq)}$$

IV Optimal parsing

1 LZ77 and LZ78 comparison

1.1 Practical uses of LZ77

Although it is based on the same approach of identifying previously seen phrases - but storing them in a fixed size dictionary instead of an evergrowing parse tree like in LZ78 - the LZ77 algorithm is the most used in practical implementations of universal compression algorithms.

Let us dive in on a few of the most recent ones :

- DivANS
- zStandard
-

1.2

2 Observations regarding the results and the proof

Corrections

- The formal definition of D_n^{LZ} in (3) is misleading because it contradicts the previous verbal definition. For w a word of size n , the formula

$$\frac{1}{M_n(w)} \sum_{j=0}^{M_n(w)-1} |u_j^{\text{LZ}}|$$

would rather be the empirical average length of a phrase in the Lempel-Ziv parsing of a word $\overline{D}_n^{\text{LZ}}$, whereas D_n^{LZ} is used in the rest of the paper as the length of a randomly selected phrase. These two aren't equal : if we build a Lempel-Ziv DST from a word, then D_n^{LZ} can be seen as the depth of a random node, which is different from the average path length computed on all the nodes.

- In *Remark 2*, I think the definition of v and t should rather be $v = a_{i'} \dots a_i$ and $t = a_{i+1} \dots a_n$.
- The result (14) should be an equality, and it is one indeed because of the flexible parsing algorithm. A proof by contradiction can show this. However, since we upperbound $g(j)$ by $+\infty$ in (23) there might be no purpose to using (14) instead of just (13).
- The proof around *Theorem 1* has several flaws. The notation X for a sequence depending on n and not a random variable is misleading. On the other hand, it should appear that both g and j are random variables, as the randomness of j is used in the end of the proof. I wrote some possible definitions [here](#), and applied them to make some computations that seemed otherwise incorrect because of their use of randomness outside of a probability measure ([here](#)).
- As for the arguments that link $|L_{g(j)-k}|$ to D_n^{LZ} , I have indicated how I think they could be developped in [this part](#). These arguments are the most controversial part right now I think.
- *Theorem 2* is false as stated: we proved *Theorem 1* using a random j . The randomness remains, so the quantifier 'for any $j < M_n$ ' should be removed. This would be true for *Theorem 1* too.
- The proof of *Theorem 2* may stop at (26) since we can directly prove this upperbound goes to 0. This yields a tighter upperbound for *Theorem 2*. I detailed this analysis in the [last part of this report](#).
- In that same proof, the step between (25) and (26) relies heavily on a result from [6]. A bit more context on this result (and why it does apply here) would make things clearer.
- The conclusion claims to use *Cramer's* theorem to link

$$\max_{0 \leq k \leq g_W(J)} \{L_{g_W(J)-k} - k\}$$

to D_n^{FLEX} , which is a sum of random variables. Since *Cramer* applies to independent random variables and the lengths of successive phrases are not independent, something must be missing there.

3 Definitions

3.1 Definitions

Notations

These are definitions and notations in order to write the proof of *Theorem 1*.

Definition 10 For all $n \in \mathbb{N}$, calling Ω_n the set of words of length n .

Definition 11 Defining $W \in \Omega_n$ to be a random variable which outputs words of length n from a memoryless source.

Definition 12 Considering $J \in \mathbb{N}^*$ to be a random variable which, in the event $\{W = w\}$, uniformly randomly picks the index of one of the phrases of w . The joint law of J with W being :

$$PW = w, J = j = \begin{cases} 1/M_n(w) & \text{if } j \leq M_n(w) \\ 0 & \text{else} \end{cases}$$

Remark 9 We might choose another randomness for J , but this one seems more natural.

Definition 13 For a given word $w \in \Omega_n$, and for all $i \in \mathbb{N}^*$, we consider $g_w(i)$ defined by

$$g_w(i) = f_w(i) + |L_{f(i)}|$$

where $f_w(i)$ is the starting index of the i^{th} phrase of the flexible parsing of w , and $|L_{f_w(i)}|$ is the length of the longest greedy phrase given by the Lempel-Ziv parsing of this same word.

Definition 14 We define $g_w(J)$ to be the random variable which outputs $g_w(j)$ during the events $\{W = w\}$ and $\{J = j\}$.

Definition 15 For all $k \in \mathbb{N}$, let $L_{g_w(J)-k}$ be the random variables which gives the $(k+1)^{\text{th}}$ possible phrase for the flexible parsing at index $g_w(J) - k$. Its only randomness comes from W and J . If $i \leq 0$, we might assume that L_i will be the empty word of size 0.

Notation 4 Denoting by

$$B_{J,W}^k = |L_{g_w(J)-k}|$$

the length of this $(k+1)^{\text{th}}$ candidate.

We can now study the random variable

$$\max_{0 \leq k \leq g_w(J)} \{B_{J,W}^k - k\}$$

Given any $(j, w) \in \mathbb{N}^* \times \Omega_n$, under the events $\{J = j\}$, $\{W = w\}$ and $\{J \leq M_n(W)\}$, this random variable is the length of the j^{th} phrase of the flexible parsing of the word w .

Definition 16 Defining the random variable D_n^{FLEX} as

$$\begin{aligned} D_n^{\text{FLEX}}(w) &= \frac{1}{M_n(w)} \sum_{j=0}^{M_n(w)-1} |u_j^{\text{FLEX}}(w)| \\ &= \frac{1}{M_n(w)} \sum_{j=0}^{M_n(w)-1} \max_{0 \leq k \leq g_w(j)} \{B_{j,w}^k - k\} \end{aligned}$$

where $D_n^{\text{FLEX}}(w)$ is the empirical average of the lengths of the phrases of the flexible parsing of a word w , contrary to $\max_{0 \leq k \leq g_w(J)} \{B_{J,w}^k - k\}$, with J a random variable and w fixed, which is the length of a uniformly randomly selected phrase of the flexible parsing of w .

Definition 17 We denote x_n the average value of D_n^{LZ} :

$$x_n = \frac{1}{h} \log_2 \left(\frac{nh}{\log_2(n)} \right)$$

Remark 10 The random variable D_n^{LZ} is the length of a phrase randomly taken from the Lempel-Ziv parsing of a memoryless-generated word. It is not the same as the empirical average length of a phrase, which we can denote, for $w \in \Omega_n$:

$$\overline{D}_n^{\text{LZ}}(w) = \frac{1}{M_n(w)} \sum_{j=0}^{M_n(w)-1} |u_j^{\text{FLEX}}|$$

Notation 5 We denote b_n^δ as

$$b_n^\delta = x_n + (c_3 x_n)^\delta$$

We will study

$$\mathbb{P} \max_{0 \leq k \leq g_W(J)} \{B_{J,W}^k - k\} > b_n^\delta$$

3.2 Clues in proving the result

Let $k \in \mathbb{N}$. Currently, the proof to show this stands on three arguments :

- (1) The first is that $L_{g_W(J)-k}$ is a random phrase from the Lempel-Ziv parsing of a word of length N , where $N \leq g_W(J) \leq n$. In the event $\{N = n'\}$, we consider $D_{n'}^{LZ}$.
- (2) The second, is that $|L_{g_W(J)-k}|$ can therefore be considered the same as D_N^{LZ} i.e at least equal in law.
- (3) The third is that, for all $n' \leq n$, $D_n'^{LZ} \leq D_n^{LZ}$.

Although they seem generally true, there are different problems with each of these arguments :

- N isn't clearly established, so D_N^{LZ} isn't really known.
- If we can identify N and, let's say, condition our probability with $\{N = n'\}$, it is not obvious that the choice of a phrase at position $g_W(J) - k$ knowing $\{N = n'\}$ is the same as the uniform choice that operates when choosing a random phrase from a word of size n' , in $D_{n'}^{LZ}$.
- As for (3), this result is true on average, but not in all cases. Indeed, since D_n^{LZ} (resp. $D_n'^{LZ}$) is concentrated around x_n (resp. $x_{n'}$), and $x_{n'} < x_n$ since $n' < n$, we can show that this result holds with high probability. To write the proof, we may condition using events of the type

$$\{|D_n^{LZ} - x_n| \leq k_n^{(1)} v_n\}$$

$$\text{and} \quad \{|D_{n'}^{LZ} - x_{n'}| \leq k_n^{(2)} v_{n'}\}$$

where

$$v_n = \sqrt{\log(nh/\log(n))}$$

A sketch of the proof is to apply concentration inequalities to these events while picking $(k_n^{(1)}, k_n^{(2)})$ such that

$$k_n^{(1)} v_n + k_n^{(2)} v_{n'} < x_n - x_{n'}$$

and having $k_n^{(1)}$ and $k_n^{(2)}$ go to $+\infty$ for n going to $+\infty$ in order for the upperbound probability to converge to zero.

3.3 Working with these probabilistic objects

Computations

With these definitions, we can do the formal computations at the beginning of the proof of *Theorem 1*. By conditionning on W and J , we obtain

$$\begin{aligned}
P_{0 \leq k \leq g_W(J)} \max \{B_{J,W}^k - k\} > b_n^\delta &= \sum_{w \in \Omega_n} \sum_{j \in \mathbb{N}^*} P_{0 \leq k \leq g_W(J)} \max \{B_{J,W}^k - k\} > b_n^\delta, W = w, J = j \\
&= \sum_{w \in \Omega_n} \sum_{j \in \mathbb{N}^*} P \bigcup_{k=0}^{g_w(j)} \{B_{w,j}^k > k + b_n^\delta\}, W = w, J = j \\
&\leq \sum_{w \in \Omega_n} \sum_{j \in \mathbb{N}^*} \sum_{k=0}^{g_w(j)} P B_{w,j}^k > k + b_n^\delta, W = w, J = j \\
&\leq \sum_{w \in \Omega_n} \sum_{j \in \mathbb{N}^*} \sum_{k=0}^{+\infty} P B_{w,j}^k > k + b_n^\delta, W = w, J = j \\
&= \sum_{k=0}^{+\infty} \sum_{w \in \Omega_n} \sum_{j \in \mathbb{N}^*} P B_{w,j}^k > k + b_n^\delta, W = w, J = j \\
&= \sum_{k=0}^{+\infty} P B_{W,J}^k > k + b_n^\delta
\end{aligned}$$

For all $k \in \mathbb{N}$, we may now prove that

$$P B_{J,W}^k > k + b_n^\delta \leq P D_n^{LZ} > k + b_n^\delta$$

4 Proof of an upperbound for the paper

Upperbound proof

This is a proof that the upperbound in (26) goes to 0. Assuming

$$\sum_{k=0}^{+\infty} P D_n^{LZ}(W) > k + b_n^\delta \leq A \alpha^{(c_3 x_n)^{\delta-1/2}} \sum_{i=0}^{+\infty} \alpha^{i/\sqrt{c_3 x_n}} \quad (26)$$

We can prove directly that the upperbound term goes to 0 for n going to $+\infty$, without resorting to another majoration. Since $\sqrt{c_3 x_n}$ goes to infinity for $n \rightarrow +\infty$, we can pick n such that $\sqrt{c_3 x_n} > 1$. Therefore $\alpha^{1/\sqrt{c_3 x_n}} < 1$ and the geometric sum gives

$$\sum_{i=0}^{+\infty} \alpha^{i/\sqrt{c_3 x_n}} = \frac{1}{1 - \alpha^{1/\sqrt{c_3 x_n}}}$$

It remains to prove that

$$\lim_{n \rightarrow +\infty} \frac{\alpha^{(c_3 x_n)^{\delta-1/2}}}{1 - \alpha^{1/\sqrt{c_3 x_n}}} = 0$$

which is done by using L'Hospital's rule. Define :

$$f: \begin{cases} \mathbb{R}_+^* \longrightarrow \mathbb{R} \\ x \longmapsto \alpha^{x^{\delta-1/2}} \end{cases} \quad \text{and} \quad g: \begin{cases} \mathbb{R}_+^* \longrightarrow \mathbb{R} \\ x \longmapsto 1 - \alpha^{\frac{1}{\sqrt{x}}} \end{cases}$$

Let $x \in]0; +\infty[$. Derivating yields :

$$f'(x) = \ln \alpha \left(\delta - \frac{1}{2} \right) x^{\delta-3/2} f(x) \quad \text{and} \quad g'(x) = \ln \alpha \frac{1}{2x\sqrt{x}} \alpha^{\frac{1}{\sqrt{x}}}$$

We proceed to show that $\frac{f'(x)}{g'(x)}$ goes to 0 as x goes to $+\infty$. We have

$$\frac{f'(x)}{g'(x)} = \frac{\left(\delta - \frac{1}{2} \right) x^{\delta-3/2} \alpha^{x^{\delta-1/2}}}{\frac{1}{2x\sqrt{x}} \alpha^{\frac{1}{\sqrt{x}}}} = 2 \left(\delta - \frac{1}{2} \right) x^\delta \cdot \frac{\alpha^{x^{\delta-1/2}}}{\alpha^{\frac{1}{\sqrt{x}}}}$$

Since $\alpha^{\frac{1}{\sqrt{x}}} \xrightarrow{x \rightarrow +\infty} 1$, we are left to study $x^\delta \alpha^{x^{\delta-1/2}}$.

Writing $x^\delta \alpha^{x^{\delta-1/2}} = e^{\delta \log x + \log \alpha \cdot x^{\delta-1/2}}$

and taking the log, since $\delta > 1/2$ and $\log \alpha < 0$ we see that

$$\delta \log x + \log \alpha \cdot x^{\delta-1/2} \xrightarrow{x \rightarrow +\infty} -\infty$$

Therefore $x^\delta \alpha^{x^{\delta-1/2}} \xrightarrow{x \rightarrow +\infty} 0$, and given that $f(0) = g(0) = 0$, L'Hospital's rule applies, proving that $\frac{f(x)}{g(x)} \xrightarrow{x \rightarrow +\infty} 0$.

5 Research at Purdue

During my internship, I worked in the Felix Haas building which notably hosts the Center for Science of Information - a research lab dedicated to information theory problematics. Early after my arrival, I started collaborating with some of my supervisor's contacts overseas: namely Philippe Jacquet in Paris Nokia Labs, and an Italian research team in Palermo - but I also had the pleasure to discuss, work and spend some time with several postdoctoral researchers at Purdue, working on Graph Dependency Theory (Abram Magner), Information Theory of Communication (Arun Padakandla) and other interesting subjects. Thanks to them I acquired precious insights into the academic environment in the United States. Speaking of which : I am sincerely and very thankful to the CSOI and its administration for all their help in hosting and organizing my trip. Any trouble I could have had was considered and I was able to focus on research and writing in very good and peaceful conditions. Overall, working at Purdue was a great opportunity which I would recommend to anyone wanting to expand their mind about academia and its power to connect people overseas. It has stirred my interest about doing research abroad - in particular in the U.S., and I look forward for another occasion to do so.

References

- [1] Philippe Jacquet and Wojciech Szpankowski. Towards LZ78 Analysis for Markov Sources: Distribution of Tail Symbols in DST. *Unpublished*.
- [2] Philippe Jacquet and Wojciech Szpankowski. On the Limiting Distribution of Lempel-Ziv'78 Redundancy for Memoryless Sources. *IEEE Transactions on Information Theory*, 60(11):6917–6930, 2014.
- [3] Philippe Jacquet and Wojciech Szpankowski. *Analytic Pattern Matching*. Cambridge University Press, Cambridge, 2015.
- [4] Philippe Jacquet, Wojciech Szpankowski, and Jing Tang. Average profile of the Lempel-Ziv parsing scheme for a Markovian source. *Algorithmica*, 31(3):318–360, 2001.
- [5] Alessio Langiu, Francesca Marzi, Filippo Mignosi, and Massimo Tivoli. On the Speed of Convergence of LZ78 with Optimal Parsing. *Unpublished*, 2018.
- [6] Kevin Leckey, Ralph Neininger, and Wojciech Szpankowski. A Limit Theorem for Radix Sort and Tries with Markovian Input. *arXiv preprint arXiv:1505.07321*, 2015.
- [7] Kevin Leckey, Nicholas Wormald, and Ralph Neininger. Probabilistic Analysis of Lempel-Ziv Parsing for Markov Sources. *Unpublished*, 2018.

- [8] Guy Louchard and Wojciech Szpankowski. Average profile and limiting distribution for a phrase size in the Lempel-Ziv parsing algorithm. *IEEE Transactions on Information Theory*, 41(2):478–488, 1995.
- [9] C E Shannon. A Mathematical Theory of Communication. page 55.

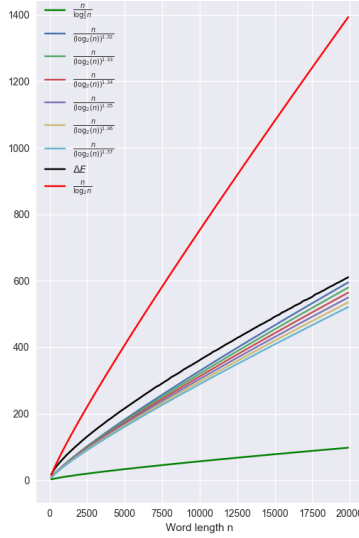


Figure 5: Difference $\mu - \frac{nh}{\log_2(n)}$ and approximation plots.

Appendices

A Asymptotics of the mean $E[M_n]$

I also tried to normalize M_n using theoretical expressions of the mean and variance. For the mean, the first order expression

$$E_n \sim \frac{nh}{\log_2(n)}$$

is, under $n \leq 10^6$, not sufficient to center the distribution. I conducted a numerical analysis of the difference between this expression and the empirical mean for growing values of n . In particular, here is how their difference, in black, compares with different approximation functions

This is not troubling as it was already predicted in the formula:

$$E_n = \frac{nh}{\log_2(n)} + \mathcal{O}\left(\frac{n}{\log_2(n)}\right)$$

B Much longer words

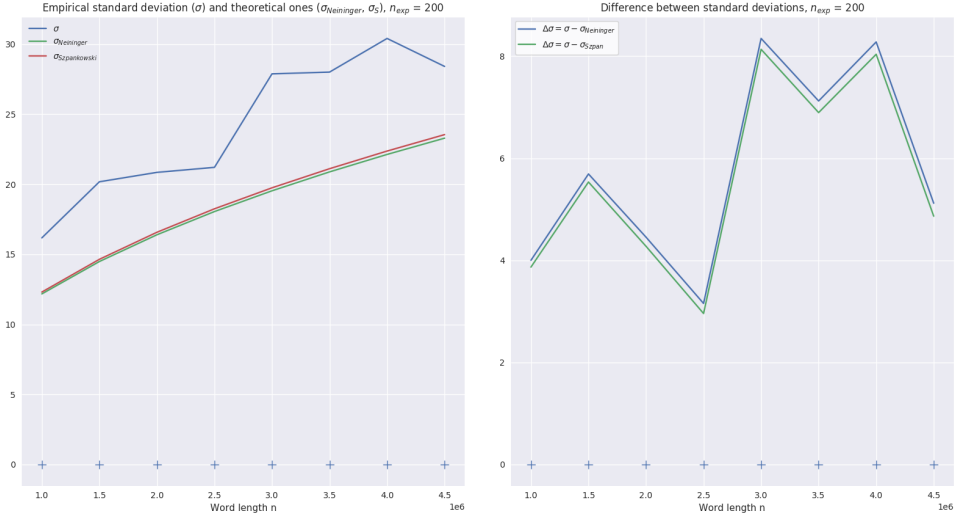


Figure 6: Standard deviations behaviors and difference with empirical
 $10^6 \leq n_{\text{word}} \leq 5 \cdot 10^6, n_{\text{exp}} = 200$



Figure 7: Standard deviations behaviors and difference with empirical
 $10^7 \leq n_{\text{word}} \leq 2 \cdot 10^7, n_{\text{exp}} = 100$

C Another (more complicated) computation of $\ddot{\lambda}(-1)$

This expression gives the same numerical results as the first one, but is more complex to compute for no apparent gain other than having yet another similar way of computing $\ddot{\lambda}(-1)$. Computing $\delta(s)$, a complex root of $\Delta(s)$, writing Δ as :

$$\begin{aligned}\Delta &= \underbrace{p_{00}^{-2 \operatorname{Re}(s)} \cos(2 \ln(p_{00}) \operatorname{Im}(s))}_{a_0(s)} \\ &+ \underbrace{p_{11}^{-2 \operatorname{Re}(s)} \cos(2 \ln(p_{11}) \operatorname{Im}(s))}_{a_1(s)} \\ &\quad \underbrace{-2(p_{00} p_{11})^{-\operatorname{Re}(s)} \cos(\ln(p_{00} p_{11}) \operatorname{Im}(s))}_{a_2(s)} \\ &+ \underbrace{4(p_{01} p_{10})^{-\operatorname{Re}(s)} \cos(\ln(p_{01} p_{10}) \operatorname{Im}(s))}_{a_3(s)} \\ &+ i \operatorname{Im}(\Delta)\end{aligned}$$

where $\operatorname{Im}(\Delta) = b_0(s) + b_1(s) + b_2(s) + b_3(s)$, with each $b_i(s)$ being the same term as $a_i(s)$ with \cos replaced by \sin . Writing

$$\Delta = \alpha(s) + i\beta(s)$$

and searching for $\delta = x(s) + iy(s)$, meaning that

$$\begin{cases} x^2 - y^2 &= \alpha \\ 2xy &= \beta \\ x^2 + y^2 &= \sqrt{\alpha^2 + \beta^2} \end{cases}$$

This yields

$$\begin{cases} x &= \pm \sqrt{\frac{1}{2}(\sqrt{\alpha^2 + \beta^2} + \alpha)} \\ y &= \pm \sqrt{\frac{1}{2}(\sqrt{\alpha^2 + \beta^2} - \alpha)} \end{cases}$$

and since $2xy = \beta$, there is $\varepsilon \in \{-1, 1\}$ such that

$$\delta = \pm(x + i\varepsilon y)$$

so
$$\lambda(s) = \frac{p_{00}^{-s} + p_{11}^{-s} \pm (x + i\varepsilon y)}{2}$$

i.e.
$$\ddot{\lambda}(-1) = \frac{p_{00} \ln^2(p_{00}) + p_{11} \ln^2(p_{11}) \pm (\ddot{x}(-1) + i\varepsilon \ddot{y}(-1))}{2}$$

where we'll have to find what is ε and which sign to pick.

But first, computing the derivatives of $x(s) = \sqrt{f(s)}$:

$$\dot{x}(s) = \frac{f'(s)}{2x(s)}$$

and
$$\ddot{x}(s) = \frac{f''(s)x(s) - f'(s) \cdot \frac{f'(s)}{2x(s)}}{2x^2(s)}$$

and then computing $f(s)$:

$$f(s) = \frac{1}{2}(\sqrt{\alpha^2 + \beta^2} + \alpha)$$

$$f'(s) = \frac{1}{2} \left[\underbrace{\frac{\dot{\alpha}\alpha + \dot{\beta}\beta}{\sqrt{\alpha^2 + \beta^2}}}_{\kappa(s)} + \dot{\alpha} \right]$$

with $\dot{\alpha} = \dot{a}_0 + \dot{a}_1 + \dot{a}_2 + \dot{a}_3$

As for $f''(s)$, it is

$$f''(s) = \frac{1}{2} \left[\frac{\dot{\gamma}(s)\kappa(s) - \gamma(s)\dot{\kappa}(s)}{\kappa^2(s)} + \ddot{\alpha}(s) \right]$$

with $\dot{\gamma}(s) = \ddot{\alpha}\alpha + \dot{\alpha}^2 + \ddot{\beta}\beta + \dot{\beta}^2$

$$\dot{\kappa}(s) = \frac{2\alpha\dot{\alpha} + 2\beta\dot{\beta}}{2\sqrt{\alpha^2 + \beta^2}}$$

Derivating according to s amounts to derivating according to $\text{Re}(s)$, so in $s = -1$:

$$\dot{\alpha}(-1) = -2 \ln p_{00} a_0(-1) - 2 \ln p_{11} a_1(-1) - \ln q_0 a_2(-1) - \ln q_1 a_3(-1)$$

and

$$\ddot{\alpha}(-1) = 4 \ln^2 p_{00} a_0(-1) + 4 \ln^2 p_{11} a_1(-1) + \ln^2 q_0 a_2(-1) + \ln^2 q_1 a_3(-1)$$

At this point we have fully determined $\ddot{x}(s)$, and we realize two things:

1. In $s = -1$, since $\text{Im}(-1) = 0$ and because of the sinus function, all the β terms, including derivatives, are equal to 0. This will simplify the expression for $\ddot{x}(-1)$.
2. Furthermore, it also means that $\ddot{y}(-1) = 0$, so

$$\boxed{\ddot{\lambda}(-1) = \frac{p_{00} \ln^2(p_{00}) + p_{11} \ln^2(p_{11}) + \ddot{x}(-1)}{2}}$$

where the $+$ comes from the fact that $\lambda(s)$ is the highest eigenvalue (and $\ddot{x}(-1) > 0$, so by continuity the expression around $s = -1$ retained the same sign)

The final expression of $\ddot{\lambda}$ (as well as $\dot{\lambda}(-1)$) can be fully expressed with $\alpha(-1)$, $\dot{\alpha}(-1)$ and $\ddot{\alpha}(-1)$. I empirically verified that $\dot{\lambda}(-1) = h$, and the final result is the same as with the first method of computation.