

Numerical simulations of LZ78 for Markovian sources

Simulation

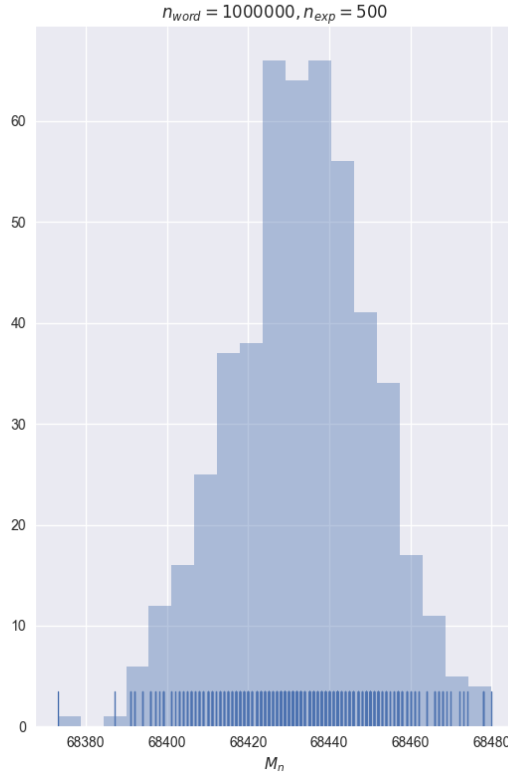
This document presents the different graphics I obtained during the following experimental process :

- Generating a random Markov chain of size 2 of matrix

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

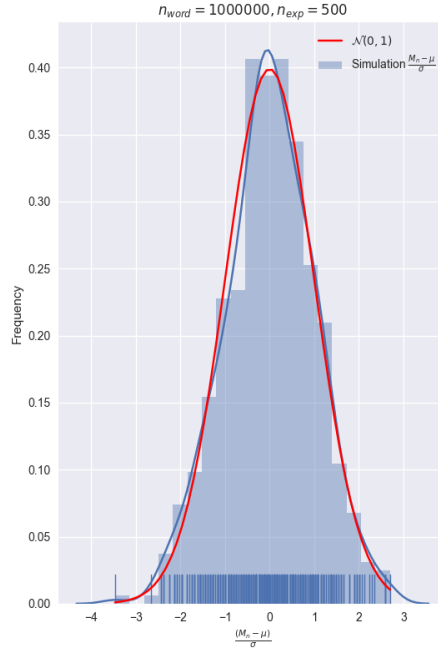
- Generating $n_{\text{exp}} \sim 500$ words of length n (or n_{word}), with $n \leq 10^6$
- Applying LZ78 on each of these words to estimate, for each n , the number of phrases M_n . A simple histogram of these values can be seen in figure 1.
- From this sampling of the random variable M_n and other parameters such as the entropy of the Markov chain, computing
 - the empirical mean (μ) and the empirical variance (σ^2)
 - different theoretical expressions of the mean and 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 mean and variance by plotting their differences for a large range of values of n , and a constant number of experiments n_{exp} .

This histogram represents the counts of the different values taken by M_n for $n = 10^6$. Each tick on the x-axis is a data point.

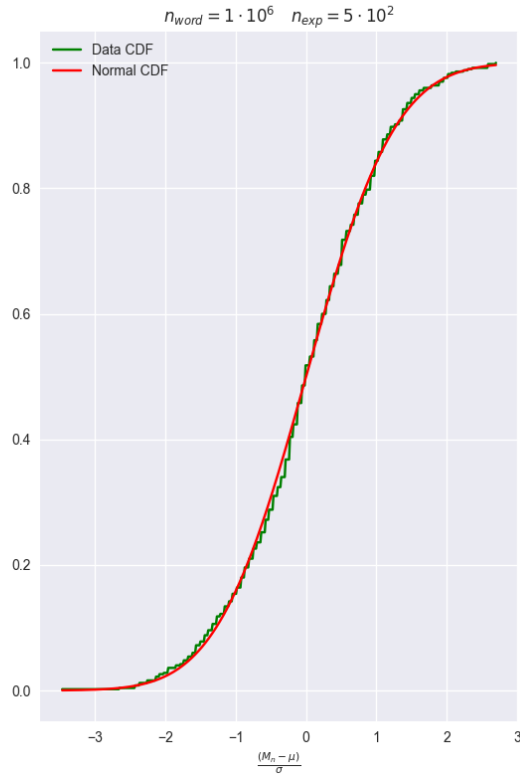


Empirical normalization

Using the empirical mean (μ) and variance (σ^2) of the dataset to normalize M_n , this is a plot of the normalized distribution, compared to the normal distribution in red :



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

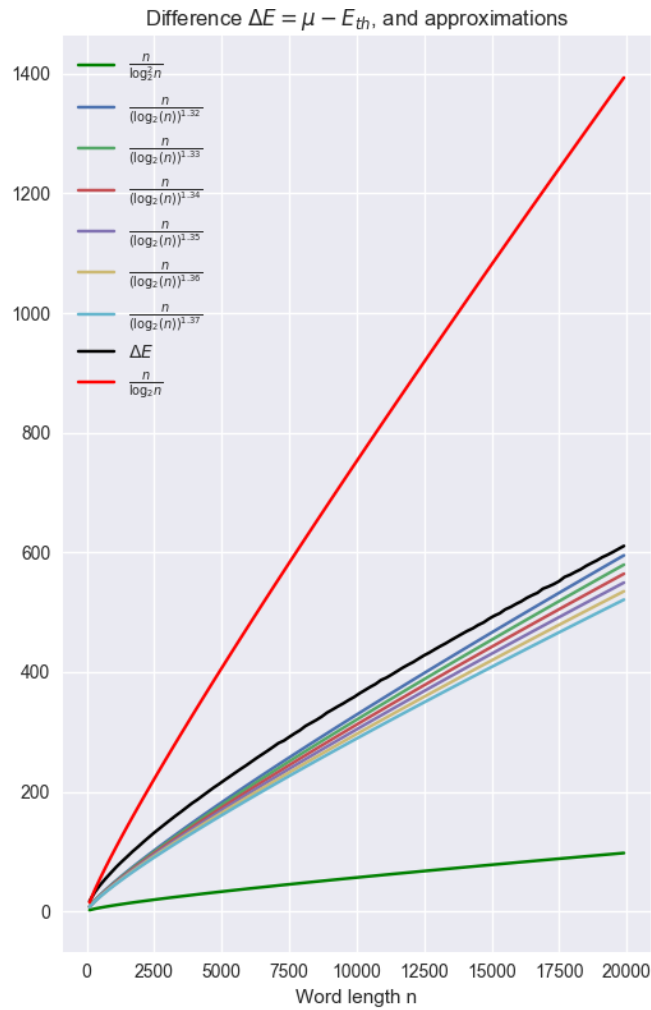


Theoretical mean

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 :

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

Theoretical variance

For the variance, I tried to use the expression of $\frac{H^3 \sigma^2}{n \log_2^2(n)}$ from K. Leckter, N. Wormald and R. Neininger's paper *Probabilistic Analysis of Lempel-Ziv Parsing for Markov Sources* :

$$\sigma^2 = \sigma_0^2 + \sigma_1^2$$

where

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

with

$$\pi_0 = \frac{p_{10}}{p_{10} + p_{01}} \quad \pi_1 = \frac{p_{01}}{p_{10} + p_{01}}$$

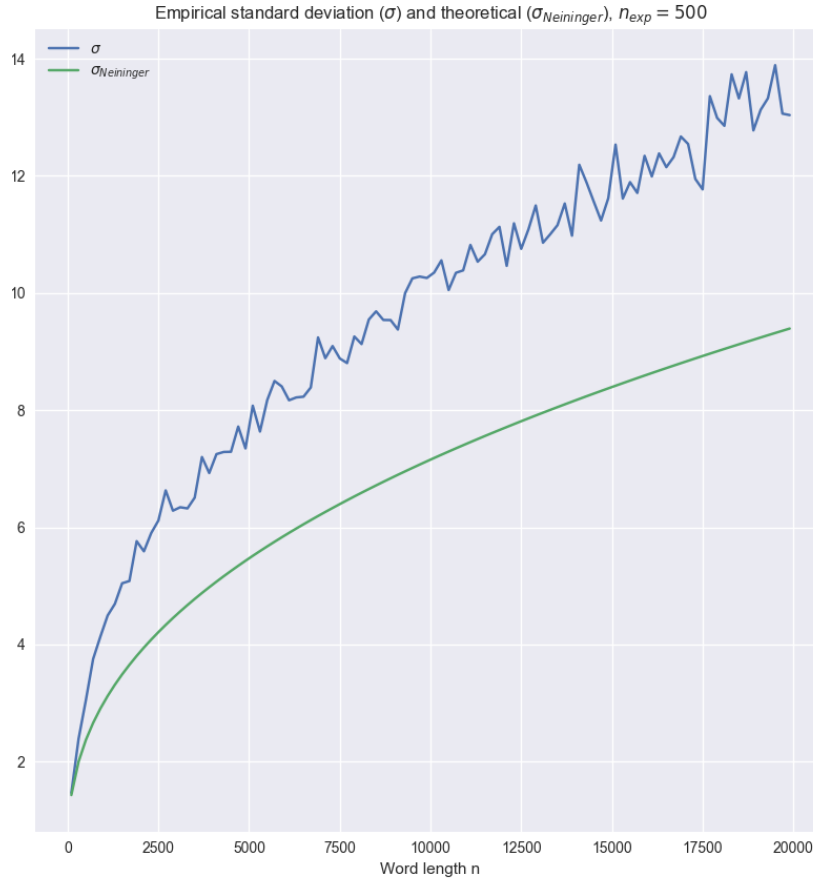
and

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

The first term in the squared part of σ_i^2 accounts for the expression of the variance for memoryless sources :

$$\begin{aligned} p_{i0} p_{i1} \log_2^2 \left(\frac{p_{i0}}{p_{i1}} \right) &= p_{i0} \log_2^2(p_{i0}) + p_{i1} \log_2^2(p_{i1}) - (-p_{i0} \log_2(p_{i0}) - p_{i1} \log_2(p_{i1}))^2 \\ &= h_2 - h^2 \end{aligned}$$

It seems, from simulations, that this variance is too small and doesn't catch up with the empirical variance. Here is how they compare when plotted together :



It seems, at first glance, that the increase would asymptotically be simply logarithmic



Now, I'm trying to compute the variance using the formula from Jacquet and Szpankowski, *Average profile of the Lempel-Ziv parsing scheme for a markovian source*, using the formula for the variance V_n :

$$V_n = \frac{1}{h^3} \left(-\frac{\beta}{\omega} - \frac{2}{\omega} \pi \dot{Q}^* \psi - h^2 \right) \log(m)$$

This formula is obtained in the Markov independent model, so m is the number of sequences with which we build a DST. Therefore, in my case, I would take

$$m \sim \frac{nh}{\ln n}$$

I computed the other terms from the general case as follows:

$$\omega = \det \begin{pmatrix} 1 & -p_{01} \\ 1 & 1 - p_{11} \end{pmatrix} = (1 - p_{11}) + p_{01}$$

And since
$$Q(s) = \begin{pmatrix} 1 - p_{00}^{-s} & 1 - p_{01}^{-s} \\ 1 - p_{10}^{-s} & 1 - p_{11}^{-s} \end{pmatrix}$$

then
$$Q'(s) = \begin{pmatrix} p_{00}^{-s} \ln p_{00} & p_{01}^{-s} \ln p_{01} \\ p_{10}^{-s} \ln p_{10} & p_{11}^{-s} \ln p_{11} \end{pmatrix}$$

and
$$Q''(s) = \begin{pmatrix} -p_{00}^{-s} \ln^2 p_{00} & -p_{01}^{-s} \ln^2 p_{01} \\ -p_{10}^{-s} \ln^2 p_{10} & -p_{11}^{-s} \ln^2 p_{11} \end{pmatrix}$$

hence
$$\det Q''(s) = (p_{00} p_{11})^{-s} \ln^2 p_{00} \cdot \ln^2 p_{11} - (p_{01} p_{10})^{-s} \ln^2 p_{01} \cdot \ln^2 p_{10}$$

therefore
$$\beta = [\det Q''(s)]_{s=-1} = p_{00} p_{11} \ln^2 p_{00} \cdot \ln^2 p_{11} - p_{01} p_{10} \ln^2 p_{01} \cdot \ln^2 p_{10}$$

After that, with

$$\mathbf{Q}^*(s) = \begin{pmatrix} 1 - p_{11}^{-s} & -(1 - p_{01}^{-s}) \\ -(1 - p_{10}^{-s}) & 1 - p_{00}^{-s} \end{pmatrix}$$

which gives

$$\dot{\mathbf{Q}}^*(s) = \begin{pmatrix} p_{11}^{-s} \ln p_{11} & -p_{01}^{-s} \ln p_{01} \\ -p_{10}^{-s} \ln p_{10} & p_{00}^{-s} \ln p_{00} \end{pmatrix}$$

then $\pi \dot{\mathbf{Q}}^* \psi = \pi_0 p_{11} \ln p_{11} - \pi_1 p_{10} \ln p_{10} - \pi_0 p_{01} \ln p_{01} + \pi_1 p_{00} \ln p_{00}$
--

However, this does not work yet since I obtain negative values when computing the coefficient

$$-\frac{\beta}{\omega} - \frac{2}{\omega} \pi \dot{\mathbf{Q}}^* \psi - h^2$$

These are the values I obtain for 10 random Markov chains :

	$\frac{1}{h^3}$	$\frac{\beta}{\omega}$	$\frac{2}{\omega} \pi \dot{\mathbf{Q}}^* \psi$	h^2	$\ln(m)$	V_n
0	5.564964	0.126471	-0.118784	0.318438	3.815542	-6.924727
1	11.551129	0.011461	0.293426	0.195697	3.572111	-20.655011
2	4.760918	0.065130	-0.003900	0.353351	3.867559	-7.633738
3	5.581253	0.082544	-0.004799	0.317818	3.814568	-8.421570
4	61.445820	-0.844743	0.109159	0.064220	3.014987	124.375693
5	3.287099	-0.068752	0.088698	0.452333	3.991037	-6.195797
6	6.472210	0.153696	-0.166402	0.287938	3.765200	-6.707164
7	9.034341	-0.082809	-0.081034	0.230534	3.654028	-2.201585
8	4.096312	0.050649	-0.026543	0.390605	3.917676	-6.655298
9	8.840450	-0.016878	-0.043772	0.233893	3.661260	-5.607404

The problem probably comes from my computation, in the case of a Markov chain of size 2, of the values of ω , β or $\dot{\mathbf{Q}}^*$. However, I did follow the general formulas from the paper, so this is intriguing.