

Numerical simulations of LZ78 for Markovian sources

Simulation

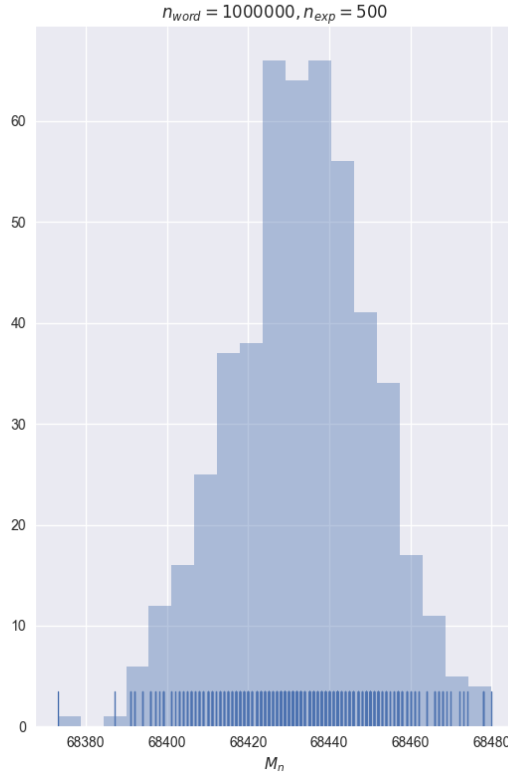
This document presents the different graphics I obtained during the following experimental process :

- Generating a random Markov chain of size 2 of matrix

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

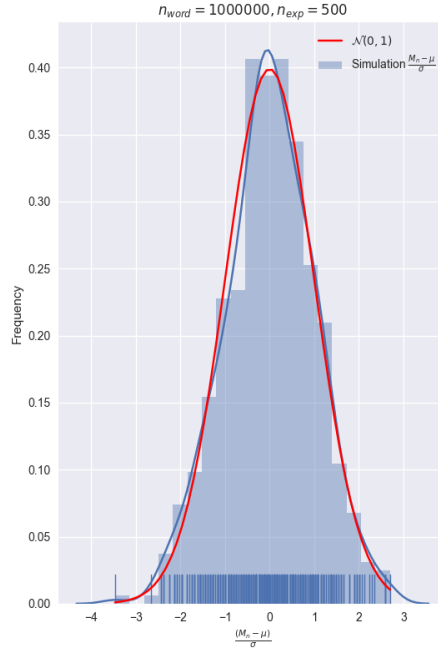
- Generating $n_{\text{exp}} \sim 500$ words of length n (or n_{word}), with $n \leq 10^6$
- Applying LZ78 on each of these words to estimate, for each n , the number of phrases M_n . A simple histogram of these values can be seen in figure 1.
- From this sampling of the random variable M_n and other 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.

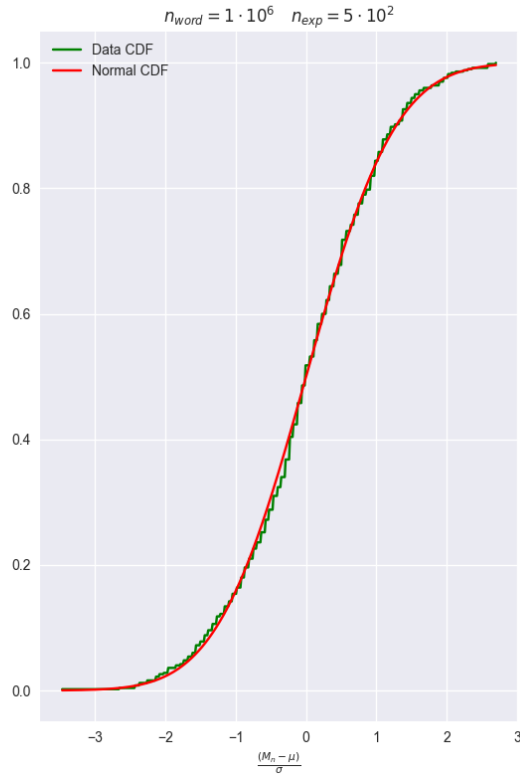


Empirical normalization

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 :

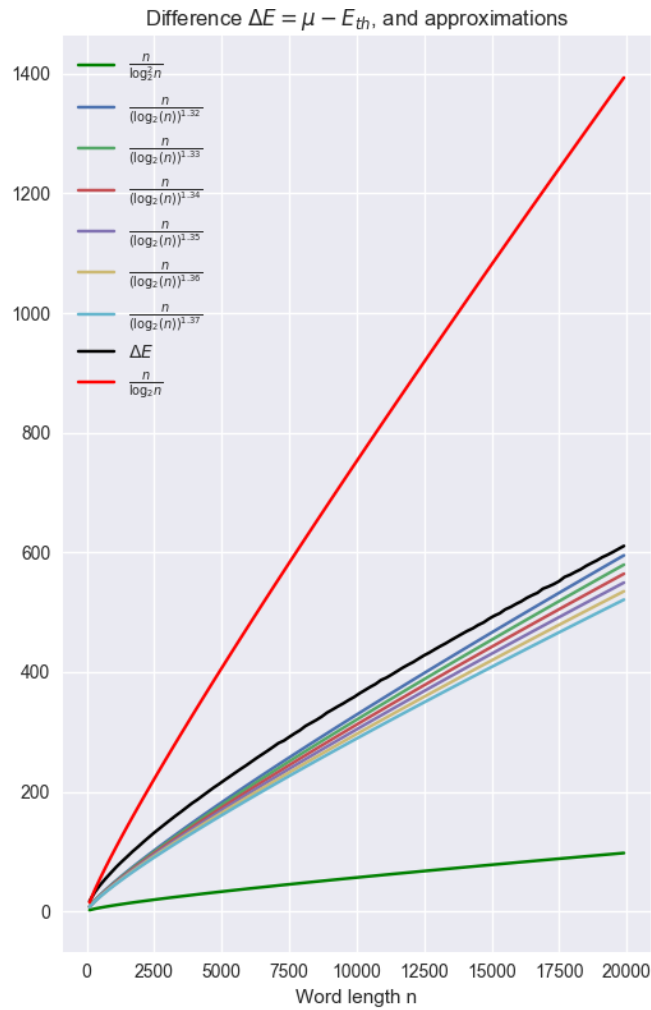


Theoretical mean

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

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

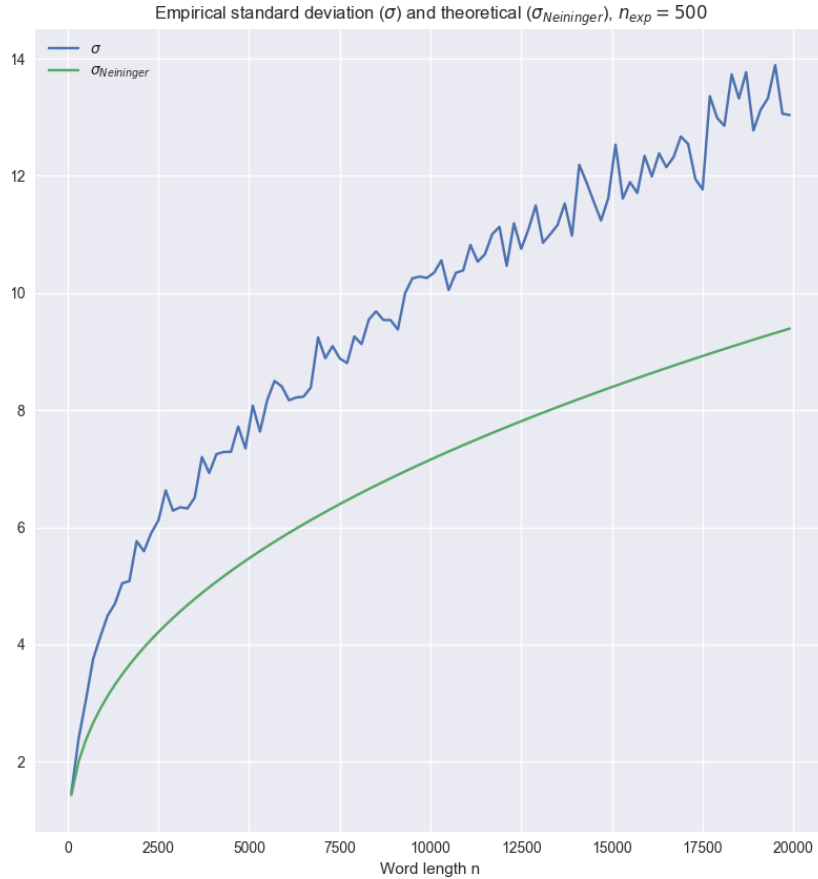
and

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

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

$$m \sim \frac{nh}{\log(n)}$$

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

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

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

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

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

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

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

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

Alternative representation

Another expression for the variance might be

$$\frac{\ddot{\lambda}(-1) - \dot{\lambda}(-1)^2}{\dot{\lambda}(-1)^3}$$

I was able to compute this coefficient for a Markov chain of size 2. It might be applied to the general case n by solving a linear system of the same size and computing an approximation of the largest eigenvalue λ .

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

First term is

$$\pi_0 p_{00} \log^2(p_{00}) + \pi_1 p_{10} \log^2(p_{10}) + \pi_0 p_{01} \log^2(p_{01}) + \pi_1 p_{11} \log^2(p_{11})$$

The second is

$$-2 [\dot{\pi}_0(-1)p_{00} \log(p_{00}) + \dot{\pi}_1(-1)p_{10} \log(p_{10}) + \dot{\pi}_0(-1)p_{01} \log(p_{01}) + \dot{\pi}_1(-1)p_{11} \log(p_{11})]$$

The third is

$$-2\dot{\lambda}(-1) [\dot{\pi}_0(-1) + \dot{\pi}_1(-1)]$$

We have to compute $\dot{\pi}(-1)$. With $\pi(s) = (\pi_0(s), \pi_1(s))$, and since

$$\pi(s)P(s) = \lambda(s)\pi(s)$$

then we have

$$\begin{cases} p_{00}^{-s}\pi_0(s) + p_{01}^{-s}\pi_1(s) = \lambda(s)\pi_0(s) \\ p_{10}^{-s}\pi_0(s) + p_{11}^{-s}\pi_1(s) = \lambda(s)\pi_1(s) \end{cases}$$

which I'm not sure how to solve formally.

Références

- [1] JACQUET, SZPANKOWSKI, TANG, *Average profile of the Lempel-Ziv parsing scheme for a Markovian source*