The data compression problem
Process evaluation
Analytic information theory
Application to covariance analysis

# Asymptotics on the Lempel-Ziv 78 compression of Markov sources

Exploring analytic information theory : from Markov source sampling to combinatorial analysis proofs

## Guillaume Duboc

Computer Science Department
Ecole Normale Supérieure de Lyon

M1 Internship, 2018

ENS DE LYON

Asymptotics on the Lempel-Ziv 78 compression of Markov source

The data compression problem
Process evaluation
Analytic information theory
Application to covariance analysis

# Table of contents

ENS DE LYON

The data compression problem
Process evaluation
Analytic information theory
Application to covariance analysis

Introduction to information sources
The LZ78 compression scheme
The compression ratio, and entropy
Algorithmic improvements

# Table des matières

The data compression problem | Introduction to information sources
Process evaluation | The LZ78 compression scheme
Analytic information theory | The compression ratio, and entropy
Application to covariance analysis | Algorithmic improvements

## Words or sequences, and memoryless sources

### Definition : word or sequence or string

Given an alphabet $\mathcal{A}$, a **word** or **sequence** or **string** is an infinite sequence of random variables $X = (X_k)_{k\mathbb{N}^*}$, each $X_k$ representing a symbol in $\mathcal{A}$.

### Definition : Bernoulli or Memoryless source

A source of information is a **Bernoulli** or **memoryless source** when all the symbols of $\mathcal{A}$ occur independently with a fixed probability. The word can be seen as an *infinite sequence of Bernoulli trials*.

ENS DE LYON

The data compression problem
Process evaluation
Analytic information theory
Application to covariance analysis

Introduction to information sources
The LZ78 compression scheme
The compression ratio, and entropy
Algorithmic improvements

## Markov sources definition

### Definition : Markov source

An information source is a **Markov source** when there is a **Markov dependency** between the consecutive symbols of a string.

### Definition : order of a Markov source

Let $V = |\mathcal{A}|$. A **Markov source** is of **order** $r$ when the dependency can be encoded in a transition matrix of size $V^r \times V$, with coefficients :

$$P(c|w) \qquad \forall\, (w, c) \in \mathcal{A}^r \times \mathcal{A}$$

Informally : *the probability that a symbols occurs depends on the previous $r$ symbols.*

The data compression problem
Process evaluation
Analytic information theory
Application to covariance analysis

Introduction to information sources
The LZ78 compression scheme
The compression ratio, and entropy
Algorithmic improvements

# Table des matières

| The data compression problem | Introduction to information sources |
| Process evaluation | The LZ78 compression scheme |
| Analytic information theory | The compression ratio, and entropy |
| Application to covariance analysis | Algorithmic improvements |

# Description of the LZ78 algorithm

### Algorithm

Given a word *w*.

- Initialize an empty dictionary
- While it is possible :

> Find longest prefix of *w* that is not in the dictionary
>
> Add it to the dictionary, cut it from *w*

The data compression problem
Process evaluation
Analytic information theory
Application to covariance analysis

Introduction to information sources
The LZ78 compression scheme
The compression ratio, and entropy
Algorithmic improvements

### Elements description

The data representation is (dictionary_reference, symbol).

### Remarks

The LZ78 algorithm builds a prefix tree from which the original word can be reconstructed.

| The data compression problem | Introduction to information sources |
| Process evaluation | The LZ78 compression scheme |
| Analytic information theory | The compression ratio, and entropy |
| Application to covariance analysis | Algorithmic improvements |

### Definition : number of phrases

After compressing a word $w$, the number of phrases in the dictionary is noted $M(w)$.
For words of size $n$, we write $M_n(w)$.

### Code length

$$C(w) = \sum_{k=0}^{M(w)} (\lceil \log_2(k) \rceil + \lceil \log_2(\mathcal{A}) \rceil)$$

ENS DE LYON

Asymptotics on the Lempel ZN 78 compression of Markov source

The data compression problem · Process evaluation · Analytic information theory · Application to covariance analysis

Introduction to information sources
The LZ78 compression scheme
The compression ratio, and entropy
Algorithmic improvements

# Table des matières

| The data compression problem | Introduction to information sources |
| Process evaluation | The LZ78 compression scheme |
| Analytic information theory | The compression ratio, and entropy |
| Application to covariance analysis | Algorithmic improvements |

### Definition : compression ratio

Let *w* a word, and *C*(*w*) its *encoding* by a compression algorithm. The **compression ratio** of *w* is $\dfrac{|C(w)|}{|w|}$.

### Main goals of compression algorithms

- Improving the compression ratio
- Fast compression/decompression speed in Mb/s

### T

he tradeoff between these two goals is a sensitive research problem. Different compression standards :

- Google (Brotli, 2015)
- Facebook (Zstandard, 2016)

The data compression problem
Process evaluation
Analytic information theory
Application to covariance analysis

Introduction to information sources
The LZ78 compression scheme
The compression ratio, and entropy
Algorithmic improvements

## Optimal encoding

### Entropy of a Markov source

Let $\pi$ be a stationary distribution. The entropy of a Markov chain is

$$h = -\sum_{i=1}^{V} \pi \sum_{j=1}^{V} p_{ij} \log(p_{ij})$$

### Optimality of LZ78

Considering words of length $n$.

$$\frac{|C(w)|}{|w|} - h \text{ goes to zero for } n \to +\infty$$

The data compression problem
Process evaluation
Analytic information theory
Application to covariance analysis

Introduction to information sources
The LZ78 compression scheme
The compression ratio, and entropy
Algorithmic improvements

# Table des matières

The data compression problem
Process evaluation
Analytic information theory
Application to covariance analysis

Introduction to information sources
The LZ78 compression scheme
The compression ratio, and entropy
Algorithmic improvements

## Optimal parsing

The data compression problem | Introduction to information sources
Process evaluation | The LZ78 compression scheme
Analytic information theory | The compression ratio, and entropy
Application to covariance analysis | Algorithmic improvements

## Flexible parsing

The data compression problem
Process evaluation
Analytic information theory
Application to covariance analysis

Theoretical models
Experimental conditions
Extracting results

# Table des matières

The data compression problem
Process evaluation
Analytic information theory
Application to covariance analysis

Theoretical models
Experimental conditions
Extracting results

# Markov Independent Model

$$X(1) = 0000000\ldots$$
$$X(2) = 1010101\ldots$$
$$X(3) = 1001101\ldots$$
$$X(4) = 001100111\ldots$$

The data compression problem
**Process evaluation**
Analytic information theory
Application to covariance analysis

Theoretical models
**Experimental conditions**
Extracting results

# Table des matières

The data compression problem
Process evaluation
Analytic information theory
Application to covariance analysis

Theoretical models
Experimental conditions
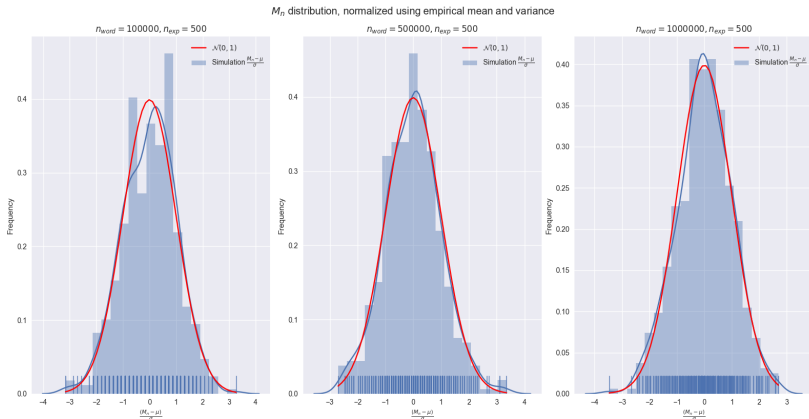Extracting results

## Coding details

- Python code $\sim$ 2000 lines

- Markov source sampling

- Optimized datastructure (digital search tree)

- Parallelization

- Reproducibility of datasets

The data compression problem
Process evaluation
Analytic information theory
Application to covariance analysis

Theoretical models
Experimental conditions
Extracting results

# Table des matières

The data compression problem
Process evaluation
Analytic information theory
Application to covariance analysis

Theoretical models
Experimental conditions
Extracting results

# Central Limit Theorem confirmation

The data compression problem
Process evaluation
Analytic information theory
Application to covariance analysis

Theoretical models
Experimental conditions
Extracting results
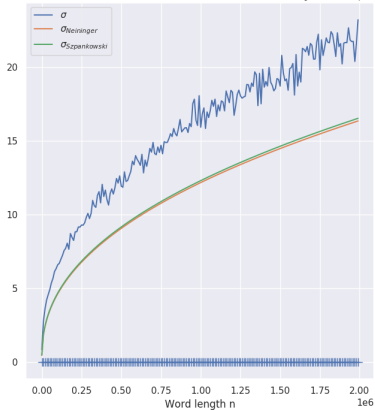
# Hypothesis testing for the variance

### Complex matrix

Defining $P(s)$ as $\begin{matrix} p_{11}^{-s} & p_{12}^{-s} \\ p_{21}^{-s} & p_{22}^{-s} \end{matrix}$

### Variance expression

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

Asymptotics on the Lempel-Ziv 78 compression of Markov source

The data compression problem
Process evaluation
Analytic information theory
Application to covariance analysis

Theoretical models
Experimental conditions
Extracting results

Empirical standard deviation ($\sigma$) and theoretical ones ($\sigma_{Neininger}$, $\sigma_S$), $n_{exp} = 400$

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

The data compression problem
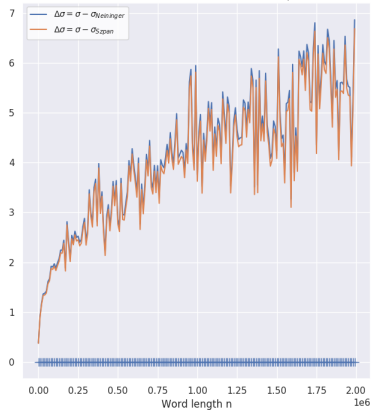Process evaluation
Analytic information theory
Application to covariance analysis

Power series
Complex analysis tools

# Table des matières

ENS DE LYON

The data compression problem
Process evaluation
Analytic information theory
Application to covariance analysis

Power series
Complex analysis tools

## Definition, usage

### Definition

$$A(z) = \sum_{n \geqslant 0} a_n z^n$$

### Remarks

- Used as an algebraic item with the convolution product
- No convergence problems

The data compression problem
Process evaluation
**Analytic information theory**
Application to covariance analysis

Power series
Complex analysis tools

# Table des matières

ENS DE LYON

The data compression problem
Process evaluation
Analytic information theory
Application to covariance analysis

Power series
Complex analysis tools

### Poissonization and Depoissonization

$$\widetilde{G}(z) = \sum_{n \geqslant 0} a_n \frac{z^n}{n!} e^{-z}$$

### Mellin transform

Make recurrence relation between random variables become linear in order to solve them more easily.

The data compression problem
Process evaluation
Analytic information theory
Application to covariance analysis

Tail symbols
Simulation results
Analytic solution

# Table des matières

The data compression problem
Process evaluation
Analytic information theory
**Application to covariance analysis**

Tail symbols
Simulation results
Analytic solution

# Tail symbols

### Illustration

$$X(1) = 0000000\ldots$$
$$X(2) = 1010101\ldots$$
$$X(3) = 1001101\ldots$$
$$X(4) = 001100111\ldots$$

### Definition

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*.

The data compression problem
Process evaluation
Analytic information theory
**Application to covariance analysis**

Tail symbols
Simulation results
Analytic solution

# Definition and relation

## Recurrence

For $n \geqslant 0$, we have :

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

## Notations

- $\delta_a = \begin{cases} 1 & \text{if } a \text{ is the tail symbol of the first sequence} \\ 0 & \text{else} \end{cases}$

- $N_a$ is the random variable giving *the size of the left subtree which contains phrases whose second letter is a*

- $\widetilde{T}_{N_a}^a$ is the number of times $a$ is a tail symbol for the sequences that were used to build the subtree with

The data compression problem
Process evaluation
Analytic information theory
**Application to covariance analysis**

Tail symbols
Simulation results
Analytic solution

# Total path lenght

### 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.

### Recurrence relation

For all $n \geqslant 0$ :

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

The data compression problem
Process evaluation
Analytic information theory
Application to covariance analysis

Tail symbols
Simulation results
Analytic solution

# Table des matières

The data compression problem
Process evaluation
Analytic information theory
Application to covariance analysis

Tail symbols
Simulation results
Analytic solution

Inconclusive, but informative

The data compression problem
Process evaluation
Analytic information theory
Application to covariance analysis

Tail symbols
Simulation results
Analytic solution

# Table des matières

The data compression problem
Process evaluation
Analytic information theory
**Application to covariance analysis**

Tail symbols
Simulation results
Analytic solution

## Recurrence

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

The data compression problem
Process evaluation
Analytic information theory
Application to covariance analysis

Tail symbols
Simulation results
Analytic solution

### Poisson transform

Defining

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

### Differential equation

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

ENS DE LYON

Asymptotics on the Lempel-Ziv 78 compression of Markov source

The data compression problem
Process evaluation
Analytic information theory
Application to covariance analysis

Tail symbols
Simulation results
Analytic solution

$M_n$ distribution, normalized with empirical mean and Szpankowski variance.

The data compression problem
Process evaluation
Analytic information theory
**Application to covariance analysis**

Tail symbols
Simulation results
**Analytic solution**

$M_n$ distribution, normalized with empirical mean and Neininger variance.

The data compression problem
Process evaluation
Analytic information theory
Application to covariance analysis

Tail symbols
Simulation results
Analytic solution

The data compression problem
Process evaluation
Analytic information theory
Application to covariance analysis

Tail symbols
Simulation results
Analytic solution

The data compression problem
Process evaluation
Analytic information theory
Application to covariance analysis

Tail symbols
Simulation results
Analytic solution