

INFORMATION THEORY:
SOURCES, DIRICHLET SERIES,
REALISTIC ANALYSIS OF ALGORITHMS

Brigitte VALLÉE
GREYC Laboratory
(CNRS and University of Caen, France)

Talk based on joint works with
Eda CESARATTO, Julien CLÉMENT, Jim FILL,
Philippe FLAJOLET, and Mathieu ROUX

Séminaire de Combinatoire énumérative et analytique,
IHP, Paris, 3 février 2011

Description of a framework which

- unifies the analyses for text algorithms and searching/sorting algorithms
- provides a general model for sources
- shows the importance of the Dirichlet generating functions
- explains the importance of tameness for sources
- defines a natural subclass of sources, the dynamical sources
- provides sufficient conditions for tameness of dynamical sources
- provides probabilistic analyses for algorithms built on tame sources.

Plan of the talk.

- General motivations
- Models for a source
- The Dirichlet generating function of the source
- Conclusion and possible extensions.

Plan of the talk.

- General motivations
- Models for a source
- The Dirichlet generating function of the source
- Conclusion and possible extensions.

The classical framework for analysis of algorithms
in two main algorithmic domains:

Text algorithms – Sorting or Searching algorithms.

– In text algorithms, algorithms deal with words

What is a word ?

.... a sequence of symbols from the same alphabet

– In sorting or searching algorithms, algorithms deal with keys.

The set of keys must be ordered.

What is a key?

.... a sequence of symbols from the same alphabet

Key or word? the same object... but,

– for comparing two words, importance of the structure of words

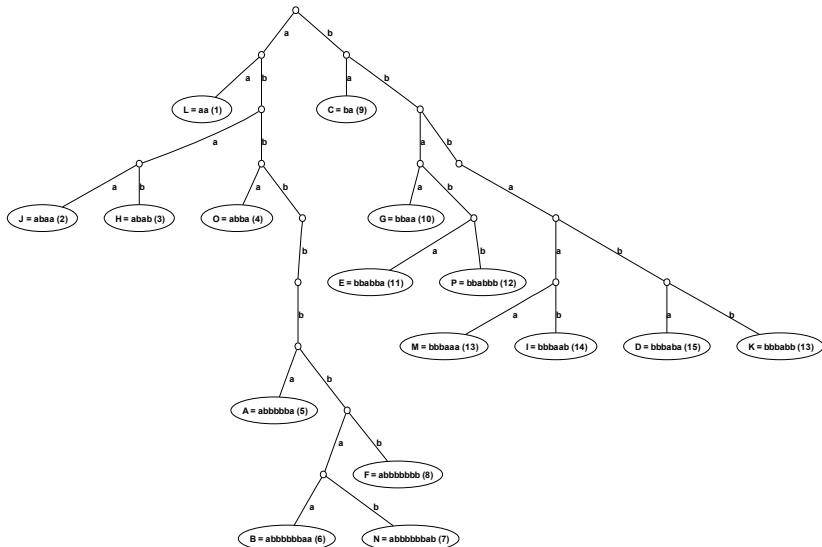
– for comparing two keys, transpance of the structure of keys

only their relative order plays a role.

Text algorithms and dictionaries : The trie structure

An example : A trie built on a set of 16 words.

A = **abbbba**aabab B = **abbbba**abaa C = **ba**abbbabba D = **bbbaba**bbbaab E = **bbabba**ababb
F = **abbbbbb**abb G = **bbaa**abbabab H = **abab**bbabbab I = **bbbaab**bbbbbb J = **abaa**bbbaabb
K = **bbbabb**bbbaa L = **aa**aabbabaaba M = **bbbaa**abbbbbb N = **abbbbbb**abaa O = **abba**bababbbb P = **bbabbb**aaaabb



Probabilistic study of the Trie structure.

Main parameter on a node n_w labelled with prefix w :

N_w := the number of words which **begin** with prefix w .

N_w := the number of words which **go through** the node n_w

The size, and the path length of a trie equal

$$R = \sum_{w \in \Sigma^*} \mathbf{1}_{[N_w \geq 2]} \qquad T = \sum_{w \in \Sigma^*} \mathbf{1}_{[N_w \geq 2]} \cdot N_w,$$

Role of p_w := the probability that a word **begins** with prefix w .

Classical analyses of the main algorithms for searching or sorting

The unit cost is the **key-comparison**.

The behaviour of the algorithm (wrt to **key-comparisons**) only depends on the **relative order** between the keys.

It is sufficient to restrict to the case when $\Omega = [1..n]$.

The input set is then \mathfrak{S}_n , with uniform probability.

Then, the analysis of all these algorithms is very well known, with respect to the **number of key-comparisons performed** in the worst-case, or in **the average case**.

A more realistic framework for sorting or searching.

Keys are viewed as words. The domain Ω of keys is a subset of $\Sigma^{\mathbb{N}}$,

$$\Sigma^{\mathbb{N}} = \{\text{the infinite words on some ordered alphabet } \Sigma\}.$$

The words are compared [wrt the lexicographic order].

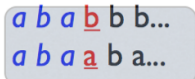
The realistic unit cost is now the symbol-comparison.

The realistic cost of the comparison between two words A and B ,

$$A = a_1 a_2 a_3 \dots a_i \dots \quad \text{and} \quad B = b_1 b_2 b_3 \dots b_i \dots$$

equals $k + 1$, where k is the length of their largest common prefix

$$k := \max\{i; \forall j \leq i, a_j = b_j\} = \text{the coincidence}$$



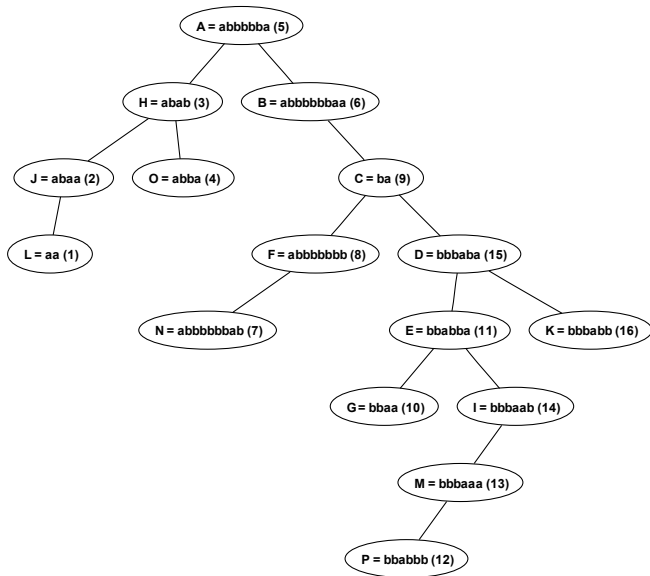
a b a b b b...
a b a a b a...

coincidence=3; #comparisons=4.

Now, sorting or searching algorithms are viewed as text algorithms.

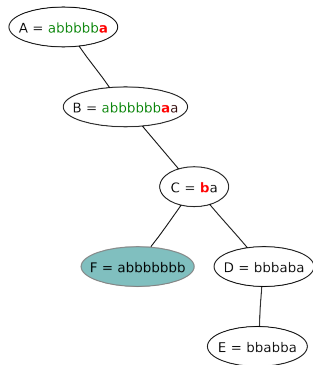
An example : The BST (binary search tree) built on a sequence of words

A = abbbbaaabab B = abbbbaabaa C = baabbabbbba D = bbbababbaab E = bbabbaababb
F = abbbbbbabab G = bbaabbababa H = ababbabbbab I = bbbaabbbbbbb J = abaaabbbbaab
K = bbbabbbbaa L = aaabbabaaba M = bbbbaabbbbbb N = abbbbbbabba O = abbaabababbb P = bbabbbbaaab



An example : The cost of the insertion of the key F into the BST

$F = \text{abbbbbb}$

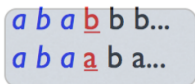


Number of symbol comparisons
needed ?

= 7 for comparing to A
+ 8 for comparing to B
+ 1 for comparing to C

Total = 16

The **realistic cost** of the comparison between two words A and B is related to their **coincidence** $c(A, B)$



The diagram shows two words, $A = a b a b b \dots$ and $B = a b a a b a \dots$, enclosed in a rounded rectangle. The first three characters of both words are aligned. In the first row, the characters 'a', 'b', and 'a' are blue, while the 'b' is red and underlined. In the second row, the characters 'a', 'b', and 'a' are blue, while the 'a' is red and underlined. Below the rectangle, there are two horizontal lines.

coincidence=3; #comparisons=4.

The coincidence $c(A, B)$ satisfies $c(A, B) \geq k$ if and only if A and B **begin** with the **same prefix** of length k

Importance of $p_w :=$ the probability that a word **begins** with prefix w .

Now, we work inside an **unifying** framework
where **searching and sorting** algorithms are viewed as **text** algorithms.

In this context, the **probabilistic behaviour** of algorithms heavily depends
on the **mechanism** which produces **words**.

A **source** := a mechanism which produces symbols from alphabet Σ ,
one for each time unit.

When (discrete) time evolves, a source produces (infinite) words of $\Sigma^{\mathbb{N}}$.

For $w \in \Sigma^*$, $p_w :=$ probability that a word **begins** with the prefix w .

The set $\{p_w, w \in \Sigma^*\}$ defines the source \mathcal{S} .

Fundamental role of the **Dirichlet generating functions** of the source

$$\Lambda(s) := \sum_{w \in \Sigma^*} p_w^s, \quad \Lambda_k(s) = \sum_{w \in \Sigma^k} p_w^s$$

Remark: $\Lambda_k(1) = 1$ for any k , $\Lambda(1) = \infty$.

- they encapsulate the main probabilistic properties of the source
- they translate them into analytic properties

For instance, the **entropy** $h_{\mathcal{S}}$, the **coincidence** $c_{\mathcal{S}}$

$$h(\mathcal{S}) := \lim_{k \rightarrow \infty} \frac{-1}{k} \sum_{w \in \Sigma^k} p_w \log p_w = - \lim_{k \rightarrow \infty} \Lambda'_k(1)$$

$$\Pr[c_{\mathcal{S}} \geq k + 1] = \sum_{w \in \Sigma^k} p_w^2 = \Lambda_k(2)$$

- they intervene in the probabilistic analyses of text algorithms and (also) sorting and searching algorithms.

Three main steps for the analysis
of the mean number S_n of symbol comparisons

(1) **First step** (algebraic).

The **Poisson model** \mathcal{P}_Z deals with a variable number N of words:
 N is a **random variable** which follows a Poisson law of **parameter** Z .

We first obtain **nice** expressions for the mean number $\tilde{S}(Z)$

(2) **Second step** (algebraic).

It is possible to return to the model where the **number** n of words is **fixed**.
We obtain a nice **exact** formula for S_n

from which it is **not easy** to obtain the asymptotics...

(3) **Third step** (analytic).

Then, the **Rice formula** provides the **asymptotics of** S_n ($n \rightarrow \infty$),
as soon as the **source** is “tame”

After the second step, an **exact formula** for S_n

$$S_n = \sum_{k=2}^n (-1)^k \binom{n}{k} \varpi(k)$$

...which involves the series ϖ at integer values k .

For the mean path length (Trie or BST),

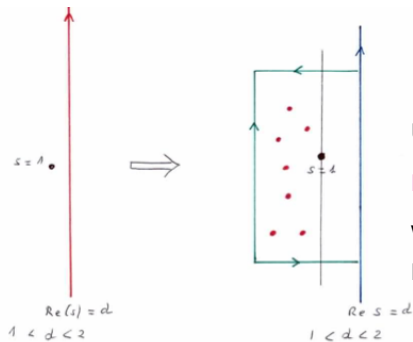
$\varpi(s)$ is closely related to the Dirichlet series of the probabilities,

$$\varpi_T(s) = s\Lambda(s) \quad \varpi_B(s) = 2 \frac{\Lambda(s)}{s(s-1)} \quad \text{where} \quad \Lambda(s) := \sum_{w \in \Sigma^*} p_w^s$$

Asymptotic analysis.

The residue formula transforms the sum into an integral with $1 < d < 2$.

$$S_n = \sum_{k=2}^n (-1)^k \binom{n}{k} \varpi(k) = \frac{1}{2i\pi} \int_{d-i\infty}^{d+i\infty} \varpi(s) \frac{n! (-1)^{n+1}}{s(s-1)\dots(s-n)} ds,$$



We **shift** the integral on the **left**,

Usually, the first singularities occur at $\Re s = 1$.

Behaviour of $\varpi(s)$ [or $\Lambda(s)$] near $\Re s = 1$?

Where are the **red singularities** closest to $\Re s = 1$?

Is $\Lambda(s)$ of polynomial growth on the **green contour**?

Importance of the existence of a **region \mathcal{R}**

– which contains only $s = 1$ as a **pole** – where $\Lambda(s)$ is of **polynomial growth**.

Tameness of the source

Case of $\text{Trie}(n)$ [Clément, Flajolet, V. 2001]

Theorem 1. For any tame source, the mean path length T_n of a trie built on n words independently drawn from the source satisfies

$$T_n \sim \frac{1}{h_{\mathcal{S}}} n \log n.$$

and involves the **entropy** $h_{\mathcal{S}}$ of the source \mathcal{S} , defined as

$$h_{\mathcal{S}} := \lim_{k \rightarrow \infty} \left[\frac{-1}{k} \sum_{w \in \Sigma^k} p_w \log p_w \right],$$

where p_w is the probability that a word **begins** with prefix w .

Case of QuickSort(n) or BST(n) [Clément, Fill, Flajolet, V. 2009]

Theorem 2. For any tame source, the mean number S_n of symbol comparisons used by QuickSort(n) (or the mean number of symbols comparisons used to built the BST) on n words of the source satisfies

$$B_n \sim \frac{1}{h_S} n \log^2 n.$$

and involves the **entropy** h_S of the source \mathcal{S} , defined as

$$h_S := \lim_{k \rightarrow \infty} \left[\frac{-1}{k} \sum_{w \in \Sigma^k} p_w \log p_w \right],$$

where p_w is the probability that a word **begins** with prefix w .

Compared to $K_n \sim 2n \log n$, there is an extra factor equal to $1/(2h_S) \log n$

Compared to $T_n \sim (1/h_S) n \log n$, there is an extra factor of $\log n$.

Plan of the talk.

- General motivations
- Models for a source
- The Dirichlet generating function of the source
- Conclusion and possible extensions.

The parametrization of a general source

A general **source** \mathcal{S} produces infinite words

on an **ordered alphabet** $\Sigma := \{a_1, \dots, a_r\}$.

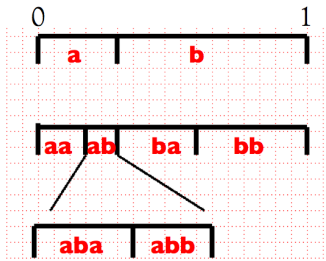
The set of infinite words produced by \mathcal{S} is $\mathcal{L}(\mathcal{S}) \subset \Sigma^{\mathbb{N}}$.

For $w \in \Sigma^*$, $p_w :=$ probability that a word **begins** with the prefix w .

The set $\{p_w, w \in \Sigma^*\}$ defines the source \mathcal{S} .

We assume: $\pi_k := \sup\{p_w, w \in \Sigma^k\} \rightarrow 0$ for $k \rightarrow \infty$

For **each length** k , we consider the **probabilities** p_w with $w \in \Sigma^k$
sorted with respect to the **lexicographic order** on Σ^k .



For each $w \in \Sigma^k$ $p_w^{(-)} := \sum_{\substack{\alpha \in \Sigma^k, \\ \alpha < w}} p_\alpha$

For any $X \in \mathcal{L}(\mathcal{S})$, let

If $X := \lim_{k \rightarrow \infty} w_k$, $F(X) := \lim_{k \rightarrow \infty} p_w^{(-)}$

Then $F(X) := \Pr[Y < X]$

F is the distribution function on $\mathcal{L}(\mathcal{S})$.

The function $F : \mathcal{L}(\mathcal{S}) \rightarrow [0, 1]$ is **continuous**
and **strictly increasing** outside an exceptional (denumerable set)

Outside this exceptional set, each infinite word X of $\mathcal{L}(\mathcal{S})$ is written as

$$X = M(u) \text{ with } M : [0, 1] \rightarrow \mathcal{L}(\mathcal{S}).$$

The real u is the **parameter** of the word X .

The map M provides a **parametrization** of the source \mathcal{S} .

Via the mapping M ,

[Drawing of words X in \mathcal{S}] \equiv [Uniform drawing of parameters u in $[0, 1]$]

For any finite prefix $w \in \Sigma^*$,

the set $\{u, M(u) \text{ begins with } w\}$ is an interval

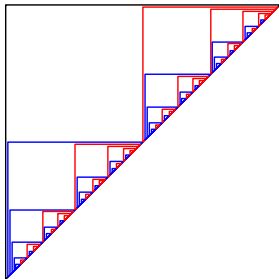
This is the **fundamental interval** of w . Its length equals p_w .

For any finite prefix $w \in \Sigma^*$,

the set $\{u, M(u) \text{ begins with } w\}$ is an interval

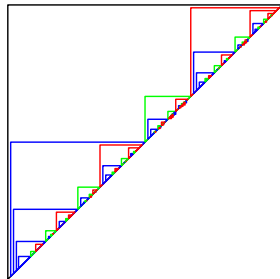
This is the **fundamental interval** of w . Its length equals p_w .

Instances of fundamental intervals for two **memoryless** sources.



Memoryless source on $\{a, b\}$

$$p_a = 1/2, p_b = 1/2$$



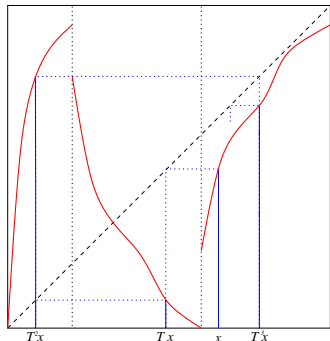
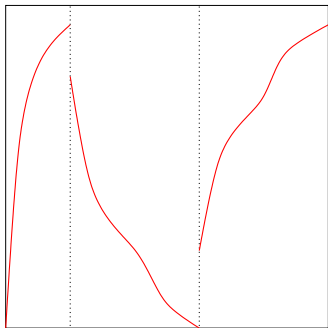
Memoryless source on $\{a, b, c\}$

$$p_a = 1/2, p_b = 1/6, p_c = 1/3$$

In the memoryless case, there are **regular** splittings.

A general class of “natural” sources: dynamical sources
associated to a “natural” parametrization

With a shift map $T : \mathcal{I} \rightarrow \mathcal{I}$ and an encoding map $\sigma : \mathcal{I} \rightarrow \Sigma$,
the emitted word is $M(u) = (\sigma u, \sigma T u, \sigma T^2 u, \dots, \sigma T^k u, \dots)$
namely, the encoded trajectory of u



A dynamical system, with $\Sigma = \{a, b, c\}$ and a word $M(u) = (c, b, a, c \dots)$.

A **dynamical source** = a source built with a dynamical system

A **dynamical system** (\mathcal{I}, S) is defined by four elements:

- a finite **alphabet** Σ ,
- a topological **partition** of $\mathcal{I} :=]0, 1[$ with open intervals $\mathcal{I}_{m, m \in \Sigma}$,
- an **encoding mapping** σ equal to m on each \mathcal{I}_m ,
- a **shift mapping** T

s.t. $T|_{\mathcal{I}_m}$ is a bijection of class \mathcal{C}^2 from \mathcal{I}_m to $\mathcal{J}_m := T(\mathcal{I}_m)$.

This gives rise to a source: on an input u of \mathcal{I} , it outputs the word

$$M(u) := (\sigma u, \sigma T u, \sigma T^2 u, \dots).$$

When an **initial density** –and an initial distribution F – is chosen on \mathcal{I} ,
this induces (via M) a **probabilistic model** on Σ^∞
= a dynamical source \mathcal{S}_F .

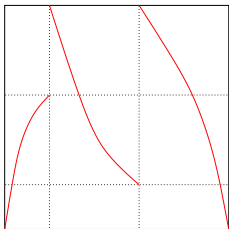
Strong relations between the geometry of the system
and the probabilistic properties of the source.

Correlations between symbols are mainly due to two geometric characteristics:
the position of the branches and the shape of the branches.

– the **position** of the branches: the position of $T(\mathcal{I}_m)$ wrt \mathcal{I}_ℓ ;
it describes the set $s(m)$ of possible successors of the symbol m .

Particular cases: – Complete systems $T(\mathcal{I}_m) = \mathcal{I}$

– Markovian systems $T(\mathcal{I}_m) = \text{union of some } \mathcal{I}_\ell$
give rise to a **finite** characterization of $s(m)$.



A markovian system

Generally speaking,
importance of **topological mixing**:
“There is a word of length n
which begins with b and ends with e ”.

Strong relations between the geometry of the system
and the probabilistic properties of the source.

Correlations between symbols are also due to the shape of the branches.

– the shape of the branches, is described by their derivatives;
it explains how the distribution evolves.

Less correlated systems correspond to systems with affine branches.

Generally speaking, importance of expansiveness:

the derivative T' satisfies $\forall u \in \mathcal{I} \quad |T'(u)| \geq \delta > 1$.

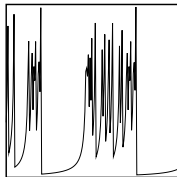
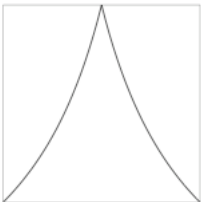
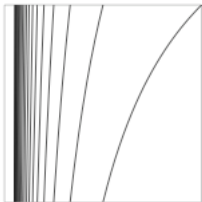
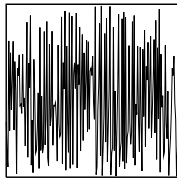
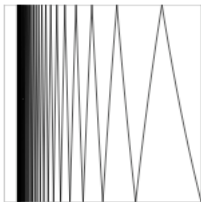
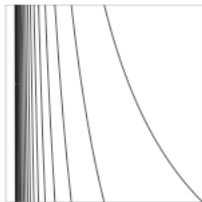
When true, this implies a chaotic behaviour for trajectories.

When this condition is violated at only one indifferent point,

$$[T(u) = u, |T'(u)| = 1]$$

this leads to intermittency phenomena.

Four Euclidean dynamical sources, Two different behaviours



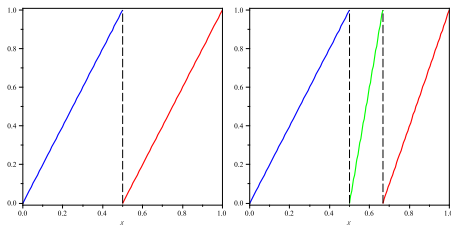
Particular cases: simple sources and affine branches

A **memoryless** source

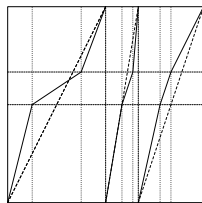
:= a complete system with affine branches and uniform initial density

A **Markov chain**

:= a Markovian system with affine branches,
with an initial density which is constant on each \mathcal{I}_m .



Two memoryless sources

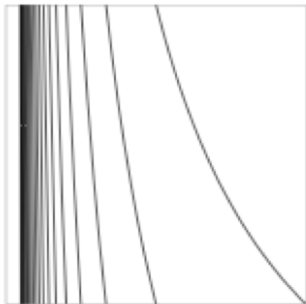


a Markov chain.

General case of interest.

- A **complete** –or a **Markovian**– system
- with a possible **infinite** denumerable alphabet
 - topologically **mixing** – and **expansive**.

Main instance: the **Euclidean source** defined with $T(x) := \frac{1}{x} - \lfloor \frac{1}{x} \rfloor$



Plan of the talk.

- General motivations
- Models for a source
- The Dirichlet generating function of the source
- Conclusion and possible extensions.

A main analytical object related to any source:

the Dirichlet series of probabilities, $\Lambda(s) := \sum_{w \in \Sigma^*} p_w^s$

Memoryless sources, with probabilities (p_i)

$$\Lambda(s) = \frac{1}{1 - \lambda(s)} \quad \text{with} \quad \lambda(s) = \sum_{i=1}^r p_i^s$$

Markov chains, defined by – the vector R of initial probabilities (r_i)
– and the transition matrix $P := (p_{i,j})$

$$\Lambda(s) = {}^t \mathbf{1} (I - P(s))^{-1} R(s) \quad \text{with} \quad P(s) = (p_{i,j}^s), \quad R(s) = (r_i^s).$$

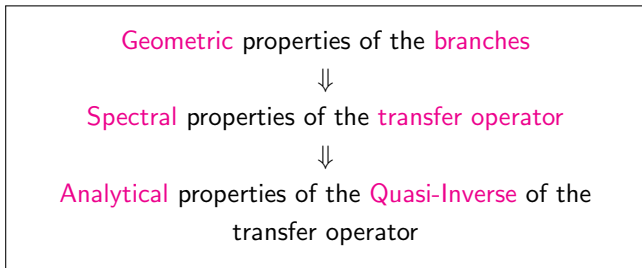
A general dynamical source

$$\Lambda(s) \text{ closely related to } (I - \mathbb{H}_s)^{-1}$$

where \mathbb{H}_s is the (secant) transfer operator of the dynamical system.

Probabilistic analysis of text algorithms
when the text is generated by a dynamical source.

A dynamical source



↓

Analytical properties of the Dirichlet generating function $\Lambda(s)$

↓

Probabilistic analysis of text algorithms

Fundamental probabilities for a dynamical source (complete case).

The fundamental probability p_w = the probability that $M(u)$ begins with w

For any $w \in \Sigma^*$, $\mathcal{I}_w := \{u; M(u) \text{ begins with the prefix } w\}$

Here, in the complete dynamical case, for any k and any $w \in \Sigma^k$:

- the restriction $T^k|_{\mathcal{I}_w}$ is a \mathcal{C}^2 bijection of \mathcal{I}_w onto \mathcal{I}
- Its inverse mapping h_w is a \mathcal{C}^2 bijection from \mathcal{I} to $\mathcal{I}_w = [h_w(0), h_w(1)]$.

With a change of variables, this provides an alternative expression for p_w :

$$p_w = \int_{\mathcal{I}_w} f(t) dt = \int_{\mathcal{I}} |h'_w(x)| \cdot f \circ h_w(x) \cdot dx.$$

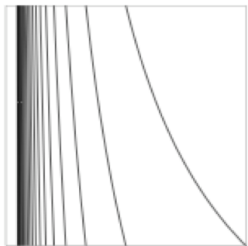
which involves the operator $\mathbf{H}_{[w]}$

$$\mathbf{H}_{[w]}[f](x) := |h'_w(x)| \cdot f \circ h_w(x)$$

via the relation

$$p_w = \int_{\mathcal{I}} \mathbf{H}_{[w]}[f](t) dt$$

The density transformer and the transfer operator



The operator $\mathbf{H} := \sum_{a \in \Sigma} \mathbf{H}_{[a]}$

with $\mathbf{H}_{[a]}[f](x) = |h'_a(x)| \cdot f \circ h_a(x)$

is the density transformer of the dynamical system.

It describes the evolution of the density after one iteration.

For a density f on $[0, 1]$,

$\mathbf{H}[f]$ is the density on $[0, 1]$ after one iteration.

Transfer operator (Ruelle)

$\mathbf{H}_s := \sum_{a \in \Sigma} \mathbf{H}_{s,[a]}$ with $\mathbf{H}_{s,[a]}[f](x) = |h'_a(x)|^s f \circ h_a(x)$.

The k -th iterate satisfies:

$\mathbf{H}_s^k = \sum_{w \in \Sigma^k} \mathbf{H}_{s,[w]}$ with $\mathbf{H}_{s,[w]}[f](x) = |h'_w(x)|^s f \circ h_w(x)$

Generation of p_w^s .

For an inverse branch h of any depth, and an initial distribution F , consider the “secants” H, L of h, F ,

$$H(x, y) := \left| \frac{h(x) - h(y)}{x - y} \right|, \quad L(x, y) = \left| \frac{F(x) - F(y)}{x - y} \right|.$$

and the component “secant” operator $\mathbb{H}_{s,[w]}$, defined as

$$\mathbb{H}_{s,[w]}[L](x, y) := H_w^s(x, y) \cdot L(h_w(x), h_w(y)).$$

This operator generates p_w^s :

$$\begin{aligned} p_w^s &= |F(h_w(1)) - F(h_w(0))|^s = \left| \frac{h_w(1) - h_w(0)}{1 - 0} \right|^s \cdot \left| \frac{F(h_w(1)) - F(h_w(0))}{h_w(1) - h_w(0)} \right|^s \\ &= \mathbb{H}_{s,[w]}[L^s](1, 0) \end{aligned}$$

“On the diagonal” $\mathbb{H}_{s,[w]}[L](x, x) = \mathbf{H}_{s,[w]}[f](x)$

The diagonal of the secant is the tangent

The secant operator is an extension of the plain (tangent) operator

For $w, w' \in \Sigma^*$, the multiplicative relation $\mathbb{H}_{[s, w \cdot w']} = \mathbb{H}_{s, [w']} \circ \mathbb{H}_{s, [w]}$ generalizes the equality $p_{w \cdot w'}^s = p_w^s \cdot p_{w'}^s$,
no longer true when the source has memory.

The Dirichlet series of fundamental probabilities

$$\Lambda_k(s) := \sum_{w \in \Sigma^k} p_w^s, \quad \Lambda(s) := \sum_{w \in \Sigma^*} p_w^s$$

are “generated” by the secant transfer operator \mathbb{H}_s [V. 2000]

$$\Lambda_k(s) = \mathbb{H}_s^k[L^s](0, 1), \quad \Lambda(s) = (I - \mathbb{H}_s)^{-1}[L^s](0, 1).$$

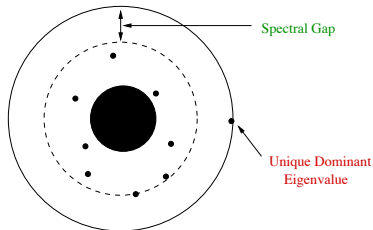
Singularities of $s \mapsto \Lambda(s)$ are essential in the analysis.

Singularities of $(I - \mathbb{H}_s)^{-1}$ are related to spectral properties of \mathbb{H}_s .

For $s = 1$, \mathbb{H}_1 is an extension of \mathbf{H} and has an eigenvalue equal to 1.

A source is decomposable if there exists a Banach space \mathcal{L} for which $\mathbb{H}_s : \mathcal{L} \rightarrow \mathcal{L}$ possesses for real $s > s_0$ (with $s_0 < 1$)

- a unique **dominant eigenvalue** $\lambda(s)$
- and a **spectral gap**.



Spectrum of \mathbb{H}_s for a decomposable source

Sufficient conditions for decomposability:

A **markovian** and **expansive** system on a finite alphabet is decomposable.

In this case, the function $s \mapsto \Lambda(s)$

- is **analytic** in the plane $\Re(s) > 1$,
- and it has a **simple pole** at $s = 1$.

$$\Lambda(s) \sim_{s \rightarrow 1} \frac{1}{1 - \lambda(s)} \sim_{s \rightarrow 1} \frac{-1}{\lambda'(1)(s - 1)}$$

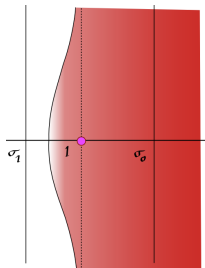
The entropy of the source h_S is related to $\lambda(s)$, namely $h_S = -\lambda'(1)$

And on the left of the vertical line $\Re s = 1$?

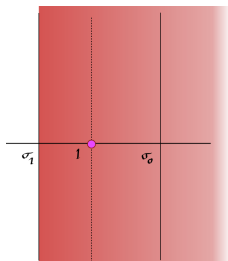
It is important for the analysis to deal with a region \mathcal{R} where $\Lambda(s)$ is **tame**

- it is analytic
- it is of polynomial growth when $\Im s \rightarrow \infty$

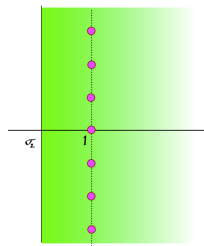
Different possible regions \mathcal{R} on the left of $\Re s = 1$ where $\Lambda(s)$ is tame.



Situation 1
Hyperbolic region

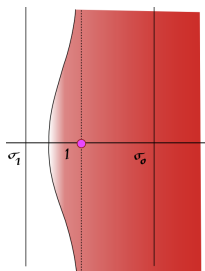


Situation 2
Vertical strip

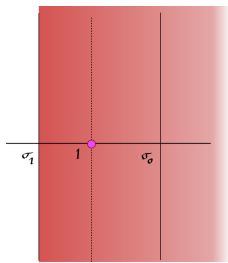


Situation 3
Vertical strip with holes

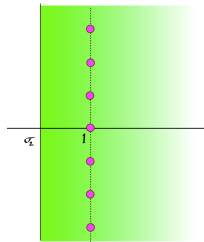
Different possible regions on the left of $\Re s = 1$ where $\Lambda(s)$ is tame.



Situation 1
Hyperbolic region



Situation 2
Vertical strip



Situation 3
Vertical strip with holes

For which simple sources do these different situations occur?

For **memoryless** sources relative to probabilities (p_1, p_2, \dots, p_r)

- S2 is **impossible**
- S3 occurs when **all** the ratios $\log p_i / \log p_j$ are **rational**
- S1 occurs if there **exists** a ratio $\log p_i / \log p_j$
which is **badly approximable by rationals**.

Memoryless sources – The periodic case

In this case

$$\Lambda(s) = \frac{1}{1 - \lambda(s)} \quad \text{with} \quad \lambda(s) = p_1^s + p_2^s = \dots p_r^s.$$

Case $r = 2$. Suppose that there exists $t \neq 0$ for which $\lambda(1 + it) = 1$.

Then $t \neq 0$ is solution of $p_1 p_1^{it} + p_2 p_2^{it} = 1$.

Then, as $p_1 + p_2 = 1$, this implies:

$$1 = p_1 |p_1^{it}| + p_2 |p_2^{it}|, \quad |p_1 p_1^{it} + p_2 p_2^{it}| = 1$$

The “converse of the triangular inequality” entails $p_1^{it} = p_2^{it} = 1$

$$\implies t \log p_1 = 2q_1 \pi, \quad t \log p_2 = 2q_2 \pi$$

$$\implies \frac{\log p_1}{\log p_2} \in \mathbb{Q}$$

$$\implies s \mapsto \lambda(s) \text{ is periodic with period } it$$

Memoryless sources – Case when there are poles close to $\Re s = 1$

Knowing $p_1 + p_2 = 1$,

Look for a solution s of $p_1^s + p_2^s = 1$ when $s = \sigma + it$, with σ close to 1

$$p_1^\sigma p_1^{it} + p_2^\sigma p_2^{it} = 1 \quad \implies \quad p_1^{it} \approx 1, \quad p_2^{it} \approx 1$$

$$\implies \quad \exists q_1, q_2 \in \mathbb{Z} \quad \text{for which} \quad t \approx \frac{2\pi}{\log p_1} q_1, \quad t \approx \frac{2\pi}{\log p_2} q_2$$

Thus

$$\frac{\log p_2}{\log p_1} \approx \frac{q_2}{q_1}$$

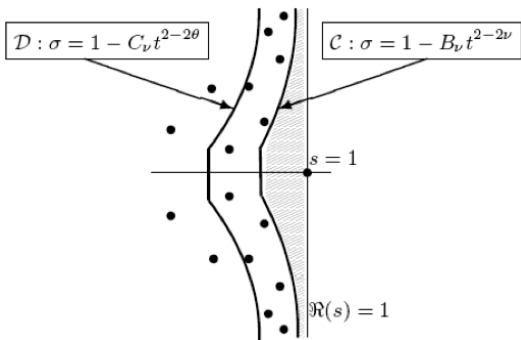
The **poles** of $\Lambda(s)$ close to $\Re s = 1$

are related to good **rational approximations** of $\log p_2 / \log p_1$

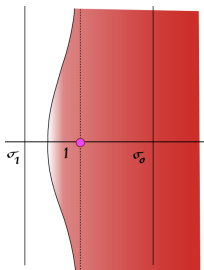
A number x is **diophantine** with exponent μ if there exists C for which

$$\left| x - \frac{a}{b} \right| > \frac{C}{b^\mu} \quad \forall a, b \in \mathbb{Z}$$

Consider $\Lambda(s) = 1/(1 - p_1^s - p_2^s)$. If $\log p_1 / \log p_2$ is μ -diophantine, then, for any θ, ν with $\theta < \mu < \nu$, the tameness region has an hyperbolic form: [Flajolet-Roux-V. 2010]



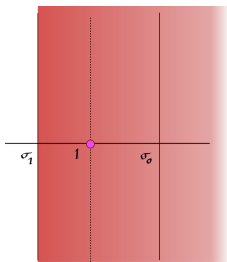
Different possible regions on the left of $\Re s = 1$ where $\Lambda(s)$ is tame.



Situation 1

Hyperbolic region

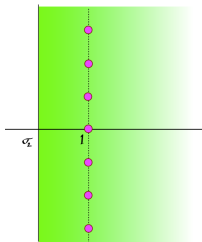
Arithmetic condition



Situation 2

Vertical strip

Geometric condition



Situation 3

Vertical strip with holes

Periodicity condition

For which **general dynamical** sources do these different situations occur?

- S3 **occurs only** if the source is conjugated to a **simple** source.
- S2 occurs when “ the branches are **not** too often of the **same shape**”.
- S1 occurs if a extension of the following condition holds:
 - “there **exists** a ratio $\log p_i / \log p_j$ which is **badly approximable** by rationals”.

Situation 2- Existence of a vertical strip where $\Lambda(s)$ is tame

There exists a condition, the condition UNI, which expresses that
“ the branches of the dynamical system are not **too often** of the **same form**”

Theorem [Dolgopyat-Baladi-Cesaratto-V]. For a **good** dynamical system [complete, expansive, with bounded distortion], which satisfies the **condition UNI**, there exists a **vertical strip** where its Dirichlet series $\Lambda(s)$ is **tame**.

Dolgopyat (98) proves the result for the **plain** transfer operator, in the case of a **finite** number of branches

- Baladi and V. (03) extend the result for an **infinite** number of branches
- Cesaratto and V. (09) extend the result to the **secant** transfer operator.

Description of the UNI Condition

A **distance** Δ between two inverse branches of the same depth:

$$\Delta(h, k) := \inf_{x \in \mathcal{I}} \Psi'_{h,k}(x), \quad \text{with} \quad \Psi_{h,k}(x) := \log \frac{|h'(x)|}{|k'(x)|}$$

Contraction ratio ρ . $\rho := \limsup (\{\max |h'(x)|; h \in \mathcal{H}^n, x \in \mathcal{I}\})^{1/n}$.

Probability \Pr_n on $\mathcal{H}^n \times \mathcal{H}^n$. $\Pr_n(h, k) := |h(\mathcal{I})| \cdot |k(\mathcal{I})|$

For a system \mathcal{C}^2 -conjugated with a piecewise-affine system :

For any $\hat{\rho}$ with $\rho < \hat{\rho} < 1$, for any n , $\Pr_n[\Delta < \hat{\rho}^n] = 1$

Condition UNI.

For any $\hat{\rho}$ with $\rho < \hat{\rho} < 1$, for any n , $\Pr_n[\Delta < \hat{\rho}^n] \ll \hat{\rho}^n$

Situation 3- Existence of a hyperbolic region where $\Lambda(s)$ is tame

The **condition DIOP** extends the arithmetic condition

“There exists a ratio $\log p_i / \log p_j$ which is **diophantine**”

to the general dynamical case.

Condition DIOP: There exists a ratio $c(h, k)$ which is **diophantine**.

For a complete system, each branch h has a fixed point denoted by h^* .

We consider the ratios between the derivatives $|h'(h^*)|$ at the fixed point

$$c(h, k) := \frac{\log |h'(h^*)|}{\log |k'(k^*)|}$$

Theorem [Dolgopyat-Roux-V.] For a **good** dynamical system [complete, expansive, with bounded distortion], which satisfies the **condition DIOP**, there exists an **hyperbolic region** where $\Lambda(s)$ is **tame**.

Dolgopyat (98) proves the result for the **plain** transfer operator, in the case of a **finite** number of branches – Roux and V. (2010) extend the result : for an **infinite** number of branches and for the **secant** transfer operator.

Plan of the talk.

- General motivations
- Models for a source
- The Dirichlet generating function of the source
- Conclusion and possible extensions.

Conclusions.

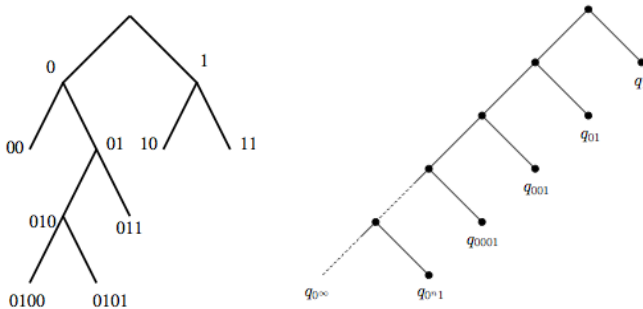
Description of a framework which

- unifies the analyses for text algorithms and searching/sorting algorithms
- provides a general model for sources
- shows the importance of the Dirichlet generating functions
- explains the importance of tameness for sources
- defines a natural subclass of sources, the dynamical sources
- provides sufficient conditions for tameness of dynamical sources
- provides probabilistic analyses for algorithms built on tame sources.

Possible extensions and work in progress

I- Classification of sources

- Place of dynamical sources amongst general sources:
- A dynamical source = limit of Markov chains with increasing order?
- Comparing dynamical sources with Markov chains of variable length



Possible extensions and work in progress

II– Realistic analyses of other algorithms and other structures

- Analysis of other sorting algorithms
 - Analysis of Insertion Sort easy
 - Analysis of QuickSelect already done
 - And Selection algorithm ?
- Analysis of the DST structure?

