# Digital Trees and Memoryless Sources: from Arithmetics to Analysis

Philippe Flajolet, Mathieu Roux, Brigitte Vallée

## AofA 2010, Wien

= a data structure for dynamic dictionaries

- TOP-DOWN construction: Set  $\mathcal{E}$  is split into  $\mathcal{E}_a, \dots, \mathcal{E}_z$ , according to initial letter; continue with next letter; stop when elements are separated.
- INCREMENTAL construction:

start with the empty tree and insert elements of E one after the other.

A = Finite alphabet W = infinite sequences E : W<sup>n</sup> -> tree



## What does a trie look like?



# What does a trie look like?

Expected size seems to be <u>asymptotically linear</u>.

Convergence to asymptotic regime seems to be <u>fast</u>.

But...Things are not quite as they seem!



### Probabilistic model: Memoryless sources

- A finite alphabet  $\mathcal{A} = \{a_1, \ldots, a_r\}$ .
- Letters drawn independently to form words from  $\mathcal{W} = \mathcal{A}^{\infty}$ :

 $\mathbb{P}(a_j) = p_j.$ 

• Words drawn independently: model is  $\mathcal{W}^n$ .

- Want fixed number, *n* items, to build the trie.
- Often use  $N = \mathcal{P}oisson(x)$  items:

$$\mathbb{P}(N=n)=e^{-x}\frac{x^n}{n!}.$$





Expect ( $\pm$ elementarily)  $\mathcal{P}(x) \approx \text{fixed-}n$ , when  $x \approx n$ .

## Memoryless sources (Bernoulli)

- I965: Knuth & De Bruijn analyse binary tries, with Pr(0)=Pr(1)=1/2, showing oscillations.
- I973: Knuth discusses biased bit models, including golden-section case [Ex 5.2.2-53]
- I986: Fayolle-F-Hofri exhibit periodicity criterion, extended by, e.g., Schachinger [2000]; Jacquet-Szpankowski-Tang [2001]
- Igo-2000: Convergence to asymptotic regime often wrongly assumed to be fast. Caveats by Schachinger (~2000).
- 2010; this paper: convergence to asymptotic regime is very slow and depends on fine arithmetic properties of probabilistic model.

### Definition

The probability vector 
$$(p_1, \ldots, p_r)$$
 is periodic if  
— all ratios  $\frac{\log p_j}{\log p_k}$  are rational. (E.g.,  $\frac{\log p_2}{\log p_1} \in \mathbb{Q}$ ; binary alph.)

Theorem (Periodic sources; folklore)
$$H = -p_1 \log p_1 - \cdots - p_r \log p_r$$
Expected size  $\overline{S_n}$  is, with  $\Phi$  a smooth periodic function: $\overline{S}_n = \frac{n}{H} + n\Phi(\log n) + O(n^{1-A}), \quad A > 0.$ 

 $\implies$  Oscillations (O(n)), plus good error term.

• These cases are **exceptional**: the  $p_j$  are *algebraic numbers*. Such families are a denumerable set; hence have measure 0.

### Definition

The probability vector 
$$(p_1, \ldots, p_r)$$
 is aperiodic if  
— at least one ratio  $\frac{\log p_j}{\log p_k}$  is irrational. (E.g.,  $\frac{\log p_2}{\log p_1} \notin \mathbb{Q}$ ; binary a.)

Theorem (Aperiodic sources; this paper)  $H = -p_1 \log p_1 - \cdots - p_r \log p_r$ 

Expected size  $\overline{S_n}$  is, for "diophantine sources" (generic case)

$$\overline{S}_n = \frac{n}{H} + O\left(n\exp(-\sqrt[\theta]{\log n})\right), \quad \theta > 1.$$

This is better than  $n/(logn)^a$ , any a; much worse than  $n^{1-\epsilon}$ , any  $\epsilon$ . • For remaining "Liouvillean sources" (rare), error term can come arbitrarily close to o(n).

 $\implies$  No oscillation, but poor error term.

• This case is **generic**: it has has measure 1.

# 1. Basics

Fundamental intervals
+ Mellin
= Formal analysis

## [Vallée 1997++]

View <u>source model</u> in terms of <u>fundamental intervals</u>:

**w** -> p<sub>w</sub>

 Size = Number of places occupied by at least two prefixes

Mellinize ->...





## The Mellin transform





(It exists in strips of  $\mathbb{C}$  determined by growth of f(x) at  $0, +\infty$ .) **Property 1.** Factors *harmonic sums*:

$$\sum_{(\lambda,\mu)} \lambda f(\mu x) \quad \stackrel{\mathcal{M}}{\leadsto} \quad \left(\sum_{(\lambda,\mu)} \lambda \mu^{-s}\right) \cdot f^*(x).$$

**Property 2.** Maps asymptotics of f on singularities of  $f^*$ :

$$f^{\star} pprox rac{1}{(s-s_0)^m} \implies f(x) pprox x^{-s_0} (\log x)^{m-1}.$$

Proof of  $P_2$  is from Mellin inversion + residues Singularities?

# Lambda(s)

$$\begin{cases} E_{\mathcal{P}(x)}[\text{Size}] = \sum_{w \in \mathcal{A}^*} g(p_w x) \\ g(x) = 1 - (1+x)e^{-x}. \end{cases}$$

$$\int S^{\star}(s) = -(s+1)\Gamma(s)\Lambda(s)$$
  
 $\Lambda(s) := \sum_{w} p_{w}^{-s}$ 



### Geometry of the poles of

$$\Lambda(s) = \frac{1}{1-\textbf{p}_1^s-\cdots-\textbf{p}_r^s}$$

# 2. Geometry of poles

Poles are associated with simultaneous approximations to logs of probabilities

### Ø Distinguish:

- -- <u>Diophantine</u> = badly approximable (generic);
- -- <u>Liouvillean</u> = unusally well approximable (rare)

Poles of  $\Lambda(s)$  near  $\Re(s) = 1$ 

$$\Lambda(s) = \frac{1}{1 - p_1^s - \dots - p_r^s}$$

• Look for s: 
$$p_1^s + p_2^s = 1$$
,  $s = \sigma + it$ .

$$p_1^{\sigma} p_1^{it} + p_2^{\sigma} p_2^{it} = 1, \qquad p_1 + p_2 = 1.$$

Implies  $p_1^{it} \approx 1$  and  $p_2^{it} \approx 1$ ; i.e.,  $t \approx \frac{2\pi}{\log p_1} q_1$  and  $t \approx \frac{2\pi}{\log p_2} q_2$ .

$$rac{\log p_2}{\log p_1} pprox rac{q_2}{q_1}.$$

**Pole** of  $\Lambda(s) \implies$  **"good" rational approximation** to  $\frac{\log p_2}{\log p_1}$ .

For general  $(p_1, \ldots, p_r)$ , must have a **common denominator**  $q_1$ :  $\forall j : q_1 \frac{\log p_j}{\log p_1}$  is a near-integer.

## Poles of $\Lambda(s)$ near $\Re(s) = 1$

 $\beta = (\beta_1, \ldots, \beta_r) \in \mathbb{R}^r$ ; fix a norm  $\|\cdot\|$  on  $\mathbb{R}^r$ .

 $\{x\}$  = centred fractional part;  $\|\{\beta\}\|$  is distance to nearest integer lattice point.

Look at "record" approximants; measure quality by f(t).

### Definition

• *Q* is a **Best Simultaneous Approximant Denominator** (BSAD), if  $\|\{Q\beta\}\| < \|\{q\beta\}\|$ , for all q < Q.

• f(t), the approximation function, is staircase and  $f(t) = \frac{1}{\|\{Q^-\beta\}\|}$ ., if  $Q^-, Q^+$  are the BSADs that frame t. Thus:



For a probability vector  $(p_1, \ldots, p_r)$ :

- **Periodic sources** (All ratios of logs are in  $\mathbb{Q}$ )
- Aperiodic sources (Some ratios  $\notin \mathbb{Q}$ ):
  - Diophantine: approximation function f(t) is polynomial;
     optimal exponent is known as *irrationality measure*;
  - Liouvillean: approximation function f(t) is superpolynomial.

— Scalars  $\pi$ , e, tan(1),  $\sqrt[3]{2}$ ,  $\zeta(3)$ , log 5, ... are Diophantine. Logs of rational and algebraic numbers are Diophantine. Also numbers with bounded continued fraction quotients, ...

— Numbers with very fast-converging sums, e.g.,  $\sum 2^{-2^n}$ , are Liouvillean.

#### Theorem

If  $(p_1, \ldots, p_r)$  is Diophantine, zeros are well-separated from  $\Re(s)$ :

- All zeros are to the *left* of a pseudo-hyperbola;
- Infinitely many zeros are to the right of a pseudo-hyperbola.

### Theorem

If  $(p_1, \ldots, p_r)$  is Liouvillean, zeros come closer to  $\Re(s) = 1$ :

- All zeros are to the left of a curve  $1 1/F_{-}(t)$ ;
- Infinitely many zeros are to the right of  $-1 + 1/F_+(t)$ .

 $F_{-}(t), F_{+}(t)$  are dictated by approximation functions of  $(\log p_j)/(\log p_k)$ .



### Proofs

• Pole of  $\Lambda(s) \implies$  "good" rational approximation to  $(\log p_j)(\log p_k)$ .

— Follow sketch above and develop properties of "ladders".

• **"Good" rational approximation** to  $(\log p_j)/(\log p_k) \implies$  **Pole** of  $\Lambda(s)$ 

— use analytic, multivariate **Implicit** Function Theorem,  $\Re(s) \approx 1$ ;  $u_j \approx 0$ :

$$1 - p_1^s p_1^{iu_1} - \cdots p_r^s p_r^{iu_r} = 0.$$



Fractal Geometry, Complex Dimensions and Zeta Functions

# 3. Inverse Mellin analysis

Make use of integration contour that avoids poles

Estimate global contribs:
 pole-free region matters

Poles are well-separated



# 4. Tries and QuickSort

Applies to size of tries

& almost anything that contains Lambda(s).

Diophantine => error terms are exp-of-root-of-log



Liouvillean => error terms are
 o(n) and very close to O(n)

#### Theorem

Consider aperiodic Diophantine probabilities with irrationality exponent  $\mu$ .

$$\begin{cases} trie \ size; & \overline{S}_n = \frac{n}{H} + \mathbf{n} \Phi(\mathbf{n}) \\ trie \ pathlength: & \overline{S}_n = \frac{1}{H} n \log n + Cn + \mathbf{n} \Phi(\mathbf{n}) \\ symbol-cost, \ Quicksort: & \overline{S}_n = \frac{2}{H} n \log^2 n + Cn \log n + C'n + \mathbf{n} \Phi(\mathbf{n}), \end{cases}$$

where error term is, for any  $\theta > \mu$ :

$$\mathbf{\Phi}(\mathbf{x}) = O\left(\exp\left[-(\log n)^{1/\theta}\right]\right),$$

## Makes precise or improves on results of Clément, Fill, Flajolet, Jacquet, Janson, Szpankowski, Vallée,...

## Source models

### memoryless

ø periodic: good error terms aperiodic: generally (very) bad error terms (us!) Diophantine versus Liouvillean Markov; cf Szpa+Jacquet+Tang: similar (?) dynamical: Vallée + Cl-F-Vallée; cf Dolgopyat, B-V.
 @ general: à la Vallée-Clément-Fill-F.

# Numerics

 (Proved for Poisson; transfers to fixed-size)
 Initial oscillations often not seen numerically, for small n; but they matter asymptotically



