

# *Algorithms Seminar, 2000–2001*

Frédéric CHYZAK, éditeur scientifique

N ° 4406

Mars 2002

THÈME 2



*R*apport  
de recherche





## Algorithms Seminar, 2000–2001

Frédéric CHYZAK, éditeur scientifique

Thème 2 — Génie logiciel  
et calcul symbolique  
Projet Algo

Rapport de recherche n° 4406 — Mars 2002 — 197 pages

**Abstract:** These seminar notes constitute the proceedings of a seminar devoted to the analysis of algorithms and related topics. The subjects covered include combinatorics, symbolic computation, asymptotic analysis, computational biology, and average-case analysis of algorithms and data structures.

**Key-words:** combinatorics, symbolic computation, analysis of algorithms, probabilistic methods, computational biology

*(Résumé : tsvp)*

## Séminaire algorithmes, 2000–2001

**Résumé :** Ces notes de séminaires constituent les actes, le plus souvent en anglais, d'un séminaire consacré à l'analyse d'algorithmes et aux domaines connexes. Les thèmes abordés comprennent : la combinatoire, le calcul symbolique, l'analyse asymptotique, la biologie computationnelle et l'analyse en moyenne d'algorithmes et de structures de données.

**Mots-clé :** combinatoire, calcul symbolique, analyse d'algorithmes, méthodes probabilistes, biologie computationnelle

# ALGORITHMS SEMINAR

## 2000–2001<sup>1</sup>

*Frédéric Chyzak*  
(*Editor*)

### Abstract

These seminar notes constitute the proceedings of a seminar devoted to the analysis of algorithms and related topics. The subjects covered include combinatorics, symbolic computation, probabilistic methods, and average-case analysis of algorithms and data structures.

This is the tenth in our series of seminar proceedings. The previous ones have appeared as INRIA Research Reports numbers 1779, 2130, 2381, 2669, 2992, 3267, 3504, 3830, and 4056. The content of these annual proceedings consists of summaries of the talks, usually written by a reporter from the audience.<sup>2</sup> The primary goal of the seminar is to cover the major methods for the average-case analysis of algorithms and data structures. Neighbouring topics of study are combinatorics, symbolic computation, asymptotic analysis, probabilistic methods, and computational biology.

The study of combinatorial objects—their description, their enumeration according to various parameters—arises naturally in the process of analysing algorithms that often involve classical combinatorial structures like strings, trees, graphs, and permutations. Beside the traditional topics of combinatorics of words and algorithmics on words, over the years an increasing interest has been given in the seminar to biological applications of combinatorics. Symbolic computation, and in particular computer algebra, plays an increasingly important role in these areas. It provides a collection of tools that allows one to attack complex models of combinatorics and the analysis of algorithms via *generating functions*; at the same time, it inspires the quest for developing ever more systematic solutions and decision procedures for the analysis of well-characterized classes of problems. Our seminar shares a large part of its audience with ALÉA, a working group dedicated to the analysis of algorithms and to the analysis of properties of discrete random structures. This year's workshop, ALEA'2001, started with a series of short courses on various aspects of probability and enumerative combinatorics. It was decided to include lecture notes for the courses in the seminar proceedings.

The thirty-one articles included in this book represent snapshots of current research in the areas mentioned above. A tentative organization of their contents is given below. Three ALÉA lecture notes follow.

### PART I. COMBINATORICS

Sand piles are integer partitions that can be obtained from a column of grains by moving grains from left to right according to a specific set of rules; their enumeration for several models is attacked in [1]. The abelian sand-pile model is a different, 2-dimensional model, with an underlying group-theoretic structure; several algorithms to determine the identity of this group, which presents fractal aspects, are considered in [2]. The “*s*-tennis ball problem” is a combinatorial model of a tennis

---

<sup>1</sup>Partially supported by the Future and Emerging Technologies programme of the EU under contract number IST-1999-14186 (ALCOM-FT).

<sup>2</sup>The summaries for the past nine years are available on the web at the URL <http://algo.inria.fr/seminars/>.

player who receives  $s$  (labeled) new balls at each of his services. The problem of the enumeration of all orders balls may be served is explored in [3] together with a generalization in which several balls are served simultaneously. The classical coupon-collector problem is extended in [4] to the case where the collector shares his harvest with other members of his phratry. In a different direction, effective manipulations of sums is a very active research topic. An old method by Mac Mahon for the evaluation of sums over indices constrained by linear homogeneous diophantine inequalities and equations is revitalized in [5] and given an algorithmic status. Another intriguing type of expansions of  $q$ -series is the topic of [6] and is another example of combinatorics that received recent nice symbolic developments. A combinatorial problem on permutation statistics is solved by computer algebra calculations in [7].

- [1] Enumeration of Sand Piles. *S. Corteel.*
- [2] On the Group of a Sandpile. *D. Rossin.*
- [3] The Tennis Ball Problem. *D. Merlini.*
- [4] Hyperharmonic Numbers and the Phratry of the Coupon Collector. *D. Foata.*
- [5] Mac Mahon's Partition Analysis Revisited. *P. Paule.*
- [6] Engel Expansions of  $q$ -Series. *P. Paule.*
- [7] Eulerian Calculus: a Technology for Computer Algebra and Combinatorics. *D. Foata.*

## PART II. ANALYSIS OF ALGORITHMS AND COMBINATORIAL STRUCTURES

Probabilistic methods are at the heart of the analysis of several combinatorial structures or processes on combinatorial structures: the asymptotic shape of “large” random partitions is studied in [8]; the covering time of random walks on graphs satisfying self-avoiding properties is addressed in [9]; various tail bounds for occupancy problems are derived in [10], with applications to the determination of the conjectured satisfiability threshold in the random  $k$ -sat problem. Dynamical systems are a different approach used in [11] to analyse parameters of the data structure of Patricia tries. Pattern matching methods and their analysis are surveyed in [12]. More recent is the interest of our seminar to biological applications of combinatorics, and in particular to the crucial problem in genomic analysis of distinguishing “biologically significant” signals in sequences from those that are part of the ground noise. This problem has been discussed in two talks this year, both from the biologist's point of view [13] and from the combinatorial point of view [14]. NP-hard problems cannot be solved exactly *and* efficiently at the same time, and polynomial-time algorithms for these problems can only return approximate solutions. A general method to design polynomial-time approximate algorithms for solving such problems is described in [15], together with a survey on applications. Scheduling loads between  $n$  agents trying to achieve a global goal is a difficult task; the case when no communication is allowed between agents is studied in [16].

- [8] Asymptotics for Random Combinatorial Structures. *A. Dembo.*
- [9] Random Walks and Heaps of Cycles. *Ph. Marchal.*
- [10] Tail Bounds for Occupancy Problems. *P. Spirakis.*
- [11] Patricia Tries in the Context of Dynamical Systems. *J. Bourdon.*
- [12] New and Old Problems in Pattern Matching. *W. Szpankowski.*
- [13] Genome Analysis and Sequences with Random Letter Distribution. *M. Termier.*
- [14] Random Sequences and Genomic Analysis. *A. Denise.*
- [15] The Primal-Dual Schema for Approximation Algorithms: Where Does It Stand, and Where Can It Go? *V. Vazirani.*
- [16] Distributed Decision Making: The Case of No Communication. *P. Spirakis.*

### PART III. COMPUTER ALGEBRA AND APPLICATIONS

After integers, whose long-studied question of factorization is surveyed in [17], the most fundamental data structures in computer algebra are (sorted) polynomials and series. Fast algorithms for them are discussed in [18] in the univariate case, and in [19] in the case of multivariate series. Effective manipulations of sums by a counterpart for recurrences to the theory of differential forms are described in [20]. A new linear algebra algorithm for matrices with entries in skew rings is presented in [21], and has many applications to the solving of linear ordinary differential equations. Examples of Hamiltonian systems are considered in [22], which is a showcase for algorithms to solve linear ordinary differential equations. Two talks discuss effective methods for control theory, a new topic in the seminar: effective tests to capture structural properties of control systems are described in [23], leading to Gröbner basis calculations, while series expansions and Newton iteration are used in [24] to address the question of observability.

[17] Thirty Years of Integer Factorization. *F. Morain.*

[18] Variations on Computing Reciprocals of Power Series. *A. Schönhage.*

[19] Fast Multivariate Power Series Multiplication in Characteristic Zero. *G. Lecerf.*

[20] A Tutorial on Closed Difference Forms. *B. Zimmermann.*

[21] Transformations Exhibiting the Rank for Skew Laurent Polynomial Matrices. *M. Bronstein.*

[22] A Criterion for Non-Complete Integrability of Hamiltonian Systems. *D. Boucher.*

[23] Effective Algebraic Analysis in Linear Control Theory. *A. Quadrat.*

[24] Effective Test of Local Algebraic Observability — Applications to Systems and Control Theory. *A. Sedoglavic.*

### PART IV. PROBABILISTIC METHODS

Brownian motion is a central tool in probability. The area under a variant, the reflected Brownian bridge, is analysed in details in [25]. Brownian motion can be viewed as the limit of some simple random walk on integers. Two talks study other kinds of random walks: conjectures on the frequency of visits of points in a planar random walk are proved in [26]; random walks on groups are viewed in [27] from the point of view of probability theory, statistical physics, ergodic theory, harmonic analysis, and group theory. A model of queues is studied in [28], together with a link to random matrices. The information-theoretic problem of source coding had already been considered in great generality over the past years in the seminar, the key question being to analyse the redundancy of a source; in [29], different models for redundancy are detailed, and a generalized Shannon code is introduced in order to solve the minimax redundancy problem for a single memoryless source.

[25] Reflected Brownian Bridge Area Conditioned on its Local Time at the Origin. *G. Louchard.*

[26] Cover Time and Favourite Points for Planar Random Walks. *A. Dembo.*

[27] Introduction to Random Walks on Groups. *Y. Guivarc'h.*

[28] Random Matrices and Queues in Series. *Y. Baryshnikov.*

[29] Information Theory by Analytic Methods: The Precise Minimax Redundancy. *W. Szpankowski.*

### PART V. ASYMPTOTICS AND ANALYSIS

Linear functional equations and their special functions solutions are a common denominator to various topics addressed in the seminar—combinatorics, the analysis of algorithms, computer algebra. Two talks in this year's seminar are analytic studies of some properties of solutions of linear functional equations: connection formulae for a  $q$ -analog to the Bessel function equation are derived

in [30], generalizing the asymptotic expansion of the Bessel  $J_\nu$  functions; the Borel summation technique is used in [31] to recover convergent representations for everywhere divergent formal power series solutions to some “irregular singular” problems coming from differential equations.

[30] On Jackson’s  $q$ -Bessel Functions. *C. Zhang.*

[31] On the Convergence of Borel Approximants. *D. Lutz.*

## PART VI. ALEA’2001 LECTURE NOTES

General algebraic methods to solve combinatorial enumeration problems with nice decomposability properties are described in [32], where a central role is played by generating functions. The modern view of combinatorial analysis also makes enumerative generating functions its central objects: the singularity structure of the latter, now regarded as analytic functions of the complex variable, contains all information essential to the asymptotic enumeration of the combinatorial objects. The basics of this approach, nowadays known by the name of *analytic combinatorics*, are introduced in [33]. Connections between Brownian motion and related processes (meander, bridge, excursion) on the one hand, and combinatorial objects like Dyck words, trees, bi-sorted permutations, combinatorial and algorithmic problems like hashing and the parking problems on the other hand make the partial review of the numerous properties of Brownian motion proposed in [34] very welcome.

[32] Enumerative Combinatorics: Combinatorial Decompositions and Functional Equations. *M. Bousquet-Mélou.*

[33] Symbolic Enumerative Combinatorics and Complex Asymptotic Analysis. *Ph. Flajolet.*

[34] Aléa discret et mouvement brownien (*Discrete Randomness and Brownian Motion*). *Ph. Chassaing.*

*Acknowledgements.* The lectures summarized here emanate from a seminar attended by a community of researchers in the analysis of algorithms, from the Algorithms Project at INRIA (the organizers are Philippe Flajolet and Bruno Salvy) and the greater Paris area. The editor expresses his gratitude to the various people who have actively supported this joint enterprise and offered to write summaries. Thanks are also due to the speakers and to the authors of summaries. Many of them have come from far away to attend a seminar and kindly accepted to write the summary. We are also greatly indebted to Virginie Collette for making all the organization work smoothly.

The editor,  
F. CHYZAK



**Part I**

**Combinatorics**



## Enumeration of Sand Piles

Sylvie Corteel

PRISM, Université de Versailles - Saint-Quentin-en-Yvelines (France)

October 16, 2000

Summary by Michel Nguyen-Thé

### Abstract

Sand piles are integer partitions that can be obtained from a column of  $n$  grains by moving grains from left to right according to rules defined by a model. We try to better understand the structure of those objects by decomposing and counting them. For the model introduced by Goles, Morvan, and Phan, we find generating functions according to area, height, and width. We establish a bound for the number of the sand piles consisting of  $n$  grains in  $\text{IPM}(k)$  for large  $n$ . We present the series according to area and height for Phan's model  $L(\theta)$ . We introduce a more general model, where grains can also go to the left, that we call Frobenius sand piles. (Joint work of S. Corteel with D. Gouyou-Beauchamps (LRI, Orsay)).

### 1. Preliminaries and $\text{SPM}(k)$ Model

After the necessary basic concepts, we present here the simplest model of sand pile, i.e., the  $\text{SPM}(k)$  model, from which all other models are derived.

**1.1. Definitions.** A *sand pile* made of  $n$  grains is a partition of the integer  $n$ . A *partition* of an integer  $n$  is a non-increasing sequence of positive integers  $\lambda = (\lambda_1, \dots, \lambda_l)$ . The  $\lambda_i$  are called the *parts* of the partition. The *area* of the sand pile is the sum  $|\lambda| = \lambda_1 + \dots + \lambda_l = n$ . The *height* of the sand pile is the number  $h(\lambda) = l$  of parts of the partition. For any partition  $\lambda$ , we will consider that  $\lambda_i = 0$  for  $i < 1$  and  $i > h(\lambda)$ . The *width*  $w(\lambda)$  of the sand pile  $\lambda$  is the largest part  $\lambda_1$ . The *Ferrers diagram* of a partition  $\lambda$  is a drawing of  $\lambda$  such that the  $i$ th column is a pile of  $\lambda_i$  packed squares (called grains). The rows are labelled from bottom to top. The *conjugate*  $\lambda'$  of  $\lambda$  is the partition whose  $i$ th part is the number of squares in the  $i$ th column of the Ferrers diagram of  $\lambda$ .

Let  $\pi = (\pi_1, \dots, \pi_l)$  be a sand pile and  $\pi' = (\pi'_1, \dots, \pi'_{\pi_1})$  be its conjugate. The moves of the sand grains are of two types (see Figure 1):

1. *Vertical rule:* a grain can move from column  $i$  to column  $i + 1$  if  $\pi'_i - \pi'_{i+1} \leq 2$ , so that

$$(\pi'_1, \dots, \pi'_i, \pi'_{i+1}) \quad \text{is replaced with} \quad (\pi'_1, \dots, \pi'_i - 1, \pi'_{i+1} + 1, \dots, \pi'_{\pi_1}).$$

2. *Horizontal rule:* a grain can move from column  $i$  to column  $j$  if  $j > i + 1$  and  $\pi'_i - 1 = \pi'_{i+1} = \dots = \pi'_j = \pi'_{j+1} + 1$ , so that

$$(\pi'_1, \dots, \pi'_i, \pi'_{i+1}, \dots, \pi'_j, \pi'_{j+1} + 1, \dots, \pi'_{\pi_1}).$$

The shift is said to have length 0 or  $j - i - 1$ , respectively.

In the  $\text{SPM}(k)$  model (*Sand Pile Model*), introduced by Goles and Kiwi [3], the initial configuration is made of one column of  $n$  grains, and the only available rule is the vertical rule.

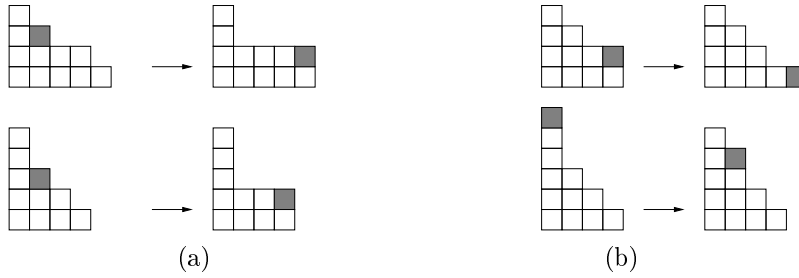


FIGURE 1. (a) Application of horizontal rule to  $(5,4,2,1)$  and  $(4,3,2,1,1)$ ; (b) application of vertical rule to  $(4,4,2,1)$  and  $(4,3,2,1,1,1)$ .

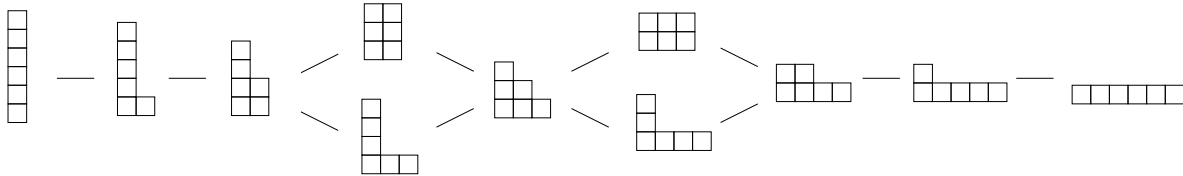


FIGURE 2.  $L_B(6)$

1.2. **Generating function.** Let  $p(n, k)$  denote the number of partitions of  $n$  of width  $k$ . Then:

$$F(q, x) = 1 + \sum_{n, k \geq 1} p(n, k)q^n x^k = xqF(q, x) + F(q, xq) = \prod_{i=1}^{\infty} \frac{1}{1 - xq^i}.$$

1.3. **Example of bijection.** There is a bijection between partitions with odd parts and partitions with distinct parts, as is reflected by the generating functions identity

$$\prod_{i=1}^{\infty} \frac{1}{1 - q^{2i+1}} = \prod_{i=1}^{\infty} \frac{1 - q^{2i}}{1 - q^i} = \prod_{i=1}^{\infty} (1 + q^i).$$

1.4. **Order on partitions.** Let  $\mu = (\mu_1, \mu_2, \dots)$  and  $\lambda = (\lambda_1, \lambda_2, \dots)$  be two partitions of  $n$ . We say that  $\mu \geq \lambda$  if and only if there exists a sequence of moves of  $n$  induced by the rules to go from  $\mu$  to  $\lambda$ . In the SPM( $k$ ) model, this order is equivalent to the dominance order  $L_B(n)$  [1] on the conjugates:  $\mu \geq \lambda$  if and only if  $\sum_{i=1}^j \mu'_i \geq \sum_{i=1}^j \lambda'_i$  for all  $j \geq 1$ . Brylawski [1] showed:

**Theorem 1.** *Let  $n$  be an integer. The set of partitions of  $n$  with the previously defined order is a lattice, where the maximal element is  $(1, 1, \dots, 1)$ , and the minimal element is  $(n)$ . Moreover, the infimum and the supremum of two partitions can be respectively defined as follows:*

1.  $\inf(\mu, \lambda) = \pi$  such that  $\pi'_j = \min\left(\sum_{i=1}^j \mu'_i, \sum_{i=1}^j \lambda'_i\right) - \sum_{i=1}^j \pi'_i$  for all  $j \geq 1$ .
2.  $\sup(\mu, \lambda) = \alpha$  such that  $\alpha'_j = \max\left(\sum_{i=1}^j \mu'_i, \sum_{i=1}^j \lambda'_i\right) - \sum_{i=1}^j \alpha'_i$  for all  $j \geq 1$ .

In Figure 2, the maximal element  $(1, 1, 1, 1, 1, 1)$  is on the left.

*Length of a maximal chain.* The length of a maximal chain is greater than  $2n - 3$  [1], and smaller than  $2\binom{l+1}{3} + lj + 1$  [3], where  $l$  and  $j$  are defined by  $n = j + l(l + 1)/2$  and  $0 \leq j \leq l$ . For  $n = 6$ , the two bounds are equal to 9, which shows that they both can be attained. The corresponding maximal chain is displayed in Figure 2.

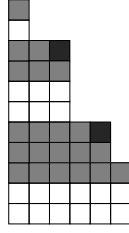


FIGURE 3. Decomposition of a sand pile in IPM(1).

## 2. IPM( $k$ ) Model

A more realistic generalization of SPM( $k$ ) model limits the lengths of the possible horizontal shifts of a grain.

**2.1. Definition.** In [4], the sand piles in IPM( $k$ ) are characterized in the following way:

**Proposition 1.** A sand pile in IPM( $k$ ) is a partition  $\pi = (\pi_1, \dots, \pi_l)$  of  $n$  such that

- for  $1 \leq i \leq l$ ,  $0 \leq \pi_i - \pi_{i+1} \leq k + 1$ ;
- for any  $i < j$  with  $\pi_i - \pi_{i+1} = k + 1$  and  $\pi_j - \pi_{j+1} = k + 1$ , there exists  $z$  with  $i < z < j$  such that  $\pi_z - \pi_{z+1} < k$ .

## 2.2. Generating functions.

### 2.2.1. Area and height.

**Theorem 2.** The generating function  $S_k(q, x)$  of IPM( $k$ ) sand piles, with  $q$  and  $x$  respectively counting area and height, satisfies

$$S_k(q, x) = 1 + \sum_{\pi \in \text{IPM}(k)} x^{l(\pi)} q^{|\pi|} = 1 + \sum_{i=1}^k \frac{xq^i}{1 - xq^i} S_k(q, xq^i) + xq^{k+1} S_k(q, xq^k).$$

*Proof.* A sand pile in IPM( $k$ ) is either the empty partition, or a partition in IPM( $k$ ) where one duplicates  $i$  times the highest column and adds to it at least one part  $i$  ( $1 \leq i \leq k$ ), or a partition in IPM( $k$ ) where one duplicates  $k$  times the highest column and adds to it one part of length  $k + 1$ . This decomposition yields the last expression for  $S_k(q, x)$  in the statement of the theorem, after noting that  $S_k(q, xq^r)$  is the generating function obtained by duplicating  $r$  times the highest column in each sand pile.  $\square$

Note the particular cases:

$$S_1(q, x) = 1 + \sum_{n \geq 1} x^n q^{n(n+1)/2} \prod_{i=1}^n \left( q + \frac{1}{1 - xq^i} \right); \quad S_\infty(q, x) = \prod_{i=1}^{\infty} \left( q + \frac{1}{1 - xq^i} \right).$$

### 2.2.2. Area and width.

**Theorem 3.** The generating function  $S_k(q, y)$  of IPM( $k$ ) sand piles, with  $q$  and  $y$  respectively counting area and width, satisfies:

$$S_k(q, y) = \left( \frac{1 - (yq)^{k+1}}{1 - yq} + y^k q^{k-1} \right) S_k(q, yq) + y^k q^{k-1} (S_k(q, yq) - S_k(q, yq^2)).$$

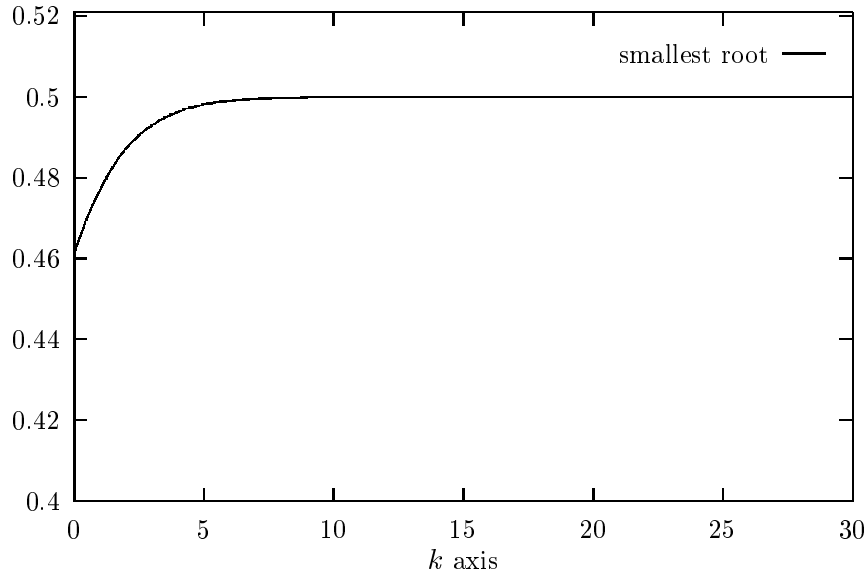


FIGURE 4. Evolution of the smallest root of the polynomial  $1 - 2x - 3x^{k+2} + x^{k+3} + x^{k+1}$ .

In particular:

$$S_1(q, y) = 1 + \sum_{n \geq 1} y^n q^{n(n-1)/2} \prod_{i=1}^n \frac{1+q-q^{i-1}}{1-q^i}.$$

2.2.3. *Height and width.* Let  $p_k(h, w)$  be the number of sand piles in  $\text{IPM}(k)$  of height  $h$  and width  $w$  and  $P_{k,h}(y)$  the generating function  $\sum_{w \geq 0} p_k(h, w)y^w$ .

**Theorem 4.** *The generating function  $P_{k,h}(y)$  follows the recurrence:*

$$P_{k,0}(y) = 1; \quad P_{k,1}(y) = y \frac{1-y^{k+1}}{1-y}; \quad P_{k,2}(y) = y \frac{1-y^{k+1}}{1-y} \frac{1-y^{k+2}}{1-y} - y^{2(k+1)};$$

$$P_{k,h}(y) = \left( \frac{1-y^{k+1}}{1-y} + y^k \right) P_{k,h-1}(y) - y^k P_{k,h-2}(y) \quad \text{for } h \geq 3.$$

Now, let  $\mathcal{P}_k(x, y)$  be the width (variable  $x$ ) and the height (variable  $y$ ) generating function  $\sum_{h \geq 0} P_{k,h}(y)x^h$ . From the previous recurrence one gets:

**Theorem 5.** *The generating function  $\mathcal{P}_k(x, y)$  is given by:*

$$\mathcal{P}_k(x, y) = \frac{1 - x(y^k + 1 - y^{k+1}) + x^2 y^k (1 - y)}{1 - x \left( \frac{1-y^{k+1}}{1-y} + y^k (1 - x) \right)}.$$

Let  $p_k(n)$  be the number of sand piles in  $\text{IPM}(k)$  of half perimeter (width + height)  $n$ , and  $\mathcal{P}_k(x) = \sum_{n \geq 0} p_k(n)x^n$  be its generating function. As  $\mathcal{P}_k(x) = \mathcal{P}_k(x, x)$ , its expression is:

$$\mathcal{P}_k(x) = \frac{(1-x)^2 (1 - x^{k+1} + x^{k+2})}{1 - 2x - 3x^{k+2} + x^{k+3} + x^{k+1}}.$$

When  $k$  grows, the quantity  $p_k(n)$ , asymptotically equal for large  $n$  to  $c_k / \rho_k^n$  with  $c_k \in \mathbb{R}$  and  $\rho_k$  the smallest root of the denominator  $1 - 2x - 3x^{k+2} + x^{k+3} + x^{k+1}$ , gets closer to  $2^n$ , the number of partitions of semi-perimeter  $n$ .

**2.3. Asymptotics.** Define

$$p_k = [q^n] \prod_{i \geq 1} \frac{1 - q^{ki}}{1 - q^i}, \quad B_k = \sqrt{\frac{k-1}{6k}}, \quad \text{and} \quad C_k = \frac{1}{2} \left( \frac{k-1}{6k^3} \right)^{1/4}.$$

Then  $p_k(n) = C_k n^{-3/4} \exp(B_k n^{1/2}) O(1 + n^{-1/4})$ . If  $I_k(n)$  is the number of partitions of  $n$  in  $\text{IPM}(k)$ , then  $p_{k+1}(n) \leq I_k(n) \leq p_{k+2}(n)$ .

### 3. The Model $L(\theta)$

The model  $L(\theta)$  generalizes the  $\text{SPM}(k)$  by restricting its vertical rule, instead of the horizontal rule as  $\text{IPM}(k)$ . Namely, the difference between the two consecutive columns involved must be greater than  $\theta$ .

**3.1. Definition.** In [4], the sand piles in  $L(\theta)$  are characterized in the following way:

**Proposition 2.** A sand pile in  $L(\theta)$  is a partition  $\pi = (\pi_1, \dots, \pi_l)$  of  $n$  such that

- for  $1 \leq i < l$ ,  $\pi_i - \pi_{i+1} \geq \theta - 1$ ;
- for any  $i < j$  with  $\pi_i - \pi_{i+1} = \theta - 1$  and  $\pi_j - \pi_{j+1} = \theta - 1$ , there exists  $z$  with  $i < z < j$  such that  $\pi_z - \pi_{z+1} > \theta$ .

Let  $\mathcal{L}_\theta(q, x) = 1 + \sum_{\pi \in L(\theta)} x^{l(\pi)} q^{|\pi|}$  the generating function of sand piles in  $L(\theta)$  according to their height and area.

**Lemma 1.**  $\mathcal{L}_\theta(q, x)$  satisfies the  $q$ -equation:

$$\mathcal{L}_\theta(q, x) = \frac{1 - (xq)^{\theta-1}}{1 - xq} + \left( \frac{(xq)^{\theta-1}}{1 - xq} + x^{\theta-1} q^\theta \right) \mathcal{L}_\theta(q, xq).$$

**Theorem 6.**  $\mathcal{L}_\theta(q, x)$  is given by:

$$\mathcal{L}_\theta(q, x) = \sum_{n \geq 0} x^{\theta n} q^{\theta n(n+1)/2} \frac{1 - (xq^{n+1})^\theta}{1 - xq^{n+1}} \prod_{i=1}^n \left( q + \frac{1}{1 - xq^i} \right).$$

Bounds can be obtained for all  $\theta$  for the number  $l_{n,\theta}$  of partitions in  $L(n, \theta)$ .

### 4. Frobenius Model

Another generalization consists in allowing the grains to move both to the left and to the right. In [2], Corteel defines such a model, called the *Frobenius sand pile*, in the following way:

**Definition 1.** Let  $l$  be an integer. A *Frobenius sand pile* is a pair consisting of a pivot indice  $p(a) \leq l$  and a sequence of integers  $(a_1, a_2, \dots, a_l)$  such that

$$a_1 \leq a_2 \leq \dots \leq a_{p(a)} \geq a_{p(a)+1} \geq \dots \geq a_l.$$

**4.1. Order on Frobenius sand piles.**

**Definition 2.** Let  $a = (p(a), (a_1, a_2, \dots, a_l))$  and  $b = (p(b), (b_1, b_2, \dots, b_l))$  be two Frobenius sand piles. Then  $a \geq_F b$  if and only if, for all  $i, j \geq 0$ ,

$$\sum_{l=p(a)-i}^{p(a)+j} a_l \geq \sum_{l=p(b)-i}^{p(b)+j} b_l.$$

**Proposition 3.** Let  $L_F(n)$  be the set of Frobenius partitions ordered by  $\geq_F$ . Then  $L_B(n)$  is a suborder of  $L_F(n)$ .

*Length of a maximal chain.* For  $n \geq 3$ , the length of a maximal chain is greater than  $2n - 4$ , and smaller than  $2\binom{l+1}{3} + lj + 1$ , where  $l$  and  $j$  are defined by  $n = j + l(l + 1)/2$  and  $0 \leq j \leq l$ .

**Definition 3.** Let  $a = (p(a), (a_1, a_2, \dots, a_l))$  be a sand pile.  $a_{<}$ ,  $a_{>}$ ,  $a_{\leq}$ , and  $a_{\geq}$  are defined by

$$\begin{aligned} a_{<} &= (a_{p(a)-1}, a_{p(a)-1}, \dots, a_1); & a_{>} &= (a_{p(a)+1}, a_{p(a)+2}, \dots, a_l); \\ a_{\leq} &= (a_{p(a)}, a_{p(a)-1}, \dots, a_1); & a_{\geq} &= (a_{p(a)}, a_{p(a)+1}, \dots, a_l). \end{aligned}$$

If we constrain horizontal shifts to be smaller than  $k$ , we can create an increasing sequence of orders IFPM( $k$ ) with the relations of order  $\geq_k$ . The Frobenius sand piles of IFPM( $k$ ) are characterized by:

**Proposition 4.** Let  $a = (p(a), (a_1, a_2, \dots, a_l))$  be a sand pile. This sand pile belongs to IFPM( $k$ ) if and only if both of  $a_{<}$  and  $a_{>}$ , and at least one of  $a_{\leq}$  and  $a_{\geq}$  belong to IPM( $k$ ).

**4.2. Generating functions.** The only available generating function is the series of  $F$ -partitions given by

$$1 + \sum_{k \geq 1} q^k \prod_{i=1}^k \frac{1}{(1 - q^i)^2}.$$

For IFPM( $k$ ), we must so far satisfy ourselves with the bound  $F_k(n) \leq |\text{IFPM}(n, k)| \leq F_{k+1}(n)$  for

$$F_k(n) = [q^n] \left( 1 + \sum_{j \geq 1} q^j \prod_{i=1}^j \frac{1 - q^{(k+1)i}}{(1 - q^i)^2} \right).$$

## 5. Conclusion and Open Questions

We have studied different sand pile models related to integer partitions, and in particular we have computed generating functions and asymptotic bounds. A question of interest would consist in getting exact asymptotics instead of asymptotic bounds only. One could start with the area generating function in the SPM case, given by

$$\sum_{n \geq 0} x^n q^{n(n+1)/2} \prod_{i=1}^n \left( q + \frac{1}{1 - q^i} \right).$$

## Bibliography

- [1] Brylawski (Thomas). – The lattice of integer partitions. *Discrete Mathematics*, vol. 6, 1973, pp. 201–219.
- [2] Corteel (Sylvie). – *Problèmes énumératifs issus de l'Informatique, de la Physique et de la Combinatoire*. – Thèse, Université Paris-Sud, 2000.
- [3] Goles (Eric) and Kiwi (Marcos A.). – Games on line graphs and sand piles. *Theoretical Computer Science*, vol. 115, n° 2, 1993, pp. 321–349.
- [4] Phan (Ha Duong). – *Structures ordonnées et piles de sable*. – Thèse, Université Paris 7, 1998.



## On the Group of a Sandpile

*Dominique Rossin*

LIX, École polytechnique (France)

October 16, 2000

*Summary by Dominique Gouyou-Beauchamps*

### Abstract

The abelian sandpile model is a cellular automaton. Its rules generalize the sandpile rules for general graphs. This model has been introduced by Bak, Tang, and Wiesenfeld [1] in 1987. Dhar [9] showed that the set of recurrent configurations of this automaton has the structure of a finite abelian group.

In this talk, we describe several algorithms to determine the identity in the group. This element presents fractal aspects that we are not able yet to explain. These algorithms allow us to introduce relationships between the sandpile group and well-known algebraic or combinatorial objects.

Details may be found in the recent works of R. Cori, D. Rossin, and B. Salvy [6], and D. Rossin [12]. The papers [1, 10], the book [2], and the thesis [13] are good introductions to sandpiles.

### 1. Introduction

Let  $G = (V, E)$  be a non-oriented and connected *multi-graph* with  $V = \{1, \dots, n\}$  its set of vertices and  $E$  a symmetric  $n \times n$  matrix whose entry  $e_{i,j}$  is the number of edges with endpoints  $i, j$ . It is assumed that for any  $i$ ,  $e_{i,i} = 0$  so that the multi-graph has no loops. Frequently,  $G$  is a graph, and hence  $e_{i,j}$  is either 0 or 1. The *degree* of vertex  $i$  in  $G$  is  $d_i = \sum_{j=1}^n e_{i,j}$ . A multi-graph is *rooted* if one of its vertices is distinguished, it is called the *sink* and is numbered  $n$ .

A *configuration*  $u = (u_1, \dots, u_n) \in \mathbb{N}^n$  of  $G$  is a vector of non-negative integers. In the context of the sandpile model, the vertices of the graph are cells, and the number  $u_i$  may be interpreted as the height of a pile of grains of sand standing in cell  $i$ . In the rest of this talk, the number of grains in the sink is not taken into account. Thus two configurations  $u$  and  $v$  which differ only in position  $n$  are considered as equal; we write  $u = v$  if  $u_i = v_i$  for all  $1 \leq i < n$ . This translates the fact that the sink collects all grains of sand getting out of the system.

A *toppling* of the vertex  $i$ ,  $1 \leq i < n$ , in configuration  $u$  consists in decreasing the number of grains in this vertex by its degree while the number of those of each of its neighbours  $j$  increases by  $e_{i,j}$ . This is equivalent to the addition to  $u$  of the vector  $\Delta_i$  such that  $(\Delta_i)_i = -d_i$  and  $(\Delta_i)_j = e_{i,j}$  for  $j \neq i$ . The notation  $u \rightarrow v$  means that  $v$  is obtained from  $u$  by *toppling* a vertex, so that there exists an  $1 \leq i < n$  such that  $v = u + \Delta_i$ . The transitive closure of the toppling operation  $\rightarrow$  is denoted  $\xrightarrow{*}$ :  $u \xrightarrow{*} v$  if  $v$  is obtained from  $u$  by a sequence of topplings. An *avalanche* is a sequence of topplings (see Figures 1, 2 and 3).

The sandpile model has been introduced by Bak, Tang, and Wiesenfeld [1] in 1987. In a recent book, Bak [2] gives an overview of many physical problems—earthquakes and solar flares for

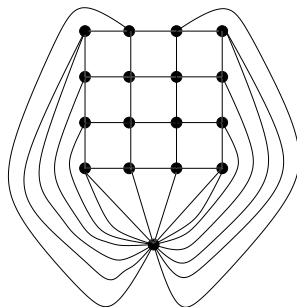


FIGURE 1. Multi-graph corresponding to the  $4 \times 4$  grid.

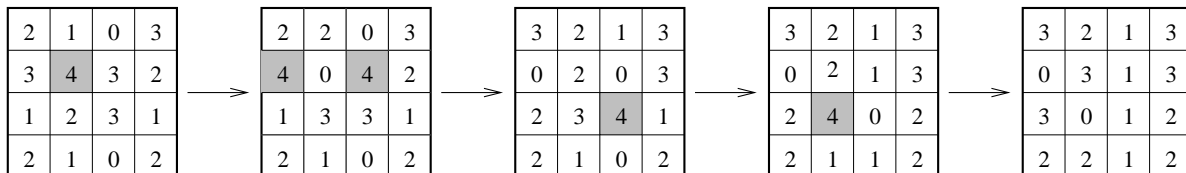


FIGURE 2. Topplings and avalanche on the  $4 \times 4$  grid.

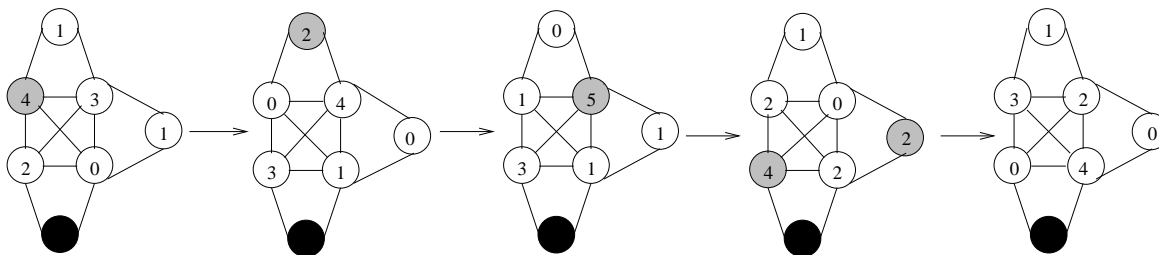


FIGURE 3. Topplings and avalanche on a graph.

example—whose models are based on the sandpile one. All these models follow the Gutenberg–Richter law:  $\log N = a - bM$ ,  $\log E = c + dM$ , and  $N \sim E^{1-\tau}$  ( $\tau \approx 2$ ) where  $M$  is the magnitude,  $N$  is the number of topplings, and  $E$  is the energy. In three dimensions,  $N \sim E^{1-\tau}$  ( $\tau \approx 2.5$ ). A very similar automaton was introduced independently by other authors under the name of the chip-firing game [4, 11]. Biggs [3] found many algebraic and combinatorial properties of the chip-firing game, some of which correspond to Dhar’s results on sandpiles [10]. In [5], we also showed a close relationship between recurrent configurations of the complete graph and the parking functions.

## 2. The Sandpile Group

A vertex is *stable* if it contains a number of grains less than its degree, otherwise this vertex is *unstable*. A *stable configuration* is a configuration where all vertices are stable. It is not difficult to prove that for every configuration  $u$  there exists a stable configuration  $\hat{u}$  such that  $u \xrightarrow{*} \hat{u}$ . Moreover this configuration is unique, and the number of topplings is independent of the way in which  $\hat{u}$  is obtained from  $u$  [9].

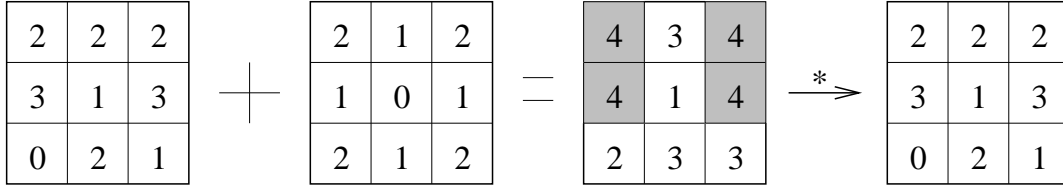


FIGURE 4. A recurrent configuration.

Let  $u, v$  be two configurations. Let  $u_i$  (resp.  $v_i$ ) be the number of grains on vertex  $i$  in configuration  $u$  (resp.  $v$ ). We will denote by  $u+v$  the configuration  $w$  such that  $w_i = u_i + v_i$ . A configuration  $u$  is *recurrent* if it is stable and if there exists a configuration  $v \neq 0$  such that  $u + v \xrightarrow{*} u$  (see Figure 4). The simplest example of a recurrent configuration is  $\delta = (d_1 - 1, d_2 - 1, \dots, d_{n-1} - 1, 0)$ . The set of recurrent configurations is isomorphic to the set of equivalence classes defined by the symmetric closure  $\equiv$  of  $\xrightarrow{*}$ .

Let  $T_G(x, y)$  be the Tutte polynomial of the graph  $G$ . Then  $T_G(1, y)$  is the the generating function (a polynomial) of the recurrent configurations according to the number of sand grains.

We can associate to the set of recurrent configurations the operator  $\oplus$  defined by  $u \oplus v = \widehat{u + v}$  where  $u$  and  $v$  are two recurrent configurations. The set of recurrent configurations with the operation  $\oplus$  is an abelian group  $\mathcal{G}$  [8], this group is equal to the product  $\mathcal{G} = \prod_{i=1}^n \mathbb{Z}/d_i\mathbb{Z}$  and the group structure does not depend on the sink choice in the graph  $G$ .

Let  $u = (u_1, \dots, u_n)$  be a recurrent configuration. We denote  $\bar{u}$  the recurrent configuration  $(d_1 - 1 - u_1, d_2 - 1 - u_2, \dots, d_{n-1} - 1 - u_{n-1}, 0)$ . Then the identity of the sandpile group is  $\text{Id} = \delta \oplus (\bar{\delta} \oplus \delta)$  and the opposite of a recurrent configuration  $u$  is  $\ominus u = \text{Id} \oplus (\bar{\delta} \oplus u)$ .

### 3. Toppling Ideal, Set Topplings and Minimal Gröbner Basis

Configurations and topplings are easily translated from the linear algebra setting into polynomial operations by associating to a configuration  $u = (u_1, u_2, \dots, u_n) \in \mathbb{N}^n$  a monomial  $x_u = x_1^{u_1} x_2^{u_2} \dots x_n^{u_n} \in \mathbb{Q}[x_1, \dots, x_n]$ . To a toppling  $\Delta_i$  is associated the binomial  $T(x_i) = x_i^{d_i} - \prod_j x_j^{e_{i,j}}$ . The addition of two configurations translates into the multiplication of the corresponding monomials and toppling vertex  $i$  in  $u$  translates into the division of  $x_u$  by  $x_i^{d_i}$  followed by the multiplication by  $\prod_{j=1}^n x_j^{e_{i,j}}$ . We define the *toppling ideal*  $\mathcal{I}_G$  as the ideal generated by  $x_n - 1$  and the toppling polynomials  $T(x_i)$  for  $i \in \{1, \dots, n\}$ .

A toppling polynomial can also be associated to a subset  $X$  of the set  $V$  of vertices as follows. For a vertex  $i$  of  $V$ , define

$$d_i(X) = \sum_{j \in X} e_{i,j},$$

the number of edges with endpoints  $i$  and a vertex of  $X$ . The *set toppling* of the set  $X$  in configuration  $u$  consists in adding the vector  $\Delta_X$  to  $u$ , where

$$(\Delta_X)_i = \begin{cases} -d_i(\bar{X}), & \text{for } i \in X, \\ d_i(X), & \text{for } i \in \bar{X}, \end{cases}$$

where  $\bar{X}$  denotes  $V \setminus X$ .

Accordingly, the *toppling polynomial* of the subset  $X$  of  $V$  is defined by

$$T(X) = \prod_{i \in X} x_i^{d_i(\bar{X})} - \prod_{i \in \bar{X}} x_i^{d_i(X)}.$$

Gröbner bases are a classical computational tool for dealing with polynomial ideals. Given an ordering on monomials which is compatible with the product (a so-called *admissible* ordering) and a set of generators of an ideal  $\mathcal{I}$ , one can compute a Gröbner basis for  $\mathcal{I}$  and from there test ideal membership and more generally compute normal forms in the quotient of the algebra by  $\mathcal{I}$ . The rest of this work makes use of the notation and basic results from [7, Chapter 2].

The *graded reverse lexicographic order* (grevlex) denoted  $\prec$ , is defined as follows. If  $A = \prod_{i=1}^n x_i^{\alpha_i}$  and  $B = \prod_{i=1}^n x_i^{\beta_i}$  are two monomials in the variables  $x_i$ ,  $i = 1, 2, \dots, n$ , then  $A \prec B$  if

$$|\alpha| = \sum_{i=1}^n \alpha_i < |\beta| = \sum_{i=1}^n \beta_i$$

or  $|\alpha| = |\beta|$  and in  $(\alpha_1, \dots, \alpha_n) - (\beta_1, \dots, \beta_n)$  the right-most non-zero entry is positive.

From there a *toppling order* is defined as follows: let  $\sigma$  be a permutation of  $\{1, \dots, n\}$  such that  $\sigma(n) = n$  and if the distance from vertex  $i$  to the sink is larger than the distance from vertex  $j$  to the sink, then  $\sigma(i) > \sigma(j)$ . The toppling order is the graded reverse lexicographic order on  $x_{\sigma(1)}, x_{\sigma(2)}, \dots, x_{\sigma(n)}$ .

**Theorem 1.** *A Gröbner basis of the ideal  $\mathcal{I}_G$  with respect to a toppling order is given by*

$$\mathcal{T} = \{T(X) \mid X \subset \{1, \dots, n\}\} \cup \{x_n - 1\}.$$

A Gröbner basis is *minimal* when its elements have leading coefficient 1 and no leading monomial divides another leading monomial in the basis. A subset  $X$  of vertices of the graph  $G = (V, E)$  is well-connected if both subgraphs of  $G$  induced by  $X$  and  $\bar{X}$  are connected.

**Theorem 2.** *The set  $S_c$  of toppling polynomials corresponding to the sets  $X \subset \{1, 2, \dots, n-1\}$  which are well-connected is a minimal Gröbner basis for the toppling order.*

In the worst case, the minimal Gröbner basis still contains  $2^{n-1}$  elements for the complete graph.

As mentioned before, the quotient  $\mathbb{Q}[x_1, x_2, \dots, x_n]/\mathcal{I}_G$  is a  $\mathbb{Q}$ -vector space whose dimension is the order of the group of recurrent configurations. From a Gröbner basis for  $\mathcal{I}_G$ , a basis of this vector space is given by the set of monomials that do not reduce to 0 by the basis. We call these *reduced* monomials. Theorem 3 below gives a simple bijection between reduced monomials for the toppling order and recurrent configurations.

Let  $\Phi$  be the mapping from the set of stable configurations onto itself given by  $\Phi(u) = \delta - u$ . We also denote  $\Phi(M) = \Phi(a_1, a_2, \dots, a_n)$  for a monomial  $M = x_1^{a_1} x_2^{a_2} \dots x_n^{a_n}$ .

**Theorem 3.** *The mapping  $\Phi$  defines a bijection between the set of reduced monomials with respect to the toppling order and the set of recurrent configurations.*

For a configuration  $u$ , let  $\rho(u)$  denote the reduced configuration obtained from the monomial associated to  $u$  by performing reductions by the Gröbner basis of  $\mathcal{I}_G$  associated with the toppling order.

**Proposition 1.** *If  $u$  is a configuration then the recurrent configuration equivalent to  $u$  is*

$$\Phi\left(\rho\left(\Phi(\rho(u))\right)\right).$$

*The identity in the group of recurrent configurations is  $\Phi(\rho(\delta))$ .*

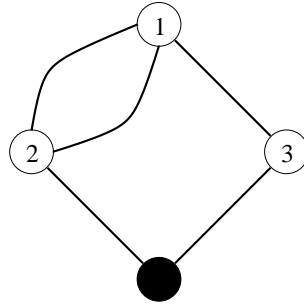


FIGURE 5. Multigraph with 4 vertices.

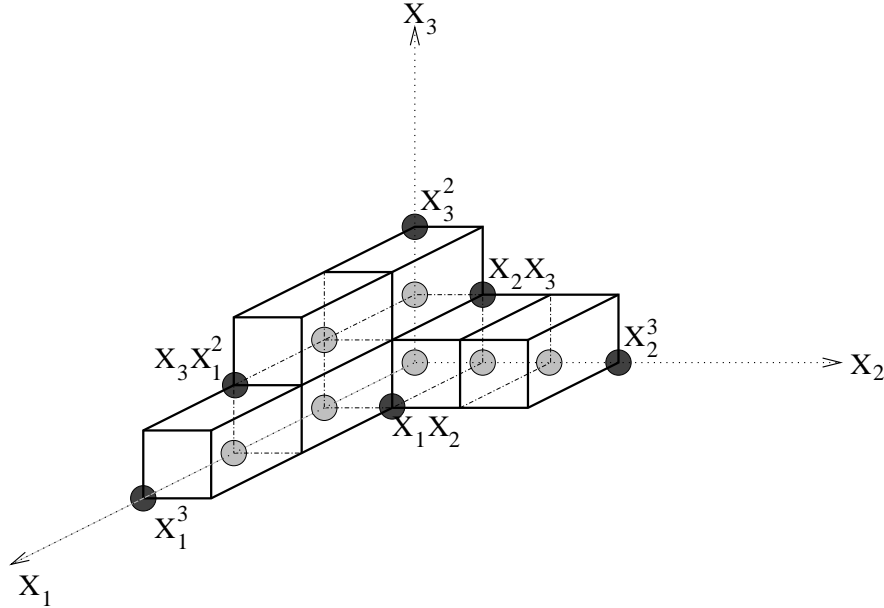


FIGURE 6. Representation of irreducible monomials.

**Corollary 1.** For two recurrent configurations  $u$  and  $v$ ,

$$u \oplus v = \Phi\left(\rho(\Phi(u) + \Phi(v))\right).$$

Proposition 1 yields the following algorithm to compute the identity on a graph  $G$  with sink  $s$ : beginning with the configuration  $\delta$ , perform the set topplings for all well-connected subgraph of  $G \setminus \{s\}$  (this is equivalent to reducing by the Gröbner basis for the toppling order). When no further set toppling can be performed, for each cell  $i$  replace its number of grains  $n_i$  with  $d_i - n_i$ . The resulting configuration is the identity.

#### 4. Examples

Our first example corresponds to the graph displayed on Figure 5. The structure of the graph is reflected by the toppling polynomials for the vertices:

$$x_1^3 - x_2^2 x_3, \quad x_2^3 - x_1^2 x_4, \quad x_3^2 - x_1 x_4, \quad x_4^2 - x_2 x_3, \quad x_4 - 1.$$

The minimal Gröbner basis for the graded reverse lexicographic order on monomials is

$$x_3^2 - x_1, \quad x_2^3 - x_1^2, \quad x_1^3 - x_2, \quad x_2x_3 - 1, \quad x_2x_1 - x_3, \quad x_3x_1^2 - x_2^2, \quad x_4 - 1.$$

Apart from the last, these polynomials correspond respectively to well-connected subgraphs with vertices

$$\{3\}, \quad \{2\}, \quad \{1\}, \quad \{1, 2, 3\}, \quad \{1, 2\}, \quad \{1, 3\}.$$

Given a Gröbner basis  $G = \{p_1, p_2, \dots, p_k\} \subset \mathbb{K}[x_1, x_2, \dots, x_n]$  for some field  $\mathbb{K}$ , it is usual to represent the leading monomials of the  $p_i$  on an integer lattice in  $n$  dimensions. Each polynomial  $p$  is associated to a point  $c(p)$  whose coordinates are the exponents of its leading monomial. The leading terms of the  $p_i$  generate the ideal of leading terms of polynomials in the ideal. These leading terms are thus exactly represented by  $\bigcup c(p_i) + \mathbb{N}^n$ . This removes from  $\mathbb{N}^n$  a staircase shape whose lattice points correspond to the quotient (see Figure 6). Their number is exactly the order of the group of recurrent configurations. Note that in our example, those seven monomials are  $\{1, x_1, x_1^2, x_2, x_2^2, x_3, x_1x_3\}$ , none of which correspond to a recurrent configuration. However, applying  $\Phi$  yields the recurrent configurations as explained above.

Our second example is the  $2 \times 2$  grid consisting of 4 cells, each connected twice to the sink. The sandpile group of this grid, computed for instance in [9], is the product of two cyclic group of orders 24 and 8.

After the computation of the Gröbner basis of the ideal generated by the toppling polynomials of vertices, it follows that  $x_4$  is of order 24 and that any element can be expressed as a product  $x_3^i x_4^j$  where  $0 \leq i \leq 7$  and  $0 \leq j \leq 23$ , which gives that the order of the group is 192. Also, since  $x_1$  and  $x_2$  can be expressed in terms of  $x_3$  and  $x_4$ , it is seen that the group has two generators.

### Bibliography

- [1] Bak (P.), Tang (C.), and Wiensfeld (K.). – An explanation of  $1/f$  noise. *Physical Review Letters*, vol. 59, n° 4, July 1987, pp. 381–384.
- [2] Bak (Per). – *How nature works*. – Copernicus, New York, 1996, xiv+212p. The science of self-organized criticality.
- [3] Biggs (N. L.). – Chip-firing and the critical group of a graph. *Journal of Algebraic Combinatorics*, vol. 9, n° 1, 1999, pp. 25–45.
- [4] Björner (Anders), Lovász (László), and Shor (Peter W.). – Chip-firing games on graphs. *European Journal of Combinatorics*, vol. 12, n° 4, 1991, pp. 283–291.
- [5] Cori (Robert) and Rossin (Dominique). – On the sandpile group of dual graphs. *European Journal of Combinatorics*, vol. 21, n° 4, 2000, pp. 447–459.
- [6] Cori (Robert), Rossin (Dominique), and Salvy (Bruno). – *Polynomial ideals for sandpiles and their Gröbner bases*. – Research Report n° 3946, Institut National de Recherche en Informatique et en Automatique, 2000. 20 pages. Submitted to Elsevier Preprint.
- [7] Cox (David), Little (John), and O’Shea (Donal). – *Ideals, varieties, and algorithms*. – Springer-Verlag, New York, 1997, second edition, xiv+536p. An introduction to computational algebraic geometry and commutative algebra.
- [8] Creutz (M.). – Abelian sandpile. *Computers in Physics*, vol. 5, 1991, pp. 198–203.
- [9] Dhar (D.), Ruelle (P.), Sen (S.), and Verma (D.-N.). – Algebraic aspects of abelian sandpile models. *Journal of Physics A*, vol. 28, n° 4, 1995, pp. 805–831.
- [10] Dhar (Deepak). – Self-organized critical state of sandpile automaton models. *Physical Review Letters*, vol. 64, n° 14, 1990, pp. 1613–1616.
- [11] Goles (Eric). – Sand piles, combinatorial games and cellular automata. In *Instabilities and nonequilibrium structures, III (Valparaíso, 1989)*, pp. 101–121. – Kluwer Acad. Publ., Dordrecht, 1991. Vol. 64 in Mathematics and its Applications.
- [12] Rossin (D.). – On the complexity of addition in the sandpile group of a graph. – Submitted to Elsevier Preprint.
- [13] Rossin (D.). – *Propriétés combinatoires de certaines familles d’automates cellulaires*. – Thèse, École polytechnique, Palaiseau, France, 2000.

## The Tennis Ball Problem

*Donatella Merlini*

DSI, Università degli Studi di Firenze (Italy)

March 19, 2001

*Summary by Cyril Banderier*

### Abstract

Our object is to explore the “ $s$ -tennis ball problem” (at each round  $s$  balls are available and we play with one ball at a time). This is a natural generalization of the case  $s = 2$  considered by Mallows and Shapiro. We show how this generalization is connected with  $s$ -ary trees and employ the notion of generating trees to obtain a solution expressed in terms of generating functions. Then, we present a variation in which at each round we have 4 balls and play with 2 balls at a time. To solve this problem we use the concepts of Riordan arrays and stretched Riordan arrays, and a generalization of generating trees. This is a joint work by D. Merlini with D. G. Rogers, R. Sprugnoli and M. C. Verri.

### 1. Introduction

Let  $1 \leq t < s$  be two integer numbers. A tennis player begins a match with 0 ball in the pocket. At each round, he is given  $s$  new balls, that he puts in the pocket, and throws away  $t$  balls, and so on until the  $n$ th round. The balls are labelled from 1 to  $sn$  and are served in increasing order. The  $tn$  balls thrown away form a sequence of  $tn$  labels. Two sequences which are equal once sorted are considered equivalent. The tennis ball problem consists in evaluating the following two quantities: the number  $f_n$  of nonequivalent configurations after  $n$  rounds and the cumulative sum  $\Sigma_n$  (i.e., the sum—over all the possible configurations—of the labels of the  $tn$  balls that the player threw away).

Turns	Balls received	Balls in the pocket	Balls thrown away
$n = 1$	1 and 2	1 and 2	1
$n = 2$	3 and 4	2, 3, and 4	3
$n = 3$	5 and 6	2, 4, 5, and 6	2
$n = 4$	7 and 8	4, 5, 6, 7, and 8	6
			sum = $1 + 3 + 2 + 6 = 12$

Turns	Balls received	Balls in the pocket	Balls thrown away
$n = 1$	1 and 2	1 and 2	2
$n = 2$	3 and 4	1, 3, and 4	3
$n = 3$	5 and 6	1, 4, 5, and 6	4
$n = 4$	7 and 8	1, 5, 6, 7, and 8	1
			sum = $2 + 3 + 4 + 1 = 10$

FIGURE 1. Two scenarios for the ( $s = 2, t = 1$ )-tennis ball player.

The configuration after 4 rounds is  $(1, 2, 3, 6)$  for the first example and  $(1, 2, 3, 4)$  for the second example. In fact, for the  $(2, 1)$ -case, one has  $f_1 = 2, f_2 = 5, f_3 = 14, f_4 = \dots$  do you guess what? There is indeed 42 different configurations (after 4 rounds), and if one adds all the sums, one gets  $\Sigma_1 = 1 + 2 = 3, \Sigma_2 = (1 + 2) + (1 + 3) + (1 + 4) + (2 + 3) + (2 + 4) = 23, \Sigma_3 = 131, \Sigma_4 = 664, \dots$

In the next section, it is shown how the  $(s, 1)$ -case can be solved in terms of  $s$ -ary trees (by symmetry, this also solves the  $(s, s - 1)$ -case). Then, the last section is dedicated to the  $(4, 2)$ -case, that the authors solved with Riordan arrays and a bilabelled generating tree technique.

Turns	Balls received	Balls in the pocket	Balls thrown away
$n = 1$	1, 2, 3, 4	1, 2, 3, 4	2, 3,
$n = 2$	5, 6, 7, 8	1, 4, 5, 6, 7, 8	1, 7
$n = 3$	9, 10, 11, 12	4, 5, 6, 8, 9, 10, 11, 12	10, 12
$n = 4$	13, 14, 15, 16	4, 5, 6, 8, 9, 11, 13, 14, 15, 16	5, 16
			$2 + 3 + 1 + 7 + 10 + 12 + 5 + 16 = 56$

FIGURE 2. A scenario for the  $(4, 2)$ -tennis ball problem.

This is the only case solved with  $t \neq 1$ . The general  $(s, t)$ -tennis ball problem remains open.

### 2. The $(s, 1)$ -Tennis Ball Problem

Generating trees are a convenient way to reexpress the problem. Consider an infinite rooted tree  $\mathcal{T}$ . The root (labelled 0 and corresponding to level 0) has  $t$  children (labelled  $1, \dots, t$ ). Each path in this tree corresponds to a scenario, thus each node at level  $n$  has a label which corresponds to the ball thrown away at round  $n$ . As we are counting the *sorted* configurations (that is, one does not care for the order of the balls thrown away), we can without loss of generality suppose that the labels increase with the depth.

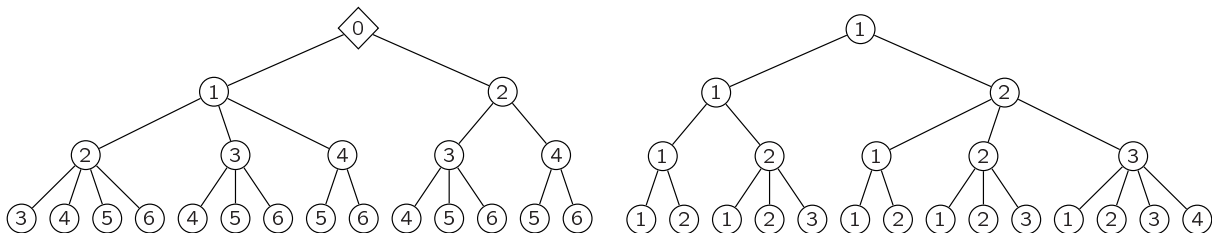


FIGURE 3. The generating tree  $\mathcal{T}$  for the  $(2, 1)$ -case and an isomorphic tree  $\tilde{\mathcal{T}}$ .

More generally, the rewriting rule  $\begin{cases} \text{root} : (1) \\ \text{rule} : (k) \mapsto (1) \dots (k + s - 2) \dots (k + s - 1) \end{cases}$  describes the

formation of a tree  $\tilde{\mathcal{T}}$  which is isomorphic to the generating tree  $\mathcal{T}$  of the  $(s, 1)$ -case: a node with label  $b$  at level  $i$  in the generating tree  $\mathcal{T}$  becomes a node with label  $si - b + 1$  in the tree  $\tilde{\mathcal{T}}$ .

**Theorem 1.** *The number  $f_n$  of configurations for the  $(s, 1)$ -tennis ball problem is the number  $T_{n+1}$  of  $s$ -ary trees with  $n + 1$  nodes. One has*

$$T_n = \frac{\binom{sn}{n}}{1 + (s - 1)n} \quad \text{and} \quad T(z) = 1 + zT(z)^s.$$



*Proof.* The problem can be seen as the enumeration of walks on the integers (with an unbounded set of jumps described by the rewriting rule), for which the generating function can be made explicit [1]. Merlini *et al.* used Riordan array techniques [4].  $\square$

**Theorem 2.** *The cumulative sum (i.e., the sum over all the configurations of the labels of the thrown balls) is*

$$\Sigma_{n-1} = \frac{sn^2 + (s-1)n + 1}{2} \frac{\binom{sn}{n}}{(s-1)n + 1} - \frac{1}{2} \sum_{k=0}^n \binom{sk}{n} \binom{s(n-k)}{n-k}.$$

*Proof.* Consider  $A_n = \sum_{i=0}^n \ell_{i,n}$ , the sum of all the labels (with multiplicity) at level  $i$  in the tree  $\mathcal{T}$ . The cumulative sum  $\Sigma_n$  satisfies

$$\Sigma_n = A_n - \frac{(sn+2)(n+1)}{2} T_{n+1}.$$

The generating function for the sequence  $A_n$  is:

$$A(z) = \frac{s(s-1)zT'(z)^2}{2T(z)} + T'(z).$$

From these two equations, one gets the almost closed form of the theorem.

Note that the asymptotics of  $\Sigma_n$  can easily be deduced from the asymptotics of  $A_n$ .  $\square$

These theorems are consistent with the fact that the (2, 1)-case leads to Catalan numbers  $f_n = \frac{\binom{2n}{n}}{n+1}$  (proven in [2]) and to  $\Sigma_n = \frac{2n^2+5n+4}{n+2} \binom{2n+1}{n} - 2^{2n+1}$  (as it was found in [3] by hand manipulations of sums of binomial coefficients).

### 3. The (4, 2)-Tennis Ball Problem

Here again, as one does not care for the order (of the balls thrown away), one can without loss of generality suppose that any configuration is represented by the smallest equivalent sequence with respect to lexicographical order. Thus the configuration (1, 4), (5, 8), (2, 10) is considered to be the same as the configuration (1, 2), (4, 5), (8, 10).

Let  $M_m^{[n]}$  be the number of pairs at level  $n$  (in the bilabelled generating tree of the (4, 2)-case) with larger element equal to  $m$ ; one has the recurrence

$$M_m^{[n+1]} = \sum_{r=2n}^{m-2} (m-r-1)M_r^{[n]}.$$

Defining  $f_{n,k} = M_{4n+1-k}^{[n]}$  gives an infinite lower-triangular array:

n/k	1	2	3	4	5	6	7	8	9
0	1								
1	3	2	1						
2	22	16	10	4	1				
3	211	158	105	52	21	6	1		
4	2306	1752	1198	644	301	116	36	8	1

One has the relation  $f_{n+1,k+2} = \sum_{j=0}^{\infty} (j+1)f_{n,k+j}$ . The sums  $f_n = \sum_{k \geq 1} f_{n,k}$  give the sequence (1, 6, 53, 554, ...), the number of configurations for the (4, 2)-tennis ball problem.

It is convenient to transform the above array into a proper Riordan array. A proper Riordan array is an infinite lower triangular array  $(D_{n,k})_{n,k \in \mathbb{N}}$  which satisfies

$$d_{n+1,k+1} = \sum_{j=0}^{\infty} a_j d_{n,k+j} \quad \text{for all } n \text{ and } k \text{ in } \mathbb{N}.$$

The generating function  $A(z) = \sum_j a_j z^j$  allows to express  $d_{n,k}$  by a Lagrangean-like formula

$$d_{n,k} = [z^n]g(z)(zh(z))^k \quad \text{where } h(z) = A(zh(z)).$$

The above array can be embedded in the array

n/k	0	1	2	3	4	5	6	7
0	1							
1	0	1						
2	1	1	1					
3	0	<b>3</b>	<b>2</b>	<b>1</b>				
4	6	6	6	3	1			
5	0	<b>22</b>	<b>16</b>	<b>10</b>	<b>4</b>	<b>1</b>		
6	53	53	53	31	15	5	1	

which satisfies  $A(z) = \frac{1}{1-z}$ ,  $h(z) = C(z)$  (the generating function of Catalan numbers), and  $g(z) = \frac{2}{2-zC(z)+zC(-z)}$ . In particular, the function  $g(z)$  generates the first column of this array and corresponds to the number of nonequivalent configurations one wants to enumerate:  $f_n = g_n = \frac{3(n+2)}{(n+3)((2n+3)} \binom{2n+4}{n+2} - \frac{4^{(n+1)/2}}{n+2} \binom{n+3}{(n+3)/2}$  (for even  $n$ ). The cumulative sum in the tree with root  $(0, 0)$  and rewriting rule  $(k_1, k_2) \mapsto (0, 0)(0, 1) \dots (k_1 + 2, k_1 + 2)$  is then given by

$$\Sigma_n = \sum_{h=0}^n \sum_r \mu_r^{[2n-2h]} w_r^{[2h]} = \sum_r \sum_{h=0}^n \mu_r^{[2n-2h]} w_r^{[2h]}.$$

Here,  $\mu_r^{[2n-2h]}$  is the number of nodes at level  $n - h$  in the subtree starting with  $(r, *)$ , and  $w_r^{[2h]}$  the total weight that the couples  $(r, *)$  have at level  $h$ . (Note that a label  $(k_1, k_2)$  at level  $n$  in the new tree corresponds to a label  $(4n - k_2 - 1, 4n - k_1)$  in the generating tree of the  $(4, 2)$ -case.) The Riordan array property yields  $\mu_r(z) = g(z)C(z)^{r+2}$  and  $w_r(z) = g(z)z^r C(z)^{r+1}(zC(z)^2 + 2r)$ , thus

$$\Sigma(z) = \frac{1}{4} \sum_r (\mu_r(z) + \mu_r(-z))(w_r(z) + w_r(-z)) = 12z^2 + 284z^4 + 5436z^6 + 96768z^8 + O(z^{10}).$$

The next nontrivial open cases are the  $(5, 2)$ - and  $(5, 3)$ -tennis ball problems. This is related to the enumeration of 2- and 3-dimensional constrained discrete random walks for which no closed form (or even recurrence) is known. Articles and slides related to this summary can be found at Donatella Merlini's web page <http://www.dsi.unifi.it/~merlini/Publications.html>.

### Bibliography

- [1] Banderier (C.), Bousquet-Mélou (M.), Denise (A.), Flajolet (P.), Gardy (D.), and Gouyou-Beauchamps (D.). – Generating functions for generating trees. *Discrete Mathematics*. – 25 pages. To appear.
- [2] Grimaldi (Ralph P.) and Moser (Joseph G.). – The Catalan numbers and a tennis ball problem. In *Proceedings of the Twenty-eighth Southeastern International Conference on Combinatorics, Graph Theory and Computing (Boca Raton, FL, 1997)*, vol. 125, pp. 65–71. – 1997.
- [3] Mallows (Colin L.) and Shapiro (Lou). – Balls on the lawn. *Journal of Integer Sequences*, vol. 2, 1999. – Article 99.1.5.
- [4] Merlini (D.), Rogers (D. G.), Sprugnoli (R.), and Verri (M. C.). – The tennis ball. – 2001. Submitted.

# Hyperharmonic Numbers and the Phratry of the Coupon Collector

*Dominique Foata*

Département de mathématique, Université Louis Pasteur (France)

May 21, 2001

*Summary by Cyril Banderier*

## Abstract

The classical coupon-collector problem is here extended to the case where the collector shares his harvest with other members of his phratry. She (!) remains the single buyer, but she gives to his brothers all the pictures that she got in double. When her album is filled, her brothers' albums have some empty places. How many in average? Dominique Foata (in a joint work with Guoniu Han and Bodo Lass) answers this question via an expression for the multivariate generating function. The problem is related to hyperharmonic numbers, that are studied here as solutions of finite differences equations.

## 1. Coupons Collector

A clever firm sells chocolate, with a picture (or “coupon”) of a famous cricket player in each bar. In total, there are  $m$  different pictures to collect and each picture appears with probability  $1/m$ . Mr. and Mrs. Brown have  $r$  sons and one daughter, chocolate and cricket addicts. The girl (she's the oldest) is the only one to buy chocolate. She tries to complete her collection. When she gets a new picture, she puts it in her album, and when she gets a double, she gives it to her oldest brother, and when this one gets a double, he gives it to the remaining oldest brother, and so on. After having bought  $T$  bars of chocolate, the girl has completed her album, and it remains  $M_T^{(i)}$  empty places in the album of the  $i$ th brother ( $i = 1, \dots, r$ ). Let  $X_n^{(k)}$  be the number of coupons which appeared exactly  $k$  times until time  $n$ . As  $M_T^{(k)} = X_T^{(1)} + \dots + X_T^{(k)}$ , the distribution of  $(T, M_T^{(1)}, \dots, M_T^{(r)})$  is then totally determined by the distribution of  $(T, X_T^{(1)}, \dots, X_T^{(r)})$ .

The question is to find formulae and asymptotics for  $M_T^{(i)}$ , or equivalently for  $X_T^{(k)}$ .

There is a lot of ways to solve coupon-collector-like problems. One can distinguish three main approaches:

- formal approach: combinatoricians indeed used a language-theory approach (shuffle products and Laplace transforms [2] or manipulation of regular expressions [6]);
- probabilistic approach: a lot of folklore results are established via basic probabilistic considerations, and more sophisticated tools such as martingales theory were also useful [5];
- matricial approach: this is in fact a mixture of the two precedent approaches, which is exploited to solve a Markovian generalization of the coupon-collector problem in [1] (with Perron–Frobenius theory and approximation of integrals).

This is with a combinatorial approach (enumeration of surjections and formal Laplace transform) that Foata *et al.* obtain in [4] the multivariate generating function of the coupon-collector problem, from which they derive the formulae for the expectations of  $\mathbf{E}[T]$  and  $\mathbf{E}[X_T^{(k)}]$ .

### 2. The Multivariate Generating Function

**Theorem 1.** *The generating function of the random variables  $T, X_T^{(1)}, \dots, X_T^{(r)}$  is*

$$\sum_{l \geq m, \mathbf{n}} \mathbf{P}\{T = l, X_T^{(1)} = n_1, \dots, X_T^{(r)} = n_r\} t^l u_1^{n_1} \dots u_r^{n_r} =$$

$$u_1 t \sum \binom{m-1}{a, b, c_1, \dots, c_r} (-1)^b \left( \prod_{k=1}^r \left( \frac{u_k - 1}{k!} \right)^{c_k} \right) \frac{(t/m)^{\sum_k k c_k}}{(1 - at/m)^{1 + \sum_k k c_k}} \left( \sum_k k c_k \right)!$$

The proof follows from setting  $u_i = 1$  (for  $i \geq r + 1$ ), expanding (with Newton multinomial formula), and applying the Laplace transform to

$$\sum_{l \geq m} \frac{t^{l-1}}{(l-1)!} \sum_{\mathbf{n}} \sum_{s \in S(l, m; \mathbf{n})} \pi(s) = m u_1 \left( \sum_{i \geq 1} u_i \frac{t^i}{i!} \right)^{m-1}.$$

The formal Laplace transform is defined as a linear map such that  $\mathcal{L}(g_n \frac{t^n}{n!}) = g_n t^n$ . This implies  $\mathcal{L}(\exp(at)t^n) = n! \frac{t^n}{(1-at)^{n+1}}$ . The set  $S(l, m; \mathbf{n})$  is defined as a subset of the surjections from  $[1, \dots, l]$  to  $[1, \dots, m]$  for which  $s \in S(l, m; \mathbf{n})$  implies  $i$  is reached  $n_i$  times and the restriction of  $s$  to  $[1, \dots, l-1]$  is still a surjection from  $[1, \dots, l-1]$  to  $[1, \dots, m] \setminus \{s(l)\}$ . The weight  $\pi$  of a surjection  $S \in S(l, m; \mathbf{n})$  is defined by  $\pi(s) = \prod u_i^{n_i}$ .

### 3. Hyperharmonic Numbers

In order to complete one collection, it is well known that the average number of needed bars is

$$\mathbf{E}[T] = m H_m \quad \text{where} \quad H_m = \sum_{k=1}^m \frac{1}{k}.$$

For example, when there are  $m = 50$  different pictures,  $\mathbf{E}[T] = 50 H_{50} \approx 50 \times 4.5 \approx 225$  and thus the daughter has a lot of doubles and we can expect that the oldest brother has almost completed his album with the 175 remaining pictures.

Pintacuda [5] proved with martingale theory that  $\mathbf{E}[M_T^{(1)}] = H_m$ . Foata *et al.* prove

**Theorem 2.** *For  $k \geq 2$ , the average number of empty places in the  $k$ th brother's album is*

$$\mathbf{E}[M_T^{(k)}] = 1 + \sum_{i=1}^k K_m^{(i)} \quad \text{where} \quad K_m^{(k)} = \sum_{i=2}^m \frac{K_i^{(k-1)}}{i}, \quad (k \geq 1, m \geq 3)$$

with the following initial conditions  $K_2^{(k)} = \frac{1}{2^k}$  (for  $k \geq 0$ ) and  $K_m^{(0)} = 1$  (for  $m \geq 2$ ).

A first derivation of this result follows of Theorem 1. Another proof is in two steps: first get the generating function for the  $K_m^{(k)}$  (end of this section) and then prove that this generating function is also the one of the coupon-collector problem (next section).

Consider the rising factorial defined by  $(a)_n = a(a+1) \dots (a+n-1)$  if  $n \geq 1$  and  $(a)_0 = 1$ . An hypergeometric function with respect to two lists  $(a_1, \dots, a_r)$  and  $(b_1, \dots, b_s)$  is defined as the function given by the series

$${}_r F_s \left( \begin{matrix} a_1, \dots, a_r \\ b_1, \dots, b_s \end{matrix}; x \right) := \sum_{n \geq 0} \frac{(a_1)_n \dots (a_r)_n x^n}{(b_1)_n \dots (b_s)_n n!}.$$

The authors prove that the numbers  $K_m^{(k)}$  (that they call “hyperharmonic numbers”) satisfy

$$K_m^{(k)} = \frac{m(m-1)}{2^{k+1}} {}_{k+2}F_{k+1} \left( \begin{matrix} -m+2, 2, \dots, 2 \\ 3, \dots, 3 \end{matrix} \middle| 1 \right).$$

Comparing the recurrences satisfied by both sides and then summing gives the generating function

$$(1) \quad \sum_{k \geq 0} K_m^{(k)} t^k = \sum_{n=2}^m \frac{(-m)_n}{(n-2)!} \frac{1}{n} \frac{1}{1-t/n} = \frac{1}{(1-t/2)(1-t/3)\dots(1-t/m)}$$

(the last equality following from a partial fraction decomposition).

Thus  $K_m^{(k)} = h_k(\frac{1}{2}, \dots, \frac{1}{k})$ , the symmetric homogeneous polynomial of (total) degree  $k$  in  $m-1$  variables. Reexpressing  $h_k$  in the basis of the power symmetric functions  $p_k := \sum x_i^k$  gives

$$K_m^{(k)} \sim \frac{p_1^k}{k!} \sim \frac{(\ln m)^k}{k!}$$

One also has explicit asymptotics (for fixed  $k$ ), e.g.,

$$K_m^{(3)} \approx 1.1666 \ln^3 m - 0.2113 \ln^2 m + 0.4118 \ln m - 0.0815.$$

#### 4. Martingales Rescue the Phratry

Let  $X_n^{(0)}$  be the number of empty places in the daughter’s album. Now, define the process  $X$  as  $X_n = (X_n^{(0)}, X_n^{(1)}, \dots, X_n^{(r)})$ . For any function  $f$ , the average increase of  $f(X)$  (knowing all the previously drawn coupons) is easy to get:

$$\mathbf{E}[f(X_{n+1}) - f(X_n) \mid Y_0, \dots, Y_n] = \sum_{k=0}^r \frac{X_n^{(k)}}{m} (f(X_n^{(0)}, \dots, X_n^{(k)} - 1, X_n^{(k+1)} + 1, \dots, X_n^{(r)}) - f(X_n));$$

this simply reflects the different possible updates ( $X_n^{(k)}/m$  is the probability to get a new coupon which was already in  $k$ -tuple).

If  $f$  is such that the sum is 0, one has also  $W_{n+1} - W_n = 0$  and thus  $W$  is a martingale, where  $W$  is the process  $f(X)$  stopped at  $T$ , that is  $W_n := f(X_n)$  (for  $n < T$ ) and  $W_n := f(X_T)$  (for  $n \geq T$ ).

More generally, suppose that for  $r$  functions  $f^{(1)}, \dots, f^{(r)}$  from  $\mathbb{N}^{k+1}$  to  $\mathbb{R}$  one has:

1.  $\sum_{i=0}^k x_i (f(x_0, \dots, x_i - 1, x_{i+1} + 1, \dots, x_k) - f^{(k)}(x_0, \dots, x_k)) = 0$  for  $x_0 \geq 1$ ;
2.  $f^{(k)}(0, x_1, \dots, x_k) = x_k$ .

Then  $\mathbf{E}[X_T^{(k)}] = f^{(k)}(m, 0, \dots, 0)$ .

*Proof.* 2. implies that  $X_T^{(k)} = f(X_T) = W_T$ ; 1. gives a martingale property for  $W$ , Doob’s theorem for stopping time of martingales gives  $\mathbf{E}[W_T] = \mathbf{E}[W_0] = W_0$ , and 2. implies that  $W_0 = f^{(k)}(m, 0, \dots, 0)$ .  $\square$

Pintacuda [5] used this result with  $k = 1$  and found  $f^{(1)}(x_0, x_1) = H_{x_0} + \frac{x_1}{1+x_0}$ . Foata *et al.* guessed the general formula:

**Proposition 1.** For  $k \geq 2$ , the function  $f^{(k)}$  defined as

$$f^{(k)}(x_0, x_1, \dots, x_k) := K_{x_0}^{(k)} + \frac{x_1 K_{x_0+1}^{(k-1)} + \dots + x_{k-1} K_{x_0+1}^{(1)} + x_k}{x_0 + 1}$$

is the only solution of 1. and 2. One also has  $f^{(k)}(x_0, 0, \dots, 0) = K_{x_0}^{(k)}$ .

They give two proofs in [4], but I prefer to explain what I heard in the meeting Random Structure and Algorithms (Poznan, August 2001), where Doron Zeilberger explained how to use a language theory argument to get a shorter proof. A complete collection of coupons can be written  $11^*2\{1, 2\}^*3\{1, 2, 3\}^*4 \dots \{1, 2, \dots, m-1\}^*m$ . Let  $\mathcal{W}$  be this set of words. This leads to the generating function

$$f(x_1, \dots, x_m) := \frac{x_1}{m} \frac{1}{1 - \frac{x_1}{m}} \cdots \frac{x_2}{m} \frac{1}{1 - \frac{x_1+x_2}{m}} \frac{x_{m-1}}{m} \frac{1}{1 - \frac{x_1+\dots+x_{m-1}}{m}} x_m.$$

Recall that  $\mathbf{E}[X_T^{(k)}]$  is the expected number of kinds of coupons in  $k$ -tuple (at time  $T$ , that is when the daughter has completed her album). Thus,

$$\begin{aligned} \sum_{k=1}^{\infty} \mathbf{E}[X_T^{(k)}] t^k &= \sum_{w \in \mathcal{W}} \sum_{k=1}^m \mathbf{P}(w) t^{|w|_k} = \sum_{k=1}^m \sum_{w \in \mathcal{W}} \left(\frac{1}{m}\right)^{|w|} t^{|w|_k} \\ &= m! (f(t, 0, \dots, 0) + f(0, t, 0, \dots, 0) + \cdots + f(0, \dots, 0, t)) \\ &= t + t \sum_{k=1}^{m-1} \frac{k!}{(2-t)(3-t) \dots (k+1-t)} = t - 1 + \frac{m!}{\prod_{j=2}^m (j-t)}. \end{aligned}$$

As words of  $\mathcal{W}$  are ordered (whereas it is in fact irrelevant for the coupon collector), there is a factor  $m!$  at the second line takes into account all the permutations. The generating function obtained at the last line shows that the hyperharmonic numbers generated by Equation (1) indeed gives the average value of Theorem 2.

## 5. Conclusion

The coupon-collector problem (like the ménage problem, the birthday paradox) belongs to the large class of problems that can be modeled by simple urns models. It is very likely that, during the next years, the symbolic method will be applied with success to all these urns problems, and analytic combinatorics will then provide enumeration, complete asymptotics expansions and limit laws. The “classical” coupon-collector problem waits for his next revisitor!

This summary is related to Foata’s article [4] (the more recent preprint [3] is also relevant). These articles are accessible at <http://www-irma.u-strasbg.fr/~foata/paper>.

## Bibliography

- [1] Banderier (Cyril) and Dobrow (Robert P.). – A generalized cover time for random walks on graphs. In *Formal power series and algebraic combinatorics (Moscow, 2000)*, pp. 113–124. – Springer, Berlin, 2000.
- [2] Flajolet (Philippe), Gardy (Danièle), and Thimonier (Loÿs). – Birthday paradox, coupon collectors, caching algorithms and self-organizing search. *Discrete Applied Mathematics*, vol. 39, n° 3, 1992, pp. 207–229.
- [3] Foata (D.) and Zeilberger (D.). – The collector’s brotherhood problem using the Newman-Shepp symbolic method. – To appear in *Algebra Universalis*. Special issue dedicated to the memory of Gian-Carlo Rota.
- [4] Foata (Dominique), Han (Guo-Niu), and Lass (Bodo). – Les nombres hyperharmoniques et la fratrie du collectionneur de vignettes. *Séminaire Lotharingien de Combinatoire*, vol. 47, n° B47a, 2001. – 20 pages. Available from <http://www.mat.univie.ac.at/~slc/>.
- [5] Pintacuda (N.). – Coupons collectors via the martingales. *Bollettino dell’Unione Matematica Italiana. A. Serie V*, vol. 17, n° 1, 1980, pp. 174–177.
- [6] Zeilberger (Doron). – How many singles, doubles, triples, etc., should the coupon collector expect? *Personal Journal of Ekhad and Zeilberger*, 2001. – 1 page. Available from <http://www.math.rutgers.edu/~zeilberg/>.

## Mac Mahon’s Partition Analysis Revisited

*Peter Paule*

RISC, Linz (Austria)

October 2, 2000

*Summary by Sylvie Corteel*

### Abstract

The purpose of this talk is to present the  $\Omega$  operator introduced by Mac Mahon in 1915 and to show its power in current combinatorial and partition-theoretic research. This operator is implemented in the Mathematica Package Omega which was developed by A. Riese. This is joint work with G. E. Andrews (Penn State University) and A. Riese (RISC-Linz).

### 1. Introduction

Mac Mahon devoted many pages of his famous book “Combinatorial Analysis” [9] to  $\Omega$ -calculus. Nevertheless this method was not used for 85 years except by Stanley in 1973 [10]. The purpose of this talk is to present the  $\Omega$  operator and to show its power in current combinatorial and partition-theoretic research [1, 2, 3, 4, 5]. In this summary, we define the  $\Omega$  operator and exhibit a few of its elimination rules, before giving two problems where this operator is a powerful tool: lecture hall partitions and  $k$ -gons of integer length.

### 2. The Omega Operator

Let us now define the operator and present a few rules.

**Definition 1.** [9] The Omega operator  $\underset{\geq}{\Omega}$  is defined as follows:

$$\underset{\geq}{\Omega} \sum_{s_1=-\infty}^{\infty} \cdots \sum_{s_r=-\infty}^{\infty} A_{s_1, \dots, s_r} \lambda_1^{s_1} \cdots \lambda_r^{s_r} = \sum_{s_1=0}^{\infty} \cdots \sum_{s_r=0}^{\infty} A_{s_1, \dots, s_r}.$$

To evaluate this operator, Mac Mahon proposed a list of elimination rules. The proof of each is straightforward as it uses the simple identity

$$\sum_{n \geq 0} x^n = 1/(1-x).$$

We list a few of them only:

$$\underset{\geq}{\Omega} \frac{\lambda^{-s}}{(1-\lambda x) \left(1 - \frac{y}{\lambda}\right)} = \frac{x^s}{(1-x)(1-xy)}, \quad s \geq 0,$$
$$\underset{\geq}{\Omega} \frac{1}{(1-\lambda x) \left(1 - \frac{y}{\lambda}\right) \left(1 - \frac{z}{\lambda}\right)} = \frac{1}{(1-x)(1-xy)(1-xz)},$$

$$\begin{aligned}\Omega_{\geq} \frac{1}{(1-\lambda x)(1-\frac{y}{\lambda^s})} &= \frac{1}{(1-x)(1-x^s y)}, \quad s > 0, \\ \Omega_{\geq} \frac{1}{(1-\lambda^s x)(1-\frac{y}{\lambda})} &= \frac{1+xy\frac{1-y^{s-1}}{1-y}}{(1-x)(1-xy^s)}, \quad s > 0.\end{aligned}$$

For example to find the generating function of the partitions with three parts and whose parts differ by at least two, we use the first rule:

$$\begin{aligned}f_3(q) &= \Omega_{\geq} \sum_{a_1, a_2, a_3 \geq 1} \lambda_1^{a_1 - a_2 - 2} \lambda_2^{a_2 - a_3 - 2} q^{a_1 + a_2 + a_3} = \Omega_{\geq} \frac{\lambda_1^{-2} \lambda_2^{-2} q^3}{(1-\lambda_1 q) \left(1 - \frac{\lambda_2 q}{\lambda_1}\right) \left(1 - \frac{q}{\lambda_2}\right)} \\ &= \Omega_{\geq} \frac{q^2 \lambda_2^{-2} q^3}{(1-q)(1-\lambda_2 q^2) \left(1 - \frac{q}{\lambda_2}\right)} = \frac{q^2 q^4 q^3}{(1-q)(1-q^2)(1-q^3)}.\end{aligned}$$

It is also possible to generalize this result for partitions with  $k$  parts and whose parts differ by at least two for any  $k > 0$ , that is

$$f_k(q) = \frac{q^{k^2}}{(1-q)(1-q^2)\dots(1-q^k)}.$$

### 3. Lecture Hall Partitions

The lecture hall partition theorem is one of the most elegant recent result in partition analysis [6, 7]. Let us state the refinement of this theorem [8].

**Theorem 1.** *The number of partitions of  $n$  of the form  $(b_j, b_{j-1}, \dots, b_1)$  with  $\frac{b_j}{j} \geq \frac{b_{j-1}}{j-1} \geq \dots \geq \frac{b_1}{1} \geq 0$  and  $b_j - b_{j-1} + \dots + (-1)^{j-1} b_1 = m$  is equal to the number of partitions of  $n$  into  $m$  odd parts less than  $2j$ .*

This theorem can also be proved with the Omega operator [1], which is what motivated G. E. Andrews to resuscitate the Omega operator. The proof mainly uses the elimination rule

$$\Omega_{\geq} \frac{1}{(1-\lambda x)(1-\frac{y}{\lambda^s})} = \frac{1}{(1-x)(1-x^s y)}$$

Let us illustrate it for  $j = 3$ .

$$\begin{aligned}\sum_{\frac{b_3}{3} \geq \frac{b_2}{2} \geq \frac{b_1}{1} \geq 0} x^{b_3 - b_2 + b_1} q^{b_3 + b_2 + b_1} &= \Omega_{\geq} \sum_{b_3, b_2, b_1 \geq 0} \lambda_1^{2b_3 - 3b_2} \lambda_2^{b_2 - 2b_1} x^{b_3 - b_2 + b_1} q^{b_3 + b_2 + b_1} \\ &= \Omega_{\geq} \frac{1}{(1-\lambda_1^2 q x) \left(1 - \frac{\lambda_2 q}{\lambda_1^3 x}\right) \left(1 - \frac{qx}{\lambda_2^2}\right)} = \Omega_{\geq} \frac{1}{(1-xq)(1-xq^3)(1-xq^5)}.\end{aligned}$$

The Omega operator can also give a bijective proof of the theorem [5]. Let us show how to proceed for  $j = 3$ :

$$\sum_{\frac{b_3}{3} \geq \frac{b_2}{2} \geq \frac{b_1}{1} \geq 0} q_3^{b_3} q_2^{b_2} q_1^{b_1} = \Omega_{\geq} \sum_{b_3, b_2, b_1 \geq 0} \lambda_1^{2b_3 - 3b_2} \lambda_2^{b_2 - 2b_1} q_3^{b_3} q_2^{b_2} q_1^{b_1} = \Omega_{\geq} \frac{1 + q_2 q_3^2}{(1-q_3)(1-q_2^2 q_3^3)(1-q_1 q_2^2 q_3^3)}.$$



From the previous equation we get that there is a bijection between the lecture hall partitions  $(b_3, b_2, b_1)$  of  $n$  and the partitions of  $n$  into parts  $\{1, 3, 5\}$  with multiplicity  $m_i$  for the part  $i$ . This bijection becomes:

$$b_3 = 3m_5 + 2m_3 - \left\lfloor \frac{m_3}{2} \right\rfloor + m_1, \quad b_2 = 2m_5 + m_3, \quad b_1 = \left\lfloor \frac{m_3}{2} \right\rfloor.$$

#### 4. $k$ -Gons with Integer Length

The problem can be defined as follows. The number  $|T_k(n)|$  of  $k$ -gons with length  $n$  is equal to the number of solutions of

$$(1) \quad a_k \geq a_{k-1} \geq \dots \geq a_1 \geq 1, \quad a_1 + a_2 + \dots + a_k = n, \quad a_1 + a_2 + \dots + a_{k-1} > a_k.$$

Let  $F_k(q) = \sum_n |T_k(n)| q^n$  be the associated generating function. For triangles ( $k = 3$ ) we get

$$F_3(q) = \sum_n |T_3(n)| q^n = \frac{q^3}{(1 - q^2)(1 - q^3)(1 - q^4)}.$$

This is easy to prove as conditions (1) give

$$\begin{aligned} F_3(q) &= \sum_{\substack{\geq \\ a_1 \geq 1 \\ a_2, a_3 \geq 0}} \lambda_1^{a_3 - a_2} \lambda_2^{a_2 - a_1} \lambda_3^{a_1 + a_2 - a_3 - 1} q^{a_1 + a_2 + a_3} \\ &= \sum_{\geq} \frac{q \lambda_1^{-1}}{\left(1 - \frac{q \lambda_2}{\lambda_3}\right) \left(1 - \frac{q \lambda_1 \lambda_3}{\lambda_2}\right) \left(1 - \frac{q \lambda_3}{\lambda_1}\right)} = \frac{q^3}{(1 - q^2)(1 - q^3)(1 - q^4)} \end{aligned}$$

We can even be more specific

$$\begin{aligned} F_3(q_1, q_2, q_3) &= \sum_{\substack{a_3 \geq a_2 \geq a_1 \geq 1 \\ a_1 + a_2 > a_3}} q_1^{a_1} q_2^{a_2} q_3^{a_3} = \sum_{\substack{\geq \\ a_1 \geq 1 \\ a_2, a_3 \geq 0}} \lambda_1^{a_3 - a_2} \lambda_2^{a_2 - a_1} \lambda_3^{a_1 + a_2 - a_3 - 1} q_1^{a_1} q_2^{a_2} \\ &= \sum_{\geq} \frac{q \lambda_1^{-1}}{\left(1 - \frac{q \lambda_2}{\lambda_3}\right) \left(1 - \frac{q \lambda_1 \lambda_3}{\lambda_2}\right) \left(1 - \frac{q \lambda_3}{\lambda_1}\right)} = \frac{q_1 q_2 q_3}{(1 - q_2 q_3)(1 - q_1 q_2 q_3)(1 - q_1 q_2 q_3^2)}. \end{aligned}$$

This shows there is a bijection between the 3-tuples  $(u_1, u_2, u_3)$  of  $\mathbb{N}^3$  and the triangles whose sides have length  $u_1 + u_2 + 1$ ,  $u_1 + u_2 + u_3 + 1$  and  $u_1 + 2u_2 + u_3 + 1$ .

Thanks to the Omega operator we can compute the generating function for larger  $k$ :

$$\begin{aligned} F_4(q) &= \frac{q^4(1 + q + q^5)}{(1 - q^2)(1 - q^3)(1 - q^4)(1 - q^6)}, \\ F_5(q) &= \frac{q^5(1 - q^{11})}{(1 - q)(1 - q^2)(1 - q^4)(1 - q^5)(1 - q^6)(1 - q^8)}, \\ F_6(q) &= \frac{q^6(1 - q^4 + q^5 + q^7 - q^8 - q^{13})}{(1 - q)(1 - q^2)(1 - q^4)(1 - q^6)(1 - q^8)(1 - q^{10})}. \end{aligned}$$

We then can see that no pattern can be found and the Omega operator was a quick tool to show that the solutions of this  $k$ -gon problem do not have “nice” generating functions.

**Bibliography**

- [1] Andrews (George E.). – MacMahon's partition analysis. I. The lecture hall partition theorem. In *Mathematical essays in honor of Gian-Carlo Rota (Cambridge, MA, 1996)*, pp. 1–22. – Birkhäuser Boston, Boston, MA, 1998.
- [2] Andrews (George E.). – MacMahon's partition analysis. II. Fundamental theorems. *Annals of Combinatorics*, vol. 4, n° 3-4, 2000, pp. 327–338. – Conference on Combinatorics and Physics (Los Alamos, NM, 1998).
- [3] Andrews (George E.) and Paule (Peter). – MacMahon's partition analysis. IV. Hypergeometric multisums. *Séminaire Lotharingien de Combinatoire*, vol. 42, n° B42i, 1999. – The Andrews Festschrift (Maratea, 1998). 24 pages.
- [4] Andrews (George E.), Paule (Peter), and Riese (Axel). – MacMahon's partition analysis: the Omega package. *European Journal of Combinatorics*, vol. 22, n° 7, 2001, pp. 887–904.
- [5] Andrews (George E.), Paule (Peter), Riese (Axel), and Strehl (Volker). – MacMahon's partition analysis. V. Bijections, recursions, and magic squares. In *Algebraic combinatorics and applications (Gößweinstein, 1999)*, pp. 1–39. – Springer, Berlin, 2001.
- [6] Bousquet-Mélou (Mireille) and Eriksson (Kimmo). – Lecture hall partitions. *The Ramanujan Journal*, vol. 1, n° 1, 1997, pp. 101–111.
- [7] Bousquet-Mélou (Mireille) and Eriksson (Kimmo). – Lecture hall partitions. II. *The Ramanujan Journal*, vol. 1, n° 2, 1997, pp. 165–185.
- [8] Bousquet-Mélou (Mireille) and Eriksson (Kimmo). – A refinement of the lecture hall theorem. *Journal of Combinatorial Theory. Series A*, vol. 86, n° 1, 1999, pp. 63–84.
- [9] MacMahon (Percy A.). – *Combinatory analysis*. – Chelsea Publishing Co., New York, 1960, xix+302+xix+340p. Two volumes (bound as one).
- [10] Stanley (Richard P.). – Linear homogeneous Diophantine equations and magic labelings of graphs. *Duke Mathematical Journal*, vol. 40, 1973, pp. 607–632.

## Engel Expansions of $q$ -Series

*Peter Paule*

RISC, Linz (Austria)

October 2, 2000

*Summary by Bruno Salvy*

### 1. Engel Expansions

A real number  $A > 0$  has a unique expansion of the form

$$A = a_0 + \frac{1}{a_1} + \frac{1}{a_1 a_2} + \frac{1}{a_1 a_2 a_3} + \dots,$$

where the  $a_i$  are positive integers with  $a_{i+1} \geq a_i$  for  $i \geq 1$ . These expansions were called *Engel expansions* by Perron and their study goes back to Lambert around 1770. Uniqueness of the expansion is not difficult to see, together with the following recurrences from which an iterative algorithm derives:

$$a_k = \lfloor r_k \rfloor + 1, \quad \frac{1}{r_k} = \frac{1}{a_k} + \frac{1}{a_k r_{k+1}}, \quad k \geq 1.$$

The initial conditions are given by  $a_0 < A \leq a_0 + 1$  and  $A - a_0 = 1/r_1$ . Rational numbers are characterized by the ultimate stationarity of the sequence  $(a_i)$ . An obvious example of Engel expansion of a non-rational number is provided by  $e = \exp(1)$  for which  $a_0 = 2$  and  $a_i = i + 1$  for  $i > 0$ .

Arnold and John Knopfmacher defined in [4, 5] an analogous notion for Laurent series.

**Definition 1.** Given a Laurent series  $A = \sum_{n \geq \nu} c_n q^n \in \mathbb{C}((q))$ , and an integer  $\rho \geq 0$ , the  $q$ -Engel sequence associated with  $A$  and  $\rho$  is the unique sequence  $(a_i)$  of polynomials in  $q^{-1}$  such that

$$A = a_0 + \sum_{n \geq 1} \frac{q^{-\rho n}}{a_1 \cdots a_n},$$

with the degrees of the  $a_i$  obeying  $\deg(a_{i+1}) \geq \deg(a_i) + \rho + 1$ .

This definition is motivated by the numerous  $q$ -identities involving such expansions. A sample is given in Table 1, using the classical notations

$$(a; q)_0 = 1, \quad (a; q)_k = (1 - a)(1 - aq) \cdots (1 - aq^{k-1}) \quad \text{for } k > 0, \quad (a; q)_\infty = \prod_{k \geq 0} (1 - aq^k).$$

Again, uniqueness is not difficult to check and an iterative algorithm follows from

$$(1) \quad A_{k+1} := q^\rho (a_k A_k - 1), \quad a_k = \left[ \frac{1}{A_k} \right], \quad k \geq 1,$$

with  $A_0 := A$ ,  $a_0 = [A]$  and  $A_1 = q^\rho (A_0 - a_0)$ . The bracket notation corresponds to the *integral part* of a Laurent series defined by  $[A] := \sum_{\nu \leq n \leq 0} c_n q^n \in \mathbb{C}[q^{-1}]$ .

$$\begin{aligned}
(2) \quad & \sum_{k \geq 0} \frac{q^{\binom{k+1}{2}}}{(q; q)_k} = \prod_{k > 0} (1 + q^k), & \text{(Euler)} \\
(3) \quad & \sum_{k \geq 0} \frac{z^k q^{k^2}}{(q; q)_k (zq; q)_k} = \frac{1}{(z; q)_\infty}, & \text{(Cauchy)} \\
(4) \quad & \sum_{k \geq 0} \frac{q^{k^2}}{(-q; q^2)_k} = \sum_{k \geq 0} \frac{(-1)^{k+1} q^k}{(-q; q)_k}, & \text{(Fine)} \\
(5) \quad & \sum_{k \geq 0} \frac{q^{k(3k-1)/2}}{(q; q)_k (q; q^2)_k} = \prod_{k \geq 1} \frac{(1 - q^{10k-6})(1 - q^{10k-4})(1 - q^{10k})}{1 - q^k}, & \text{(Rogers)} \\
(6) \quad & \sum_{k > 0} \frac{q^{k(3k-1)/2}}{(q; q)_{k-1} (q; q^2)_k} = \prod_{k \geq 1} \frac{(1 - q^{10k-8})(1 - q^{10k-2})(1 - q^{10k})}{1 - q^k}, & \text{(Rogers)} \\
(7) \quad & \sum_{k \geq 0} \frac{q^{k(2k-1)}}{(q; q)_{2k}} = \prod_{k > 0} (1 + q^k), & \text{(Slater)} \\
(8) \quad & \sum_{k \geq 0} \frac{q^{k^2}}{(q; q)_k} = \frac{1}{(q; q^5)_\infty (q^4; q^5)_\infty}, & \text{(1st Rogers–Ramanujan)} \\
(9) \quad & \sum_{k \geq 0} \frac{q^{k^2+k}}{(q; q)_k} = \frac{1}{(q^2; q^5)_\infty (q^3; q^5)_\infty}, & \text{(2nd Rogers–Ramanujan)} \\
(10) \quad & \sum_{k \geq 0} \frac{q^{2k^2}}{(q; q)_{2k}} = \prod_{k > 0, \quad k \equiv \pm 2, \pm 3, \pm 4, \pm 5 \pmod{16}} \frac{1}{1 - q^k}, & \text{(Slater)} \\
(11) \quad & \sum_{k \geq 0} \frac{q^{2k^2+2k}}{(q; q)_{2k+1}} = \prod_{k > 0, \quad k \equiv \pm 1, \pm 4, \pm 6, \pm 7 \pmod{16}} \frac{1}{1 - q^k}, & \text{(Slater)}
\end{aligned}$$

TABLE 1.  $q$ -identities involving  $q$ -Engel expansions.

## 2. Engel Guessing

Equipped with (1), it is very natural to implement a package computing  $q$ -Engel sequences of Laurent series. Such a package opens the way to experimental mathematics with  $q$ -Engel expansions [2]. For instance, starting from a truncation of the series expansion of the right-hand side of (2) (a special case of an identity due to Euler) and using  $\rho = 0$ , the package outputs

$$1 + \frac{q}{(q; q)_1} + \frac{q^3}{(q; q)_2} + \frac{q^6}{(q; q)_3} + O(q^{10}),$$

from which the left-hand side is easily guessed. The task of *proving* such an identity still requires human work.

Using  $\rho = 1$  on the same series does not reveal any pattern. However, with  $\rho = 2$ , one gets

$$1 + \frac{q}{(q; q)_2} + \frac{q^6}{(q; q)_4} + \frac{q^{15}}{(q; q)_6} + O(q^{28}),$$

from which it is easy to conjecture the general formula (7).

### 3. Identities of Rogers–Ramanujan Type

In one of his independent proofs of the Rogers–Ramanujan identities (8–9), Schur introduced two sequences of polynomials

$$d_m = \sum_k (-1)^k q^{k(5k-3)/2} \left[ \begin{matrix} m-1 \\ \lfloor \frac{m+1-5k}{2} \rfloor \end{matrix} \right], \quad e_m = \sum_k (-1)^k q^{k(5k+1)/2} \left[ \begin{matrix} m-1 \\ \lfloor \frac{m-1-5k}{2} \rfloor \end{matrix} \right], \quad m \geq 1,$$

with  $e_0 = 0$  and  $d_0 = 1$  in terms of the *Gaussian polynomials*

$$\left[ \begin{matrix} n \\ k \end{matrix} \right] = \begin{cases} \frac{(q; q)_n}{(q; q)_k (q; q)_{n-k}}, & \text{if } 0 \leq k \leq n, \\ 0, & \text{otherwise.} \end{cases}$$

The sequences  $d_m$  and  $e_m$  appear in the recent generalization of the Rogers–Ramanujan identities due to Garrett, Ismail and Stanton [3]:

$$(12) \quad \sum_{n=0}^{\infty} \frac{q^{n^2+mn}}{(q; q)_n} = \frac{(-1)^m q^{-\binom{m}{2}} d_m}{(q; q^5)_{\infty} (q^4; q^5)_{\infty}} - \frac{(-1)^m q^{-\binom{m}{2}} e_m}{(q^2; q^5)_{\infty} (q^3; q^5)_{\infty}}.$$

Setting  $m = 0$ ,  $m = 1$  in this formula yields (8) and (9).

The left-hand side of (12) is the  $q$ -Engel expansion of the right-hand side for  $\rho = 0$ , which motivates [1] in looking for a  $q$ -Engel “proof” of this identity. For this, it is sufficient to prove that the sequence  $a_n = q^{-(2n+m-1)} - q^{-(n+m-1)}$  is the corresponding  $q$ -Engel sequence. Defining

$$A_0 = A, \quad A_n = (-1)^m q^{-\binom{m}{2} - (m-1)(n-1)} \sum_{j>m} q^{jn} (d_m e_j - d_j e_m) \quad \text{for } n \geq 1,$$

the proof consists in showing that  $a_n A_n = 1 + A_{n+1}$  and  $a_n = [1/A_n]$ , together with correct initial conditions. In view of (14) below, this is not too difficult, but technical (see [1] for details).

Schur proved that both  $d_m$  and  $e_m$  satisfy the recurrence

$$(13) \quad c_{m+2} = c_{m+1} + q^m c_m, \quad m \geq 0.$$

Nowadays, this identity is proved automatically by invoking the  $q$ -WZ algorithm [7] and this leads to the first purely automatic *elementary* proof of the Rogers–Ramanujan identity [6]. In view of this recurrence,  $d_m$  and  $e_m$  are nothing but  $q$ -analogues of the Fibonacci numbers. It turns out that a generalization of the Cassini identity, namely

$$F_{m-1} F_{m+k} - F_{m+k-1} F_m = (-1)^m F_k,$$

admits a  $q$ -analogue in terms of  $e_m$  and  $d_m$ :

$$(14) \quad d_m e_{m+k} - d_{m+k} e_m = (-1)^m q^{\binom{m}{2}} \sum_{j \geq 0} \left[ \begin{matrix} k-1-j \\ j \end{matrix} \right] q^{j^2+mj}.$$

This identity can be proved automatically from (13) by univariate D-finite closure properties ( $m$  being fixed). In fact, a non-Engel proof of (12) follows from letting  $k$  tend to infinity in (14) in view of Schur’s limit formulae

$$d_{\infty} = \frac{1}{(q; q)_{\infty}} \sum_k (-1)^k q^{k(5k-3)/2} = \frac{1}{(q^2; q^5)_{\infty} (q^3; q^5)_{\infty}},$$

$$e_{\infty} = \frac{1}{(q; q)_{\infty}} \sum_k (-1)^k q^{k(5k+1)/2} = \frac{1}{(q; q^5)_{\infty} (q^4; q^5)_{\infty}}.$$

The infinite products are obtained by Jacobi's triple product identity, which also admits a simple computer proof [6].

#### 4. A New Identity Discovered by Engel Guessing

The identities (10) and (11) can be conjectured by Engel guessing after first multiplying the product by  $1 - q$ . An Engel proof is also available [2] using the *Santos polynomials* defined by

$$S_m = \sum_j q^{4j^2-j} \left[ \begin{matrix} m \\ \lfloor \frac{m+1-4j}{2} \rfloor \end{matrix} \right], \quad T_m = \sum_j q^{4j^2-3j} \left[ \begin{matrix} m \\ \lfloor \frac{m+2-4j}{2} \rfloor \end{matrix} \right],$$

whose limits  $S_\infty$  and  $T_\infty$  when  $m \rightarrow \infty$  are precisely the right-hand sides of (10) and (11).

In view of (12), a natural idea consists in experimenting with  $q$ -Engel expansions of  $S_n T_\infty - T_n S_\infty$  or variations of it. It turns out that a pattern readily emerges leading to conjecturing the following generalization of (10) and (11):

$$S_n T_\infty - T_n S_\infty = q^n (q; q^2)_n \sum_{k \geq 0} \frac{q^{2k^2+2(n+1)k}}{(q; q)_{2k+1}}.$$

Again, a possible proof [2] consists in relying on a finite version, namely

$$S_n T_{n+m} - T_n S_{n+m} = q^n (q; q^2)_n \sum_{k \geq 0} \left[ \begin{matrix} m \\ 2k+1 \end{matrix} \right] q^{2k^2+2(n+1)k}.$$

#### 5. Conclusion

Engel expansions are a new way of looking at  $q$ -identities which allows for easy computer experiments and hence should lead to many discoveries. A pending issue is to make  $q$ -Engel proving into an algorithmic task.

#### Bibliography

- [1] Andrews (George E.), Knopfmacher (Arnold), and Paule (Peter). – An infinite family of Engel expansions of Rogers-Ramanujan type. *Advances in Applied Mathematics*, vol. 25, n° 1, 2000, pp. 2–11.
- [2] Andrews (George E.), Knopfmacher (Arnold), Paule (Peter), and Zimmermann (Burkhard). – Engel expansions of  $q$ -series by computer algebra. In Alladi (K.) (editor), *q-Series. Kluwer Series Developments in Mathematics*. – Kluwer, 2000. Proceedings of the Conference “ $q$ -Series”, Gainesville, November 1999.
- [3] Garrett (Kristina), Ismail (Mourad E. H.), and Stanton (Dennis). – Variants of the Rogers-Ramanujan identities. *Advances in Applied Mathematics*, vol. 23, n° 3, 1999, pp. 274–299.
- [4] Knopfmacher (Arnold) and Knopfmacher (John). – Inverse polynomial expansions of Laurent series. *Constructive Approximation*, vol. 4, n° 4, 1988, pp. 379–389.
- [5] Knopfmacher (Arnold) and Knopfmacher (John). – Inverse polynomial expansions of Laurent series. II. In *Proceedings of the 3rd International Congress on Computational and Applied Mathematics (Leuven, 1988)*, vol. 28, pp. 249–257. – 1989.
- [6] Paule (Peter). – Short and easy computer proofs of the Rogers-Ramanujan identities and of identities of similar type. *Electronic Journal of Combinatorics*, vol. 1, 1994. – Research Paper 10. 9 pages.
- [7] Wilf (Herbert S.) and Zeilberger (Doron). – An algorithmic proof theory for hypergeometric (ordinary and “ $q$ ”) multisum/integral identities. *Inventiones Mathematicae*, vol. 108, n° 3, 1992, pp. 575–633.

# Eulerian Calculus: a Technology for Computer Algebra and Combinatorics

*Dominique Foata*

Département de mathématique, Université Louis Pasteur (France)

May 21, 2001

*Summary by Dominique Gouyou-Beauchamps*

## Abstract

Babson and Steingrímsson have introduced pairs of permutation statistics that they conjectured were all Euler–Mahonian, i.e., equidistributed with the pair  $(\text{des}, \text{maj})$  where  $\text{des}$  is the number of descents and  $\text{maj}$  is the major index. How to prove their conjecture? We use the so-called “Umbral Transfer Matrix Method” implemented by Zeilberger and specific combinatorial constructions leading to new transformations on the symmetric group. Details may be found in the recent work of D. Foata and D. Zeilberger [2].

## 1. Introduction

We use the Babson–Steingrímsson notation [1] for “atomic” permutation statistics. Given a permutation  $w = x_1x_2 \dots x_n$  of  $1, 2, \dots, n$  they denote  $(a - bc)(w)$  the number of occurrences of the pattern  $a - bc$ , i.e., the number of places  $1 \leq i < j < n$  such that  $x_i < x_j < x_{j+1}$ . Similarly, the pattern  $(b - ca)(w)$  is the number of occurrences of  $x_{j+1} < x_i < x_j$ , and in general, for any permutation  $\alpha, \beta, \gamma$  of  $a, b, c$ , the expression  $(\alpha - \beta\gamma)(w)$  is the number of pairs  $(i, j)$ ,  $1 \leq i < j < n$ , such that the orderings of the two triples  $(x_i, x_j, x_{j+1})$  and  $\alpha, \beta, \gamma$  are identical. The statistic  $(ab - c)$  is defined in the same way by looking at the occurrences  $(x_i, x_{i+1}, x_j)$  such that  $i + 1 < j$  and  $x_i < x_{i+1} < x_j$ . Of course,  $(ba)(w)$  denotes the number  $\text{des } w$  of *descents* of  $w$  (i.e., the number of places  $1 \leq i < n$  such that  $x_i > x_{i+1}$ ) and  $(ab)(w)$  denotes the number  $\text{rise } w$  of *rises* of  $w$  (i.e., the number of places  $1 \leq i < n$  such that  $x_i < x_{i+1}$ ).

The classical permutation statistics  $\text{inv}$  and  $\text{maj}$  may be written as  $(bc - a) + (ca - b) + (cb - a) + (ba)$  and  $(a - cb) + (b - ca) + (c - ba) + (ba)$ , respectively. This inspired Babson and Steingrímsson to perform a computer search for all statistics that could be thus written, and look for those that appear to be Mahonian. They came up with a list of 18. Some of them turned out to be well-known, and some were new. Yet eight new conjecturally Mahonian statistics were left open. Here we prove four of them.

## 2. Notations

Recall the usual notations

$$(a; q)_n = \begin{cases} 1 & \text{if } n = 0, \\ (1 - a)(1 - aq) \dots (1 - aq^{n-1}) & \text{if } n \geq 1, \end{cases}$$
$$(a; q)_\infty = \lim_{n \rightarrow \infty} (a; q)_n = \prod_{n \geq 0} (1 - aq^n),$$

$$[n]_q = \frac{1 - q^n}{1 - q} = \sum_{i=0}^{n-1} q^i, \quad [n]_q! = \frac{(q; q)_n}{(1 - q)^n} = \prod_{i=1}^n [i]_q.$$

A statistic  $\text{stat}$  on the symmetric group  $\mathcal{S}_n$  is said to be *Mahonian*, if for every  $n \geq 0$  we have

$$\sum_{w \in \mathcal{S}_n} q^{\text{stat } w} = [n]_q!$$

A sequence  $(A_n(t, q))_{n \geq 0}$  of polynomials in two variables  $t$  and  $q$ , is said to be *Euler–Mahonian*, if one of the following equivalent conditions holds:

1. For every  $n \geq 0$ ,

$$\frac{1}{(t; q)_{n+1}} A_n(t, q) = \sum_{s \geq 0} t^s ([s + 1]_q)^n.$$

2. The exponential generating function for the fractions  $\frac{A_n(t, q)}{(t; q)_{n+1}}$  is given by

$$\sum_{n \geq 0} \frac{u^n}{n!} \frac{A_n(t, q)}{(t; q)_{n+1}} = \sum_{s \geq 0} t^s \exp(u[s + 1]_q).$$

3. The sequence  $(A_n(t, q))$  satisfies the recurrence relation:

$$(1) \quad (1 - q)A_n(t, q) = (1 - tq^n)A_{n-1}(t, q) - q(1 - t)A_{n-1}(tq, q).$$

4. Let  $A_n(t, q) = \sum_{s \geq 0} t^s A_{n,s}(q)$ . Then the coefficients  $A_{n,s}(q)$  satisfy the recurrence:

$$A_{n,s}(q) = [s + 1]_q A_{n-1,s}(q) + q^s [n - s]_q A_{n-1,s-1}(q).$$

Now a pair of statistics  $(\text{stat}_1, \text{stat}_2)$  defined on each symmetric group  $\mathcal{S}_n$  ( $n \geq 0$ ) is said to be *Euler–Mahonian*, if for every  $n \geq 0$  we have

$$\sum_{w \in \mathcal{S}_n} t^{\text{stat}_1 w} q^{\text{stat}_2 w} = A_n(t, q).$$

### 3. Results

Our results are the following:

**Theorem 1.** *The permutation statistic  $S11 = (a - cb) + 2(b - ca) + (ba)$  is Mahonian.*

**Theorem 2.** *The permutation statistic  $S13 = (a - cb) + 2(b - ac) + (ab)$  is Mahonian.*

**Theorem 3.** *Let  $S5 = (b - ca) + (c - ba) + (a - bc) + (ab)$ . Then, the pair  $(\text{rise}, S5)$  is Euler–Mahonian.*

**Theorem 4.** *Let  $S6 = (ba - c) + (c - ba) + (ac - b) + (ba)$ . Then, the pair  $(\text{des}, S6)$  is Euler–Mahonian.*

Our Theorems 1, 2, and 4 are the three parts of Conjecture 8 of [1], while Theorem 3 is Conjecture 10 of [1]. It turns out that, thanks to Zeilberger’s recent theory of the *Umbral Transfer Matrix Method* [4], the proofs of the first three theorems are completely automatic, using the general Maple package `ROTA`, together with a new interfacing package `PERCY` that computes the appropriate Rota operators for what we will call *Markovian Permutation Statistics*.

However, `ROTA` is useless in the case of  $S6$ . So proving Theorem 4 still requires the traditional combinatorial method: construct a bijection  $w \mapsto w'$  of  $\mathcal{S}_n$  onto itself which has the property that

$$(\text{des}, S6) w' = (\text{des}, \text{maj}) w$$



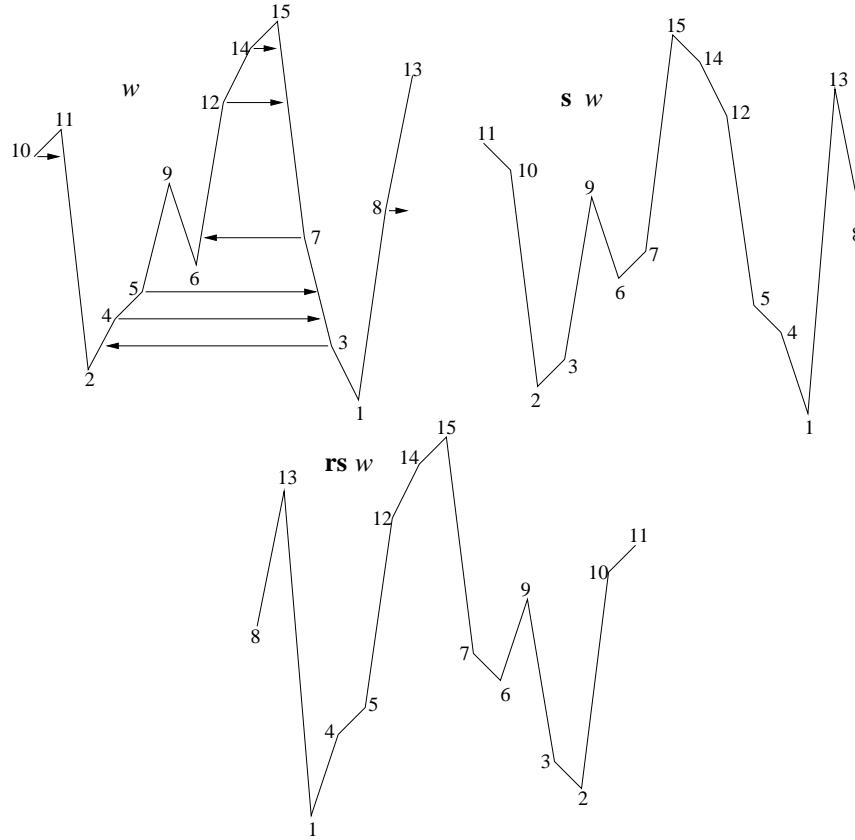


FIGURE 1.  $\mathbf{rs}(10, 11, 2, 4, 5, 9, 6, 12, 14, 15, 7, 3, 1, 8, 13) = (8, 13, 1, 4, 5, 12, 14, 15, 7, 6, 9, 3, 2, 10, 11)$ .

holds for every  $w \in \mathcal{S}_n$ .

#### 4. Proof of Theorem 4

Instead of the pair  $(\text{des}, \text{maj})$  we will take another Euler–Mahonian pair  $(\text{des}, \text{mak})$ , where  $\text{mak}$  is a Mahonian statistic that was introduced by Foata and Zeilberger in [3]. In the Babson–Steingrímsson notation  $\text{mak}$  reads

$$\text{mak} := (a - cb) + (cb - a) + (ba) + (ca - b).$$

First, the *descent bottom* of a permutation  $x_1x_2 \dots x_n$  is defined to be the set  $\text{desbot } w$  of all the  $x_i$ 's such that  $2 \leq i \leq n$  and  $x_{i-1} > x_i$ . Its cardinality is the number  $\text{des } w$  of descents of  $w$ .

Next, the word statistics  $U$  and  $V$  are introduced as follows. Let  $y = x_i$  be a letter of the permutation  $w = x_1x_2 \dots x_n$ . Define

$$U_y(w) = (ca - b)|_{b=y} w; \quad V_y(w) = (b - ac)|_{b=y} w.$$

Thus,  $U_y(w)$  is the number of adjacent letters  $x_jx_{j+1}$  to the left of  $y = x_i$  such that  $x_j > x_i > x_{j+1}$ . The word statistics  $U$  and  $V$  are then:

$$U(w) = U_1(w)U_2(w) \dots U_n(w); \quad V(w) = V_1(w)V_2(w) \dots V_n(w).$$

Now, recall the traditional reverse image  $\mathbf{r}$ , which is an involution that maps each permutation  $w = x_1x_2 \dots x_n$  onto  $\mathbf{r}w = x_nx_{n-1} \dots x_1$ . We shall introduce another involution  $\mathbf{s}$  of  $\mathcal{S}_n$ , called the

*rise-des-exchange*, which exchanges the rises and the descents of a permutation and keeps peaks and troughs in their original ordering. The involution  $\mathbf{s}$  is not explained here, but can be immediately visualized in Fig. 1.

**Proposition 1.** *The involution  $\mathbf{r s}$  of  $\mathcal{S}_n$  has the following properties:*

1.  $\text{desbot } \mathbf{r s } w = \text{desbot } w$ ,
2.  $(U, V) \mathbf{r s } w = (V, U) w$ .

Let  $\Sigma$ - $\text{desbot } w$  be the sum of all the letters  $x_i$  of the permutation  $w = x_1 x_2 \dots x_n$  which belong to the descent bottom set  $\text{desbot } w$ .

**Proposition 2.** *For each permutation  $w$  we have:*

$$\Sigma\text{-desbot } w = ((a - cb) + (cb - a) + (ba)) w.$$

Next, we introduce the *complement* to  $(n + 1)$ , denoted by  $\mathbf{c}$ , that maps each permutation  $w = x_1 x_2 \dots x_n$  onto  $\mathbf{c} w = (n + 1 - x_1)(n + 1 - x_2) \dots (n + 1 - x_n)$ . Thus the statistic  $S6 \mathbf{r c}$  reads

$$S6 \mathbf{r c} = (a - cb) + (cb - a) + (ba) + (b - ca).$$

Taking Proposition 2 into account, we get the expressions:

$$\text{mak } w = \Sigma\text{-desbot } w + U_1(w) + \dots + U_n(w),$$

$$S6 \mathbf{r c} = \Sigma\text{-desbot } w + V_1(w) + \dots + V_n(w).$$

Therefore, Proposition 1 implies the following corollary.

**Corollary 1.** *The involution  $\mathbf{r s}$  is an involution of  $\mathcal{S}_n$  having the property:*

$$(\text{des}, \text{mak}) w = (\text{des}, S6 \mathbf{r c}) \mathbf{r s } w.$$

But  $(\text{des}, \text{mak})$  is Euler–Mahonian, as proved in [3]. Therefore, the pair  $(\text{des}, S6 \mathbf{r c})$  is Euler–Mahonian, as well as  $(\text{des}, S6)$ , since we always have  $\text{des } \mathbf{r c } w = \text{des } w$ . Hence Theorem 4 is proved.

## 5. Markovian Permutation Statistics

The *reduction* of a sequence  $w$  of  $n$  distinct integers, denoted by  $\text{red}(w)$ , is the permutation obtained by replacing the smallest member by 1, the second-smallest by 2,  $\dots$ , and the largest by  $n$ . For example  $\text{red}(5 \ 8 \ 3 \ 7 \ 4) = 3 \ 5 \ 1 \ 4 \ 2$ .

A permutation statistic  $F : \mathcal{S}_n \rightarrow \mathbb{Z}$  is said to be *Markovian*, if there exists a function  $h(j, i, n)$  such that

$$F(x_1 \dots x_n) = F(\text{red}(x_1 \dots x_{n-1})) + h(x_{n-1}, x_n, n).$$

A Markovian permutation statistic  $F : \mathcal{S}_n \rightarrow \mathbb{Z}$  is said to be *nice Markovian* if the above  $h(j, i, n)$  can be written as

$$h(j, i, n) = \begin{cases} f(j, i, n) & \text{if } j < i, \\ g(j, i, n) & \text{if } j > i, \end{cases}$$

where  $f$  and  $g$  are *affine linear functions* of their arguments, i.e., can be written as  $ai + bj + cn + d$ , for some integers  $a, b, c, d$ .

We, and the Maple package PERCY, will only consider nice Markovian statistics. We will denote them by  $[f, g, j, i, n]$ . For example,  $\text{inv} = [n - i, n - i, j, i, n]$ ,  $\text{maj} = [0, n - 1, j, i, n]$ ,  $\text{des} = [0, 1, j, i, n]$ ,  $\text{rise} = [1, 0, j, i, n]$ .

Given a permutation statistic  $F$  we are interested in the sequence of polynomials

$$\text{gf}(F)_n(q) = \sum_{w \in \mathcal{S}_n} q^{F(w)} \quad (n \geq 0).$$

However, in order to take advantage of Markovity, we need to consider the more refined

$$\text{GF}(F)_n(q, z) = \sum_{w=x_1 \dots x_n \in \mathcal{S}_n} q^{F(w)} z^{x_n} \quad (n \geq 0)$$

that also keeps track of the last letter  $x_n$ . Now, by using *Rota operators* [4], it is easy to express  $\text{GF}(F)_n$  in terms of  $\text{GF}(F)_{n-1}$ . Let  $w' = x'_1 \dots x'_{n-1} = \text{red}(x_1 \dots x_{n-1})$ ; then

$$\begin{aligned} \text{GF}(F)_n(q, z) &= \sum_{i=1}^n z^i \sum_{w \in \mathcal{S}_n; x_n=i} q^{F(w)} \\ &= \sum_{j=1}^{n-1} \sum_{w' \in \mathcal{S}_{n-1}; x'_{n-1}=j} \left( \sum_{i=1}^j q^{g(j+1, i, n)} z^i + \sum_{i=j+1}^n q^{f(j, i, n)} z^i \right) q^{F(w')}. \end{aligned}$$

Now for  $i \leq j \leq n-1$  we introduce the *umbra*  $\mathcal{P}$ ,

$$\mathcal{P}(z^j) = \left( \sum_{i=1}^j q^{g(j+1, i, n)} z^i + \sum_{i=j+1}^n q^{f(j, i, n)} z^i \right),$$

and we extend by linearity, so that  $\mathcal{P}$  is defined on all polynomials of degree less than or equal to  $n-1$ . In terms of  $\mathcal{P}$ , we have the very simple recurrence:

$$\text{GF}(F)_n(q, z) = \mathcal{P}(\text{GF}(F)_{n-1}(q, z)).$$

Maple can compute the umbra automatically. All the users have to enter is  $f$  and  $g$ , and PERCY would convert it to the Markovian notation.

### 6. Proof of Theorem 1

Using PERCY and ROTA we get that the umbra  $\mathcal{P}$  linking  $\text{GF}(\text{S11})_{n-1}(q, z)$  to  $\text{GF}(\text{S11})_n(q, z)$  maps the polynomial  $a(z)$  onto

$$\frac{z^{n+1}a(1) - za(z)}{z-1} + \frac{z(a(qz) - a(q^2))}{z-q}.$$

Hence  $b_n(z) = \text{GF}(\text{S11})_n(q, z)$  satisfies the functional recurrence

$$b_n(z) = \frac{z^{n+1}b_{n-1}(1) - zb_{n-1}(z)}{z-1} + \frac{z(b_{n-1}(qz) - b_{n-1}(q^2))}{z-q},$$

with the initial condition  $b_1(z) = z$ . But if we guess (and if we check) that the sequence

$$c_n(z) = z \frac{z^n - q^n}{z - q} [n-1]_q!$$

satisfies the same recurrence, we obtain that  $b_n(z) = c_n(z)$ , and finally that  $b_n(1) = c_n(1) = [n]_q!$

### 7. Proof of Theorem 2

Using PERCY and ROTA we get that the umbra  $\mathcal{P}$  linking  $\text{GF}(\text{S13})_{n-1}(q, z)$  to  $\text{GF}(\text{S13})_n(q, z)$  maps the polynomial  $a(z)$  onto

$$\frac{z(a(zq) - a(1))}{qz-1} + \frac{zqa(z) - q^{2n+1}z^{n+1}a(q^{-2})}{1-zq^2}.$$

Hence  $d_n(z) = \text{GF}(\text{S13})_n(q, z)$  satisfies the functional recurrence

$$d_n(z) = \frac{z(d_{n-1}(zq) - d_{n-1}(1))}{qz - 1} + \frac{zqd_{n-1}(z) - q^{2n+1}z^{n+1}d_{n-1}(q^{-2})}{1 - zq^2},$$

with the initial condition  $d_1(z) = z$ . But if we guess (and if we check) that the sequence

$$e_n(z) = z \frac{(1 - z^n q^n)}{1 - qz} [n - 1]_q!$$

satisfies the same recurrence, we obtain that  $d_n(z) = e_n(z)$ , and finally that  $d_n(1) = e_n(1) = [n]_q!$ .

## 8. Proof of Theorem 3

PERCY can compute the Umbra multi-statistics, when the generating function is the weight-enumerator of  $\mathcal{S}_n$  according to the weight

$$\text{weight}(w) = z^{x_n} \prod_{j=1}^r q_j^{F_j(w)},$$

where  $w = x_1 \dots x_n$  and  $F_1(w), \dots, F_r(w)$  are several nice Markovian permutation statistics. Define

$$A_n(t, q; z) = \sum_{w \in \mathcal{S}_n} t^{\text{des } w} q^{\text{maj } w} z^{x_n}, \quad B_n(t, q; z) = \sum_{w \in \mathcal{S}_n} t^{\text{rise } w} q^{\text{S5 } w} z^{x_n}.$$

PERCY and ROTA compute the following functional recurrences

$$(2) \quad A_n(t, q; z) = \frac{z(1 - tq^{n-1})A_{n-1}(t, q; z) - z(z^n - tq^{n-1})A_{n-1}(t, q; 1)}{1 - z},$$

$$B_n(t, q; z) = \frac{z(1 - tq^n)B_{n-1}(t, q; z) - z(1 - tz^n)B_{n-1}(t, q; q)}{z - q}.$$

By comparing the two functional recurrences, we guess and we verify that

$$B_n(t, q; z) = q^{-n} z^{n+1} A_n(tq, q; q/z).$$

Hence  $B_n(t, q; 1) = q^{-n} A_n(tq, q; q)$ . By plugging  $t = tq$ ,  $z = q$  into Eq. (2), we get that

$$A_n(tq, q; q) = q^n \frac{(1 - tq^n)A_{n-1}(t, q; 1) - q(1 - t)A_{n-1}(tq, q; 1)}{1 - q}.$$

But, this equals  $q^n A_n(t, q)$  by Eq. (1). And we have proved that  $B_n(t, q; 1) = A_n(t, q; 1) = A_n(t, q)$ .

The input and output files of PERCY can be downloaded from

<http://www.math.temple.edu/~zeilberg/programs.html>.

## Bibliography

- [1] Babson (Eric) and Steingrímsson (Einar). – Generalized permutation patterns and a classification of the Mahonian statistics. *Séminaire Lotharingien de Combinatoire*, vol. 44, n° B44b, 2000. – 18 pages. Available from <http://www.mat.univie.ac.at/~slc/>.
- [2] Foata (D.) and Zeilberger (D.). – Babson-Steingrímsson statistics are indeed Mahonian (and sometimes even Euler-Mahonian). – To appear in *Advances in Applied Mathematics*.
- [3] Foata (Dominique) and Zeilberger (Doron). – Denert's permutation statistic is indeed Euler-Mahonian. *Studies in Applied Mathematics*, vol. 83, n° 1, 1990, pp. 31–59.
- [4] Zeilberger (Doron). – The umbral transfer-matrix method. I. Foundations. *Journal of Combinatorial Theory. Series A*, vol. 91, n° 1-2, 2000, pp. 451–463. – In memory of Gian-Carlo Rota.

## Part II

# Analysis of Algorithms and Combinatorial Structures



## Asymptotics for Random Combinatorial Structures

*Amir Dembo*

Mathematics and Statistics Department, Stanford University (USA)

June 18, 2001

*Summary by Philippe Flajolet*

### Abstract

What does a random partition of a large integer look like? The talk presents asymptotic results and variational problems for this question, obtained in a work of A. Dembo jointly with A. Vershik and O. Zeitouni [4]. The techniques involve some combinatorics and mostly probability theory. Other applications concern asymptotics of various random combinatorial structures, such as permutations, forests of trees, and convex polygons with integer vertices. This summary is intended as a casual introduction to the reading of the paper [4].

### 1. A Bit of Paleontology

A *partition* of the integer  $n$  is an additive decomposition of the integer  $n$  into some number  $r$  of integer summands,

$$n = x_1 + x_2 + \cdots + x_r, \quad x_j \geq x_{j+1}, \quad x_r > 0.$$

The quantity  $r$  is called the number of summands (or parts). A partition is said to be *strict* if all its summands are distinct. A partition is naturally represented by a diagram resembling a staircase and called diversely its Ferrers graph or its Young diagram. We shall let  $\mathcal{P}_n$  and  $\mathcal{P}_n^s$  denote the collections of all partitions and strict partitions summing to  $n$ , and denote with  $P_n$ ,  $P_n^s$  the corresponding cardinalities.

Euler started the analytic theory of partitions by providing the explicit generating functions

$$P(z) = \sum_n P_n z^n = \prod_{k \geq 1} \frac{1}{1 - z^k}, \quad P^s(z) = \sum_n P_n^s z^n = \prod_{k \geq 1} (1 + z^k),$$

and a good deal more. The next century mostly focussed on the corresponding theta function identities and their elliptic-modular aspects. Andrews' classic [2] is still a pretty good reference for many of these aspects.

The asymptotic theory starts 150 years after Euler, with the first letters of Ramanujan to Hardy in 1913; see [7]. There, Ramanujan stated:

“The coefficient of  $x^n$  in  $(1 - 2x + 2x^4 - 2x^9 + \cdots)^{-1}$  is the integer nearest to

$$\frac{1}{4n} \left( \cosh \pi \sqrt{n} - \frac{\sinh \pi \sqrt{n}}{\pi \sqrt{n}} \right).”$$

This assertion (that in fact needs to be mildly amended) is, in view of Euler's pentagonal number theorem, directly relevant to our subject. In a celebrated series of memoirs published in 1917

and 1918, Hardy and Ramanujan found very precise estimates for the partition numbers implying in particular:

$$(1) \quad P_n \sim \frac{1}{4n\sqrt{3}} e^{\pi\sqrt{\frac{2}{3}n}}, \quad P_n^s \sim \frac{1}{4\sqrt{3}n^{3/4}} e^{\pi\sqrt{\frac{1}{3}n}}.$$

See [7, Ch VIII] for an insightful discussion and [2, Ch. 5] for the reexamination of the subject by Meinardus.

As far as dominant asymptotics goes, it may be worth pointing out that the simply stated estimates (1) plainly derive from a saddle-point approximation of the Cauchy coefficient integral,

$$(2) \quad [z^n]C(z) = \frac{1}{2i\pi} \int_{|z|=\zeta} C(z) \frac{dz}{z^{n+1}}.$$

The saddle points to be used here are at the real points  $\zeta$  (for  $P(z)$ ) and  $\zeta^s$  (for  $P^s(z)$ ) such that

$$\zeta \sim 1 - \frac{\pi}{\sqrt{6n}}, \quad \zeta^s \sim 1 - \frac{\pi}{\sqrt{12n}},$$

the reason being that  $P(z)$  and  $P^s(z) = \frac{P(z)}{P(z^2)}$  tend to infinity like  $\exp\left(\frac{\pi^2}{6(1-z)}\right)$  and  $\exp\left(\frac{\pi^2}{12(1-z)}\right)$ , as  $z$  tends radially to 1.

Later in the last century, Erdős and Lehner [6] launched the study of various characteristics of random partitions. In particular, they showed that almost all partitions of  $\mathcal{P}_n$  have a number a summands in an interval

$$(3) \quad \frac{1}{C} \sqrt{n} \log n \pm o(\sqrt{n} \log n), \quad C := \pi \sqrt{\frac{2}{3}},$$

while for strict partitions, the interval is

$$(4) \quad \frac{2\sqrt{n}}{D} \log 2 \pm o(\sqrt{n}), \quad D := \pi \sqrt{\frac{1}{3}}.$$

The limit law is an extreme value distribution in the first case, a Gaussian distribution in the second case. (Erdős and Lehner use a mostly elementary recurrence argument induced by generating functions together with the Hardy–Ramanujan estimates.) Note the similarities between the saddle-point constants and the normalization constants  $C$ ,  $D$  in (3) and (4). Also, the scaling factor  $\sqrt{n}$  is ubiquitous in all such analyses. Roughly put, these estimates inform us that a random partition of  $n$  is expected to fit in a rectangle with sides about  $n^{1/2+o(1)}$ .

## 2. The Shape of Random Partitions

Around 1977, Vershik and Kerov [10] and, independently, Logan and Shepp [8] studied the shape of the Young tableau(s) associated to a random permutation or a random involution.<sup>1</sup> Thus, in contrast to what happens in the talk, we are momentarily dealing with a non-uniform distribution on  $\mathcal{P}_n$ . Indeed, the enumerative formulae relative to Young tableaux under these statistics (the “hook formula,” also called the Robinson–Frame–Thrall formula) renormalize in the scale of  $\sqrt{n}$  in such a way that the probability of a continuous shape  $f(t)$  (in the asymptotic limit) occurring in tableaux of size  $n$  is found to be of the rough form (see (7) below for a precise statement)

$$(5) \quad e^{-n\Theta(f)}, \quad \Theta(f) := 2 \iint_{t < s} \log\left(2e^{1/2}(s-t)\right) (1 - \dot{f}(s))(1 + \dot{f}(t)) dt ds$$

<sup>1</sup>These are “filled” Young diagrams—the filling rule corresponds to entries increasing by line and column.



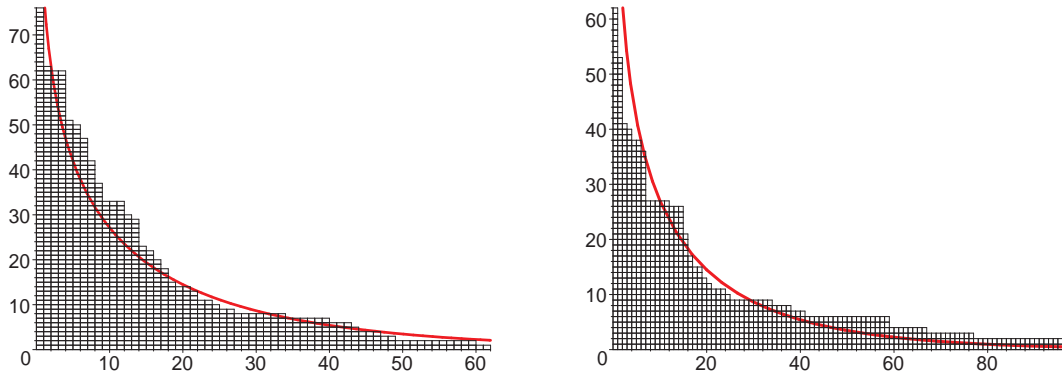


FIGURE 1. Two partitions of  $\mathcal{P}_{1000}$  drawn at random against the limiting shape  $\Psi(t)$ .

with  $\dot{f}$  the derivative of  $f$ . Thus, the most likely shape  $f_0$  solves the variational problem of minimizing the functional  $\Theta$ , and “most” tableaux are expected to be close to this particular shape  $f_0$ . From the methodological standpoint, the contributions [8, 10] are especially important. They led to a much wanted solution of Golomb’s conjecture to the effect that the average length of the longest increasing subsequence in a random permutation of size  $n$  is asymptotic to  $2\sqrt{n}$ ; see [1, 3] for recent developments in rather different directions.

We now return to partitions and let  $Q_n$  and  $Q_n^s$  represent the uniform probability models on  $\mathcal{P}_n$  and  $\mathcal{P}_n^s$ . A partition (or diagram)  $\lambda$  can be written under the form  $\lambda = 1^{r_1}2^{r_2}3^{r_3}\dots$ . Graphically, we define the “contour” or “shape,”

$$\varphi_\lambda(t) := \sum_{k=\lceil t \rceil}^{\infty} r_k, \quad t \geq 0,$$

so that  $\varphi_\lambda$  is a monotone decreasing function whose integral over  $\mathbb{R}^+$  equals  $n$ . We normalize any such  $\varphi$  by

$$\tilde{\varphi}_n(t) = \frac{1}{\sqrt{n}}\varphi_\lambda(\lceil t\sqrt{n} \rceil).$$

Under the models induced by  $Q_n$  and  $Q_n^s$ , Vershik [9] proved (in the sense of uniform convergence on compact sets) that contours tend to converge to deterministic limits,

$$\tilde{\varphi}(\cdot) \xrightarrow[n \rightarrow \infty]{} \Psi(\cdot), \quad \tilde{\varphi}^s(\cdot) \xrightarrow[n \rightarrow \infty]{} \Psi^s(\cdot), \quad \text{where}$$

$$(6) \quad \Psi(t) := \int_t^\infty \frac{du}{e^{\alpha u} - 1} du, \quad \Psi^s(t) := \int_t^\infty \frac{du}{e^{\beta u} + 1} du, \quad \alpha = \frac{\pi}{\sqrt{6}}, \quad \beta = \frac{\pi}{\sqrt{12}}.$$

Thus, a random partition under  $Q_n$  or  $Q_n^s$  tends to have a limiting shape given by the curves  $\Psi(t)$  or  $\Psi^s$ ; see Fig. 1 obtained with Maple and combstruct. (Observe that  $\Psi$  has a logarithmic singularity at 0, while  $\Psi^s$  is regular there.) Alternatively, the limit contours are the curves satisfying respectively  $e^{-\alpha x} + e^{-\alpha y} = 1$ , and  $e^{\beta y} - e^{-\beta y} = 1$ , with the first one being symmetrical, as it should.

The main objective of the talk is to consider deviations from the limit shapes. What is proved is a full large deviation principle, of speed  $\sqrt{n}$ , much in the spirit of (5). We recall that a sequence of measures  $\mu_n$  over a (completely regular Hausdorff topological) space  $\mathcal{X}$  is said to satisfy the *large deviation principle [LDP]* with speed  $b_n$  and a rate function  $I : \mathcal{X} \rightarrow [0, \infty)$  is lower

semicontinuous, and for any measurable set  $X \subset \mathcal{X}$ , there holds:

$$(7) \quad - \inf_{x \in X^\circ} I(x) \leq \liminf_{n \rightarrow \infty} \frac{1}{b_n} \log \mu_n(X) \leq \limsup_{n \rightarrow \infty} \frac{1}{b_n} \log \mu_n(X) \leq - \inf_{x \in \overline{X}} I(x).$$

There,  $X^\circ$  and  $\overline{X}$  denote the interior and the closure of  $X$ . It can be recognized that the informally stated estimate in (5) is of this type (with speed  $b_n = n$  and rate function  $\Theta$ ).

For our purposes, the set  $\mathcal{X}$  will consist of functions that are left continuous and of right limits equipped with the topology of uniform convergence. Let also  $\mathcal{AC}_\infty^{[-1,0]}$  be the subset of non-increasing absolutely continuous functions  $f(\cdot)$  satisfying  $\lim_{t \rightarrow \infty} f(t) = 0$ —and hence  $f(t) = \int_t^\infty (-\dot{f}(u)) du$ —with derivatives belonging Lebesgue-almost everywhere to the interval  $[-1, 0]$ . This last set represents the collection of all potential “shapes” of partitions considered (after normalization). By want of space, we refer to the original paper [4] for complete topological and measure-theoretic definitions and state:

**Theorem 1.** *Under the laws  $Q_n^s$ , the random variable  $\tilde{\varphi}(\cdot)$  satisfies the LDP with speed  $\sqrt{n}$  and a rate function that, for  $f \in \mathcal{AC}_\infty^{[-1,0]}$  and  $\int_0^\infty (-t) df(t) \leq 1$ , is expressed by*

$$I^s(f) = 2\beta - \int_0^\infty h(-\dot{f}_{ac}(t)) dt,$$

with  $h(x) = \log(x^{-x}(1-x)^{-(1-x)})$  the entropy function and  $g_{ac}$  the absolutely continuous part of  $g$ .

**Theorem 2.** *Under the laws  $Q_n$ , the random variable  $\tilde{\varphi}(\cdot)$  satisfies the LDP with speed  $\sqrt{n}$  and a rate function that, for  $f$  in a suitable space and  $\int_0^\infty (-t) df(t) \leq 1$ , is expressed by*

$$I^s(f) = 2\alpha - \int_0^\infty \left(1 - \dot{f}_{ac}(t)\right) h\left(\frac{-\dot{f}_{ac}(t)}{1 - \dot{f}_{ac}(t)}\right) dt.$$

The paper also states some equivalent forms that are expressed in terms of a “distance” to the most likely contours of (6). That distance involves various entropy functions.

### 3. Boltzmann Models of Combinatorics

The first step in the proof of Theorems 1 and 2 is the introduction of a family of models over the classes  $\mathcal{P}$  and  $\mathcal{P}^s$  and large deviations are first established under these models. Since the principles are of an applicability that goes well beyond the probabilistic theory of partitions, we depart a bit from the original paper [4] and discuss them first at a fair level of generality.

Let generally  $\mathcal{C}$  be a class of combinatorial objects endowed with its size function  $|\cdot|$ . What we call here, by virtue of a vague analogy with statistical mechanics, the *Boltzmann model* of parameter  $x$  (over  $\mathcal{C}$ ) is the model that assigns to any object  $\gamma \in \mathcal{C}$  the probability

$$\frac{x^{|\gamma|}}{C(x)} \quad \text{with} \quad C(x) = \sum_{\gamma \in \mathcal{C}} x^{|\gamma|},$$

the counting generating function of  $\mathcal{C}$ . There  $x$  is to be restricted to real values less than the radius  $\rho$  of convergence of  $C(x)$ .

The class  $\mathcal{C}$  being fixed, we shall let  $Q_n$  denote the uniform probability model over the subclass  $\mathcal{C}_n$  of objects of size  $n$  and, with a slight abuse of notations,  $Q_x$  represents the Boltzmann model of parameter  $x$ . Clearly,  $Q_x$  is a mixture of the family of models  $\{Q_n\}$  in the following sense:

$$(8) \quad Q_x \cong Q_N \quad \text{where } N \text{ is a random integer selected with } \mathbf{P}(N = n) = \frac{C_n x^n}{C(x)}.$$

In other words, a randomly chosen object under  $Q_x$  has a random size  $N \equiv N_x$  distributed according to the probability in (8); once the value of size has been drawn according to its distribution, say,  $N = n$ , a random element of  $\mathcal{C}_n$  is chosen uniformly at random, that is, according to  $Q_n$ . (Accordingly,  $Q_n$  is  $Q_x$  conditioned upon size, irrespective of the value of  $x \in (0, \rho)$ .) The distribution of the random size  $N$  according to  $Q_x$  is itself given by a simple generic calculation that we now explain. The probability generating function of  $N$  is

$$\sum_n \mathbf{P}(N = n) z^n = \frac{C(xz)}{C(x)}.$$

Next, the mean and second moment of  $N$  are found to be

$$(9) \quad \mathbf{E}(N) = x \frac{C'(x)}{C(x)}, \quad \mathbf{E}(N^2) = \frac{x^2 C''(x) + x C'(x)}{C(x)}.$$

The mean size increases as  $x$  approaches  $\rho^-$ , with  $\rho$  the radius of convergence of  $C$ . In particular, if the additional condition  $C'(\rho^-) = +\infty$  is met, the Boltzmann model must give preponderance to objects of larger and larger sizes. (Work in progress by Duchon, Flajolet, Louchard, and Schaeffer shows that similar considerations are otherwise of great interest for the random generation of combinatorial structures.)

We now specialize the Boltzmann model to partitions, with the Boltzmann models  $Q_x$ ,  $Q_x^s$ , and the fixed-size models  $Q_n$ ,  $Q_n^s$  taken in association to the combinatorial classes  $\mathcal{P}$ ,  $\mathcal{P}_n^s$ . The generating functions  $P(z)$ ,  $P^s(z)$  have radius of convergence  $\rho = 1$  and both blow up exponentially as  $z \rightarrow 1^-$ . Thus, the models  $Q_x$ ,  $Q_x^s$  must have something to say on the limiting behaviours of objects in  $\mathcal{P}$ ,  $\mathcal{P}_n^s$ . As it is easy to see, the Boltzmann models  $Q_x$  and  $Q_x^s$  correspond to infinite sequences of *independent* integer valued random variables  $R_k$  ( $k = 1, 2, \dots$ ), with laws as follows:

$$(10) \quad \begin{aligned} Q_x &: & R_k \in \mathbb{Z}_{>0}, & \mathbf{P}(R_k = \ell) = x^{k\ell}(1 - x^k) \\ Q_x^s &: & R_k \in \{0, 1\}, & \mathbf{P}(R_k = 1) = x^k/(1 + x^k). \end{aligned}$$

In other words, the non-identically distributed (but independent)  $R_k$  are Bernoulli in the case of  $Q_x^s$  and geometric in the case of  $Q_x$ .

A simple calculation based on Equation (9), on Chebyshev's inequalities, and on the usual approximation techniques for partition functions shows that a window narrowly centred around size  $N = n$  is obtained by fixing  $x = x_n$ ,  $x = x_n^s$  given by

$$x_n = 1 - \frac{\alpha}{\sqrt{n}}, \quad x_n^s = 1 - \frac{\beta}{\sqrt{n}},$$

for  $Q_x$  and  $Q_x^s$ , respectively. (Note that these values coincide with the saddle points of the complex-analytic approach in Section 1! This fact is general since the equations  $\mathbf{E}_x(N) = n$  and the saddle-point condition for (2) precisely coincide.) Large deviations of sums of Bernoulli or geometric random variables involve the entropy function. The Boltzmann models for partitions then provide a first hint as to the natural occurrence of entropy functions in the statements of Theorems 1 and 2.

#### 4. The Spirit of Complete Proofs

In this short abstract, we cannot do more than presenting a broad (and vague) outline of what the full proof of Theorems 1 and 2 requires.

First, under the continuous-parameter models  $Q_x, Q_x^s$ , it is easy to determine information on single parameters of partitions. The paper under review recovers for instance the analogues of Erdős and Lehner's estimates when  $x = x_n$  and  $x = x_n^s$ . It then proceeds by proving the LDP for these models. What is required is showing that, for any fixed  $m$ , and any fixed "instants"

$t_1, t_2, \dots, t_m$ , the random vectors  $(\tilde{\varphi}_n(t_1), \dots, \tilde{\varphi}_n(t_m), n^{-1}N)$  satisfy a large deviation principle. The proof bases itself on the independence granted by the models: one needs to estimate the probabilities of “slices” of summands in the scale of  $\sqrt{n}$  to be away from what is expected; this is largely based on the approximation of Riemann sums by integrals. As a snapshot of the latter technique, we offer the simple estimate

$$\mathbf{E}_{x_n}(\tilde{\varphi}_n(t)) = \sum_{k=tn^{1/2}} \frac{1}{\sqrt{n}} \frac{x_n^k}{1-x_n^k} \rightarrow \int_t^\infty \frac{du}{e^{\alpha u} - 1} du =: \Psi(t).$$

Last but not least, the treatment relies on an intensive use of large deviation techniques as exposed in [5].

In a second step, a Tauberian type of process needs to be applied. Indeed, the models  $Q_{x_n}$  are a sort of weighted average of various models of a size  $N$ , which is only controlled to lie in the vicinity of  $n$  but still fluctuates randomly. However, results at  $N = n$  exactly are wanted. Contour integration is one common way of achieving this, but the authors of [4] opt for a more combinatorial path. One of the ideas is to appeal to the following *area transformation*: given a diagram  $\lambda$  of area  $N$  at most  $n$ , form a new diagram of area  $n$  exactly, by completing the last row of  $\lambda$  by  $n - N$  elements. This establishes a mapping from  $\bigcup_{N=1}^n \mathcal{P}_N$  to  $\mathcal{P}_n$  that does not affect shape and various other characteristics of partitions too much. In this way, large deviation properties established for values of  $N$  slightly smaller than  $n$  (as given by the family of  $Q_{x_n}$  models) can be “transferred” to partitions of exact size  $n$ , that is, to the model  $Q_n$ .

The paper under discussion concludes by noting that several such large deviation principles should hold for various types of partitions with multiplicities and constrained partitions, as well as labelled trees and set partitions. In the last case, the objects at stake are enumerated by exponential generating functions, and suitable adaptations of the Boltzmann models (with the Poisson distribution replacing the geometric or Bernoulli distribution) are lurking in the background.

### Bibliography

- [1] Aldous (David) and Diaconis (Persi). – Longest increasing subsequences: from patience sorting to the Baik-Deift-Johansson theorem. *Bulletin of the American Mathematical Society (New Series)*, vol. 36, n° 4, 1999, pp. 413–432.
- [2] Andrews (George E.). – *The theory of partitions*. – Addison-Wesley Publishing Co., Reading, Mass.-London-Amsterdam, 1976, xiv+255p. *Encyclopedia of Mathematics and its Applications*, Vol. 2.
- [3] Baik (Jinho), Deift (Percy), and Johansson (Kurt). – On the distribution of the length of the longest increasing subsequence of random permutations. *Journal of the American Mathematical Society*, vol. 12, n° 4, 1999, pp. 1119–1178.
- [4] Dembo (A.), Vershik (A.), and Zeitouni (O.). – Large deviations for integer partitions. *Markov Processes and Related Fields*, vol. 6, n° 2, 2000, pp. 147–179.
- [5] Dembo (Amir) and Zeitouni (Ofer). – *Large deviations techniques and applications*. – Jones and Bartlett Publishers, Boston, MA, 1993, xiv+346p.
- [6] Erdős (Paul) and Lehner (Joseph). – The distribution of the number of summands in the partitions of a positive integer. *Duke Mathematical Journal*, vol. 8, 1941, pp. 335–345.
- [7] Hardy (G. H.). – *Ramanujan: Twelve Lectures on Subjects Suggested by his Life and Work*. – Chelsea Publishing Company, New-York, 1978, third edition. Reprinted and Corrected from the First Edition, Cambridge, 1940.
- [8] Logan (B. F.) and Shepp (L. A.). – A variational problem for random Young tableaux. *Advances in Mathematics*, vol. 26, n° 2, 1977, pp. 206–222.
- [9] Vershik (A. M.). – Statistical mechanics of combinatorial partitions, and their limit configurations. *Funktsional’nyiĭ Analiz i ego Prilozheniya*, vol. 30, n° 2, 1996, pp. 19–39, 96.
- [10] Vershik (A. M.) and Kerov (S. V.). – Asymptotics of the Plancherel measure of the symmetric group and the limiting form of Young tables. *Soviet Mathematics Doklady*, vol. 18, 1977, pp. 527–531.

## Random Walks and Heaps of Cycles

*Philippe Marchal*

Département de mathématiques et applications, École normale supérieure (France)

April 23, 2001

*Summary by Cyril Banderier*

### Abstract

The problem addressed here is the covering time of random walks on a graph satisfying “self-avoiding” properties. Appealing to the combinatorics of heaps of cycles, the author derives explicit expressions of the laws for several algorithms related to loop-erased random walks (and thus to spanning trees and Hamiltonian cycles samplings), Łukasiewicz walks, and taboo random walks.

### 1. Spanning Trees, Hamiltonian Cycles, Spanning Heaps of Cycles

Combinatorial tools (such as generating functions, context-free grammars) generally have too little “memory” to deal with “self-avoiding” walks, and thus their enumeration remains a widely open problem. However, for a few years, an approach via loop-erased random walks has seemed promising (see [3] and also the summary of R. Kenyon’s talk in the proceedings of years 99–00). Philippe Marchal exploits here the theory of determinants related to properties of heaps of cycles<sup>1</sup> and then gives the average time needed to generate self-avoiding walks of several kinds.

Define a *cycle* as a path beginning and ending at the same point, and not containing any subcycle. Given a connected graph  $G$  (where each edge is oriented and weighted), one wants to find

- a spanning tree  $T$  of this graph (i.e., a tree  $T$  whose each edge is an edge from  $G$  and each vertex of  $G$  a node of  $T$ );
- a Hamiltonian cycle  $C$  (i.e., a cycle  $C$  whose each edge is an edge from  $G$ , and each vertex of  $G$  is visited exactly one time by  $C$ );
- a spanning heap  $H$  of cycles (i.e., a heap  $H$  of cycles whose each edge of is an edge from  $G$ , and each vertex of  $G$  is visited by at least one of the cycles of  $H$ ).

In order to get a spanning tree or a Hamiltonian cycle of the graph, it is interesting to use probabilistic algorithms, since these problems are NP-complete.

On the connected edge-weighted oriented graph  $G$  (the weights are given by a matrix  $P$ ), one considers the Markov chain  $(X_n)_{n \in \mathbb{N}}$ , defined by

$$\mathbf{P}(X_{n+1} = i \mid X_n = j) = P_{ij}.$$

This means that the probability to go from vertex  $i$  (where you are at time  $n$ ) to vertex  $j$  is the weight  $P_{ij}$  of edge  $(i, j)$ . In this talk, one considers irreducible Markov chains only (i.e., the random walk visits each of the  $m$  vertices of the graph  $G$  an infinite number of times with probability 1) so that there always exists a vertex-stationary distribution  $(\pi_1, \dots, \pi_m)$ , where  $\pi_j$  is the probability

---

<sup>1</sup>The cycle decomposition was introduced by Cartier and Foata [2] and the modelling via heaps of cycles is due to Viennot [8].

to be at vertex  $j$ , after a long enough time. Similarly, there is an edge-stationary distribution for the edges.

Define the weight of the tree  $T$  (resp. the cycle  $C$ , the heap  $H$ ) as the product of the weights of its edges. Consider now a trajectory (a realization) of the Markov chain  $(X_n)_{n \in \mathbb{N}}$ , and whenever one performs a cycle, one erases this cycle from the walk and one puts this cycle on a heap (by construction, this cycle has no subcycle). If one stops at time  $n$ , one gets a “loop-erased random walk” (which is a self-avoiding walk of length less than or equal to  $n$ ) and a heap of cycles. It will be explained in Section 3 how to use this loop-erased random walk to get a spanning tree, a Hamiltonian cycle or a spanning heap of cycles.

## 2. Generating Functions

Let  $N_{ij}$  be the number of visits through the edge  $(i, j)$  and  $t_{ij}$  a formal variables associated to the edge  $(i, j)$ . Note that  $N_{ij}$  takes also into account the visits in the cycles that get erased, thus  $\sum N_{ij} = n$  is the length of the walk. Then, define the formal weight function  $\tilde{w}$  as the function which transforms a path (i.e., a sequence of edges)  $\gamma = ((x_0, x_1), \dots, (x_{n-1}, x_n))$  into the polynomial

$$\tilde{w}(\gamma) = \prod_{i=1}^n P_{x_{i-1}x_i} t_{x_{i-1}x_i}.$$

This definition (as a product of the formal weights of each edge) is easily extended to trees, cycles, graphs. Define now the formal transition matrix  $\tilde{P}$  by  $\tilde{P}_{ij} = P_{ij} t_{ij}$  and, for a subset  $S$  of the edges of the graph  $G$ , define  $\tilde{P}_S$  as equal to  $\tilde{P}$  excepted that  $(\tilde{P}_S)_{ij} := 0$  whenever  $i \notin S$  or  $j \notin S$ .

Let  $\mathcal{C}$  be the set of cycles and  $\mathcal{H}$  the set of heaps of cycles from  $\mathcal{C}$ , then

$$\sum_{H \in \mathcal{H}} \tilde{w}(H) = \left( \sum_{k \geq 1} \sum_{C_1 \dots C_k \in \mathcal{C}} (-1)^k \tilde{w}(C_1) \dots \tilde{w}(C_k) \right)^{-1} = \frac{1}{\det(\text{Id} - \tilde{P})}$$

where  $C_1, \dots, C_k$  are disjoint cycles belonging to  $\mathcal{C}$ . The proof comes from an expansion of the determinant as a sum over all permutations and then decompose each permutation in a product of cycles (each  $(-1)^k$  is nothing but an avatar of the signature of each permutation).

If  $\mathcal{H}$  stands for the set of heap of cycles avoiding a subset  $S$  of the edges of the graph  $G$ , one has

$$\sum_{H \in \mathcal{H}} \tilde{w}(H) = \frac{1}{\det(\text{Id} - \tilde{P}_S)}.$$

Whereas if  $\mathcal{H}$  stands for the set of heaps of cycles intersecting a set  $S$ , one has

$$\sum_{H \in \mathcal{H}} \tilde{w}(H) = \frac{\det(\text{Id} - \tilde{P}_S)}{\det(\text{Id} - \tilde{P})}.$$

For example, if one stops the random walk  $X$  as soon as it reaches a given point  $v$  and one considers the associated loop-erased walk  $\gamma$ , one has the following probability generating function

$$\mathbf{E} \left( \prod_{(i,j)} t_{ij}^{N_{ij}}, \gamma \right) = \frac{\tilde{w}(\gamma)}{\det(\text{Id} - \tilde{P}_v)}.$$

The right member has to be read as a generating function in several variables (the number of edges in  $G$ ) whose coefficient, e.g.,  $[t_{1,2}^4 \dots t_{1,4}^0 \dots t_{3,5}^7] = \frac{\tilde{w}(\gamma)}{\det(\text{Id} - \tilde{P}_v)}$ , gives the probability that the random walk  $X$  visits edge  $(1, 2)$  4 times, edge  $(1, 4)$  0 time, edge  $(3, 5)$  7 times, ... while the associated loop-erased random walk finally gives  $\gamma$ .

Remark: For queuing theory, assurance, etc., a usual model is left-continuous random walks (walks on  $\mathbb{Z}$  with a finite set of jumps where the only negative jump is  $-1$ ). These walks are sometimes called Łukasiewicz walks, due to their correspondence with simple families of trees, their nice combinatorial and analytic properties are well understood, see [1]. Let  $p_i$ ,  $i \geq -1$ , be the probability to do a jump  $i$  and let  $P_n$  be the transition matrix restricted to  $[0, n]$ . Then  $D_n(t) := \det(\text{Id} - tP_n)$  can easily be computed by the following recurrence:

$$D_0(t) = D_{-1}(t) = 1, \quad D_k(t) = D_{k-1}(t) - \sum_{n=0}^k p_n t (p_{-1}t)^n D_{k-n-1}(t).$$

### 3. Wilson's Algorithm and Some Variants

By convention, one considers spanning trees whose edges are all oriented to the root. Let  $\mathcal{T}_r$  be the set of spanning trees rooted at  $r$ ; a well-known result, the matrix-tree theorem implies that

$$\frac{\sum_{T \in \mathcal{T}_r} w(T)}{\pi_r} = \text{constant}.$$

The striking fact is that the quotient does not depend on  $r$ .

Wilson's algorithm [7] allows to construct a random spanning tree with a given root  $r$ . Specify an arbitrary order on  $G$ . Start the loop-erased random walk from the first point (with respect to the above order) until it reaches  $r$ . It gives a self-avoiding walk  $T_1$ . Then, restart from the first remaining point until one reaches  $T_1$ , one got a subtree  $T_2$ , etc. Finally, one gets a random spanning tree, rooted at  $r$ .

The probability to get this tree  $T$  is proportional to its weight  $w(T)$  and does not depend on the chosen order. The proof relies on the correspondence between trajectories and heap of cycles as explained above. The probability generating function is

$$\mathbf{E} \left( \prod_{(i,j) \in G^2} t_{ij}^{N_{ij}}, T \right) = \frac{\tilde{w}(T)}{\det(\text{Id} - \tilde{P}_r)}$$

and thus the average time is  $\text{tr}((\text{Id} - P_r)^{-1})$ .

Similarly, one can get a Hamiltonian cycle. Start the loop-erased random walk from a point  $r \in G$  and stop the first time one gets a Hamiltonian cycle  $C$  in the heap of cycles. Let  $\mathcal{C}$  be the set of Hamiltonian cycles. Then the probability generating function is independent from  $r$ :<sup>2</sup>

$$\mathbf{E} \left( \prod_{(i,j) \in G^2} t_{ij}^{N_{ij}}, C \right) = \frac{\tilde{w}(C)}{\det(\text{Id} - \tilde{P}) + \sum_{C' \in \mathcal{C}} \tilde{w}(C')}.$$

Finally, one gets also a sampling algorithm for a spanning heap of cycles. Choose an arbitrary order on  $G$ . Start the loop-erased random walk from  $a_1$ , stop when it returns to  $a_1$ . Then consider the first remaining non-visited point  $a_2$  and start a loop-erased random walk from  $a_2$  and stop when it returns to  $a_2$ , etc. Stop when all the points have been visited. Here again, the occupation measure does not depend on the chosen order. The proof relies on the fact that one gets a minimal

---

<sup>2</sup>Consider a nearest neighbor random walk on a cyclic graph with  $m$  vertices, and stop the walk when it comes back to the starting point, after having covered all the graph. Then, the occupation measure does not depend on the starting point. This phenomenon was observed by Pitman in 1996 for Brownian motion.

spanning heap. The probability generating function is

$$\mathbf{E} \left( \prod_{(i,j) \in G^2} t_{ij}^{N_{ij}}, H \right) = \det(\text{Id} - \tilde{P}) \sum_{F \subset G} \frac{(-1)^{|F|}}{\det(\text{Id} - \tilde{P}_F)}.$$

The waiting time  $W$  of the algorithm is stochastically less than the first time  $W_v$  that the walk returns to vertex  $v$ , *after* having visited all the vertices of the graph:

$$\forall v \in G, \forall n \in \mathbb{N} \quad \mathbf{P}(W \leq n) \geq \mathbf{P}(W_v \leq n)$$

The proof follows from the fact that any spanning pyramid (see [6]) contains a minimal spanning tiling. The author also derives this inequality:

$$\frac{1}{\inf_{v \in G} \pi_v} \leq \mathbf{E}(W) \leq \sum_{v \in G} \frac{1}{\pi_v}.$$

#### 4. Killed Random Walks

Let  $q \in (0, 1)$ , to kill  $X$  with a probability  $1 - q$  means to add a sink  $s$  and to put some probabilities of transition  $P'_{ij} = qP_{ij}$ ,  $P'_{is} = 1 - q$ . Then, if one runs Wilson's algorithm (rooted at  $s$ ), one gets a random heap with a probability proportional to  $\tilde{w}(H)q^{|H|}$  where  $|H|$  is the number of edges in  $H$ . The following process also provides a random heap (with the same distribution): construct an infinite random heap and then color each edge in red with probability  $q$ . Drop the red cycles. Then one gets a red heap with the wanted probability and another heap whose all minimal cycles have at least one non-colored edge. Let  $q$  vary continuously and thus obtain an increasing family of heaps. At a given value  $q$ , an upside-down pyramid falls. The probability that an upside-down pyramid  $P$  falls between  $q$  and  $q + dq$  equals

$$\tilde{w}(P)q^{|P|-1}dq.$$

Some generalisations of this idea allow to generate walks constrained to avoid a specified set, known as taboo random walks.

This summary is related to Marchal's articles [4, 5, 6]. The readers who want to learn more about "Perfectly Random Sampling with Markov Chains" can have a look at the web site maintained by David Wilson at <http://dimacs.rutgers.edu/~dbwilson/exact/>.

#### Bibliography

- [1] Banderier (Cyril). – *Combinatoire analytique des chemins et des cartes*. – Thèse universitaire, Université de Paris VI, 2001.
- [2] Cartier (P.) and Foata (D.). – *Problèmes combinatoires de commutation et réarrangements*. – Springer-Verlag, Berlin, 1969, iv+88p.
- [3] Lawler (Gregory F.). – Loop-erased random walk. In *Perplexing problems in probability*, pp. 197–217. – Birkhäuser Boston, Boston, MA, 1999.
- [4] Marchal (Philippe). – Cycles hamiltoniens aléatoires et mesures d'occupation invariantes par une action de groupe. *Comptes Rendus de l'Académie des Sciences. Série I. Mathématique*, vol. 329, n° 10, 1999, pp. 883–886.
- [5] Marchal (Philippe). – Loop-erased random walks, spanning trees and Hamiltonian cycles. *Electronic Communications in Probability*, vol. 5, 2000, pp. 39–50.
- [6] Marchal (Philippe). – *Loop-erased random walks and heaps of cycles*. – Prépublications n° DMA-01-07, École normale supérieure, Paris, 2001.
- [7] Propp (James Gary) and Wilson (David Bruce). – How to get a perfectly random sample from a generic Markov chain and generate a random spanning tree of a directed graph. *Journal of Algorithms*, vol. 27, n° 2, 1998, pp. 170–217. – 7th Annual ACM-SIAM Symposium on Discrete Algorithms (Atlanta, GA, 1996).
- [8] Viennot (Gérard Xavier). – Heaps of pieces. I. Basic definitions and combinatorial lemmas. In *Combinatoire énumérative (Montréal, Québec, 1985)*, pp. 321–350. – Springer, Berlin, 1986.



## Tail Bounds for Occupancy Problems

*Paul Spirakis*

Computer Technology Institute, Patras University (Greece)

November 20, 2000

*Summary by Stéphane Boucheron*

### Abstract

The talk was based on [9] and consisted in a presentation of various tail bounds for occupancy problems and applications to the determination of the conjectured satisfiability threshold in the random  $k$ -sat problem.

### 1. Bins and Balls and Occupancy Problems

In bins and balls games,  $m$  balls are placed independently and uniformly at random among  $n$  bins. Henceforth, a generic allocation will be denoted by  $\omega \in \{1, \dots, n\}^m$ :  $\omega_k = j$  if the  $k$ -th ball is located in the  $j$ -th bin. Let  $X_n(\omega, m)$  denote the number of empty bins when  $m$  balls have been assigned a position. The piecewise constant interpolation is defined by  $X_n(\omega, t) = X_n(\omega, \lceil tn \rceil)$ . To alleviate notations, we omit  $\omega$  when this is not a source of confusion. The behavior of the process  $X_n(\cdot)$  as  $n$  becomes large has been the subject of many investigations in random combinatorics. The lecture is concerned with different derivations of tail bounds for  $X_n(\cdot)$  and their application to the analysis of the threshold phenomenon for the (random)  $k$ -satisfiability problem.

**1.1. Approaches to random allocations.** There are many approaches to random allocation problems. Many early successes of analytic combinatorics have been reported in the monograph by Kolchin, Sevast'yanov and Chystiakov [11].

Probabilistic (Martingale-theoretical) approaches have been successful as well. Let  $\mathcal{F}_t$  denote the  $\sigma$ -algebra generated by the first  $\lfloor nt \rfloor$  allocations (we do not mention  $n$  to alleviate notations). Then it is straightforward to check the relation

$$\mathbf{E} \left[ X_n \left( t + \frac{1}{n} \right) \middle| \mathcal{F}_t \right] = \left( 1 - \frac{1}{n} \right) X_n(t).$$

From this, one immediately deduces that  $(1 - \frac{1}{n})^{-\lfloor nt \rfloor} X_n(t)$  is an  $\mathcal{F}_t$ -Martingale. Moreover it has bounded increments, and its quadratic variation process converges in probability towards  $t \mapsto e^t - (1 + t)$ . Applying Martingale limit theorems [8], one easily deduces:

- a law of large numbers:  $X_n(\cdot)/n$  converges in probability towards  $t \mapsto e^{-t}$ ,
- a functional central limit theorem:  $t \mapsto (X_n(t) - ne^{-t})/\sqrt{n}$  converges towards a rescaled time-changed Brownian motion, namely  $t \mapsto e^{-t}B[e^t - (1 + t)]$ .

Unfortunately, results on convergence in distribution tell little about asymptotic probability of rare events: the convergence rate cannot be better than  $O(1/\sqrt{n})$ , and probability of rare events are especially relevant to the analysis of extreme values that often constitute the core of applications.

Nevertheless, central limit theorems suggest that the tail probabilities of the empty cell statistics might be Gaussian-like. In computer science, sharp upper bounds on tail probabilities are often desirable.

If instead of throwing a fixed number  $\lfloor nt \rfloor$  of balls into the  $n$  bins, one first draws  $N$  according to a Poisson distribution with parameter  $\lfloor nt \rfloor$ , and then throws  $N$  balls into the  $n$  bins, the bin occupancies become independent Bernoulli random variables with success probability  $\approx \exp(-t)$ .  $X_n(t)$  is now distributed according to a binomial random variable with parameters  $n$  and  $\exp(-t)$ . Let  $\mathbb{P}$  denote the original probability distribution on allocations and let  $\mathbb{Q}$  denote this alternate probability distribution on  $N$  and allocations. Note that conditionally on  $N = \lfloor nt \rfloor$ , the distributions of  $X_n(t)$  under  $\mathbb{P}$  and  $\mathbb{Q}$  are identical (the multinomial distribution is a conditioned Poisson process). Then

$$(1) \quad \mathbb{P}\{X_n(t) \in A\} = \frac{\mathbb{Q}\{X_n(t) \in A \wedge N = \lfloor nt \rfloor\}}{\mathbb{Q}\{N = \lfloor nt \rfloor\}} \leq \sqrt{2\pi nt} \mathbb{Q}\{X_n(t) \in A\}$$

Inequality (1) provides with an easy tail upper bound for rare events under  $\mathbb{Q}$ , i.e., for large deviations of  $X_n(t)$  around its expectation. If  $A = \{\omega \mid X_n(\omega, t) > ne^{-t} + n\epsilon\}$ , then

$$\mathbb{P}\{X_n(t) \in A\} \leq \sqrt{2\pi n} \exp\left(-nh(e^{-t} + \epsilon, e^{-t})\right)$$

where  $h(x, y) = x \log \frac{x}{y} + (1-x) \log \frac{1-x}{1-y}$ . It obviously raises two questions: Is the order of the exponent correct? Can we get rid of the  $\sqrt{n}$  factor?

**1.2. Known results.** As allocation are performed independently, a very straightforward yet useful bound comes from the Azuma–Mc Diarmid inequality. Namely note that if  $\omega$  and  $\omega'$  are two allocation schemes that differ only in one position  $\omega_j = \omega'_j$  for all  $j \leq k = \lfloor tn \rfloor$  except for  $j = i$ , then  $|X_n(\omega, t) - X_n(\omega', t)| \leq \frac{1}{n}$ . As a matter of fact, if the space of allocations is equipped with the Hamming distance, the empty bin statistics is 1-Lipschitz. This implies that

$$(2) \quad \mathbb{P}\left\{|X_n(t) - \mathbf{E}[X_n(t)]| > n\epsilon\right\} \leq 2 \exp\left(-\frac{2n\epsilon^2}{t^2}\right).$$

Inequality (2) is obtained by a Martingale embedding argument. Namely  $X_n(t) = \mathbf{E}[X_n(t) \mid \mathcal{F}_t]$  and the process  $M_n(s) = \mathbf{E}[X_n(t) \mid \mathcal{F}_s]$  is an  $\mathcal{F}_s$ -martingale, as

$$\mathbf{E}[M_n(s+h) \mid \mathcal{F}_s] = \mathbf{E}\left[\mathbf{E}[X_n(t) \mid \mathcal{F}_{s+h}] \mid \mathcal{F}_s\right] = \mathbf{E}[X_n(t) \mid \mathcal{F}_s] = M_n(s).$$

One may wonder what the best way to apply Azuma's inequality is.

**1.3. Painless tail bounds.** The first bound presented in [9] is:

$$(3) \quad \mathbb{P}\left\{|X_n(t) - \mathbf{E}[X_n(t)]| > n\epsilon\right\} \leq 2 \exp\left(-\frac{(n-1/2)n^2\epsilon^2}{n^2 - \mathbf{E}[X_n(t)^2]}\right).$$

When  $n$  becomes large, the exponent on the right-hand side is equivalent to

$$-\frac{n\epsilon^2}{1 - e^{-2t}}.$$

The trivial Poisson estimates (1) clearly shows that this exponent is rather poor as soon as  $t$  becomes non-negligible. This is not a denial of the merits of Martingale approach. Indeed, this method provides nearly optimal bounds for smooth Gaussian functionals and for many discrete problems. The apparent flaw in Equation (3) comes from the fact that we did not use tight enough bounds on the quadratic variation process associated with  $\mathbf{E}[X_n(t) \mid \mathcal{F}_s]$ .

Next the authors of [9] proceed to establish what they call a Chernof bound for the occupancy problem. It shows that the Poisson tail estimate (1) is correct even if we do not resort to a conditioning argument, i.e., that the  $\sqrt{n}$  factor is spurious.

### 2. The Large Deviation Approach

The large deviation approach (see [2, 5, 7] for recent presentations) aims at identifying the right exponents for tail probability. It provides the right touchstone for the occupancy problem. Rather than using the martingale structure of the occupancy problem, the large deviation approach relies on the Markovian structure of the occupancy problem: conditionally on  $X_n(t)$ ,  $X_n(t + 1/n)$  does not depend on  $\mathcal{F}_{t-1/n}$ . The large deviation principle invoked in [9] comes from a contraction of a functional large deviation principle derived by Azencott and Ruget. The latter shows that asymptotically, the exponent in large deviation probabilities can be represented as a the solution of a variational problem, namely

$$(4) \quad \lim_{n \rightarrow \infty} \frac{1}{n} \log \mathbb{P} \{ X_n(t) \geq nx \} = - \inf_{\xi(0)=1, \xi(t)=x} \int_0^t h(-\dot{\xi}(s), \xi(s)) ds.$$

The article [9] solves the associated variational problem and provides a closed form for the exponent, confirming the intuition that the exponent obtained by Poissonization is not optimal.

### 3. Satisfiability Problems

The second part of the paper presents an application of tail bounds for occupancy problems to the analysis of the random 3-sat problem. An instance of the 3-sat problem is a boolean formula in conjunctive normal form, where each clause has at most 3 literals. For each number  $n$  of variables, and each problem size  $k$ , the set of instances of the 3-sat problem is provided with the uniform probability over the  $m$ -tuples of 3-clauses over the  $n$  variables. At the time of writing [9], it was conjectured that as  $n$  goes to infinity while  $k/n$  remains constant, a phase transition occurs. For  $k/n < c_3$ , random 3-sat formulas are satisfiable with overwhelming probability, while for  $k/n > c_3$  random 3-sat formulas are not satisfiable formulas with overwhelming probability.

The paper [9] proposes an upper-bound on the conjectured satisfiability threshold:  $c_3 \leq 4.758$ . This result came in a series of improvement starting from the straightforward  $c_3 \leq 5.19$ , through  $c_3 \leq 5.08$  [6],  $c_3 \leq 4.64$  [3],  $c_3 \leq 4.601$  [10], and recently culminating with  $c_3 \leq 4.506$  [4].

In the sequel,  $n$  and  $k$  are supposed to be fixed.  $F$  denotes a random 3-sat formula,  $\#F$  denotes the number of assignments of the  $n$  boolean variables that satisfy  $F$ .  $F$  is satisfiable if  $\#F \geq 1$ .  $T(F)$  equals 1 if  $F$  is satisfiable, 0 otherwise. Let  $\sigma$  denote a generic truth assignment.  $F(\sigma)$  equals 1 if  $\sigma$  satisfies  $F$ , 0 otherwise.  $\mathbf{1}$  denotes the truth assignment where all variables are set to 1. Then, we have

$$(5) \quad \mathbf{E}_F [T(F)] = \mathbf{E}_F \left[ \sum_{\sigma: F(\sigma)=1} \frac{1}{\#F} \right] = \sum_{\sigma} \mathbf{E}_F \left[ \frac{F(\sigma)}{\#F} \right] = 2^n \mathbf{E}_F \left[ \frac{F(\mathbf{1})}{\#F} \right],$$

where the second equality comes from the fact that the number of formulae that satisfy a particular truth assignment does not depend on the truth assignment. Hence, to get an upper bound on the probability of satisfiability, it is enough to get an upper bound on

$$\left( \frac{7}{8} \right)^{cn} \mathbf{E}_{F'} \left[ \frac{1}{\#F} \right],$$

where  $F$  is now picked at random among the  $\left(\frac{7}{8}\right)^{cn} \binom{n}{3}^{cn}$  formulae that are satisfied by  $\mathbf{1}$ . This distribution among formulae is a product distribution where each clause is picked uniformly at random among the clauses where at least one literal is not negated.

The main idea of the proof is to establish that conditionally on the fact that it is satisfiable, a 3-sat formula with sufficiently many clauses has exponentially many satisfying truth assignments with overwhelming probability.

What is proved in [9] is actually the following. Let  $\#F_1$  denote the number of truth assignments  $\sigma$  of  $F$  where for each clause in  $F$ , there exists a non-negated variable that evaluates to 1 in  $\sigma$ . Obviously  $1/\#F \leq 1/\#F_1$ . Now to lower bound  $\#F_1$ , it is enough to determine a minimum family of variables  $\mathcal{I}(F)$  such that any truth assignment where all variables in  $\mathcal{I}(F)$  evaluates to 1 satisfies the formula  $F$  ( $\mathcal{I}(F)$  is sometimes called a prime implicant of  $F$ ). As a matter of fact, we have  $\#F_1 \geq 2^{n-\#\mathcal{I}}$ , and hence

$$(6) \quad \mathbb{P}\{F \text{ is satisfiable}\} \leq \left(\frac{7}{8}\right)^{cn} \mathbf{E}_{F'} \left[ 2^{\#\mathcal{I}} \right].$$

Since the publication of [9], improved upper bounds on  $c_3$  have been derived by refining estimations on the fluctuations of  $\#F$  for random formulae. Those estimations still rely on statistics for random allocations. But the empty bins statistics are no more sufficient. The best known upper bounds [4] rely on a statistics that have sometimes been called empirical occupancy measures. As a matter of fact, an allocation  $\omega$  defines a probability measure on  $\mathbb{N}$ ,  $\bar{X}_n(i, t)$  denotes the fraction of bins that contain  $i$  balls for  $i \in \mathbb{N}$ . The large deviations of this measure-valued random variable may be studied in different ways: by resorting to Azencott–Ruget results and projective limit arguments [2], or directly as in [1].

### Bibliography

- [1] Boucheron (S.), Gamboa (F.), and Léonard (C.). – *Bins and balls: large deviations of the empirical occupancy process*. – Rapport de recherche du LRI n° 1255, Université Paris-Sud, 2000.
- [2] Dembo (Amir) and Zeitouni (Ofar). – *Large deviations techniques and applications*. – Springer-Verlag, New York, 1998, second edition, xvi+396p.
- [3] Dubois (O.) and Boufkhad (Y.). – A general upper bound for the satisfiability threshold of random  $r$ -SAT formulae. *Journal of Algorithms*, vol. 24, n° 2, 1997, pp. 395–420.
- [4] Dubois (O.), Boufkhad (Y.), and Mandler (J.). – Typical random 3-sat formulae and the satisfiability threshold. In *Proceedings of SODA'2000*. ACM, pp. 126–127. – 2000.
- [5] Dupuis (Paul) and Ellis (Richard S.). – *A weak convergence approach to the theory of large deviations*. – John Wiley & Sons, New York, 1997, xviii+479p. A Wiley-Interscience Publication.
- [6] El Maftouhi (A.) and Fernandez de la Vega (W.). – On random 3-sat. *Combinatorics, Probability and Computing*, vol. 4, n° 3, 1995, pp. 189–195.
- [7] Feng (Jin) and Kurtz (Thomas G.). – *Large deviations for stochastic processes*. – 2000. 194 pages. Available from <http://www.math.wisc.edu/~kurtz/feng/ldp.htm>.
- [8] Hall (P.) and Heyde (C. C.). – *Martingale limit theory and its application*. – Academic Press, New York, 1980, xii+308p. Probability and Mathematical Statistics.
- [9] Kamath (Anil), Motwani (Rajeev), Palem (Krishna), and Spirakis (Paul). – Tail bounds for occupancy and the satisfiability threshold conjecture. *Random Structures & Algorithms*, vol. 7, n° 1, 1995, pp. 59–80.
- [10] Kirousis (Lefteris M.), Kranakis (Evangelos), Krizanc (Danny), and Stamatiou (Yannis C.). – Approximating the unsatisfiability threshold of random formulas. *Random Structures & Algorithms*, vol. 12, n° 3, 1998, pp. 253–269.
- [11] Kolchin (Valentin F.), Sevast'yanov (Boris A.), and Chistyakov (Vladimir P.). – *Random allocations*. – V. H. Winston & Sons, Washington, D.C., 1978, xi+262p. Translated from the Russian.

## Patricia Tries in the Context of Dynamical Systems

Jérémie Bourdon

GREYC, Université de Caen (France)

March 19, 2001

Summary by Michel Nguyen-Thé

### Abstract

Tries, a generalized form of digital trees, are a data structure widely used in numerous domains: algorithms for searching words, compression, dynamical hashing, ... Their interest and construction lie in the partitioning of a set of words. We present a compact form of tries, called Patricia tries, in which all unary nodes are suppressed (and thus do not intervene in the partitioning). We then study the means of the memory occupation and of the cost of inserting a word for that data structure when words are produced by a probabilistic source for which the dependencies between the emitted symbols can be very important.

### 1. Size and Path Length of Tries and Patricia Tries: Expressions for Expectations

We define the notions of tries and Patricia tries. We find general expressions for the expectations of the size and path length of tries and Patricia tries in the Bernoulli model, valid for any source.

**1.1. Operations on infinite words.** For a finite alphabet  $\Sigma = \{a_1, a_2, \dots, a_r\}$ , let  $\Sigma^\infty$  be the set of infinite words on that alphabet,  $\underline{\sigma} : \Sigma^\infty \rightarrow \Sigma^\infty$  the map that returns the first letter of a word, and  $\underline{T} : \Sigma^\infty \rightarrow \Sigma^\infty$  the shift that returns the first suffix of a word. Let  $\underline{T}_{[a]}$  denote the restriction of  $\underline{T}$  to the set  $\sigma^{-1}(\{a\})$  of words beginning with symbol  $a$  and, for a finite prefix  $w = a_1 \dots a_k$ , let  $\underline{T}_{[w]}$  denote the composition  $\underline{T}_{[a_k]} \circ \underline{T}_{[a_{k-1}]} \circ \dots \circ \underline{T}_{[a_1]}$ . The notations  $\sigma$  and  $T$  are kept for operators acting on reals which will be used later.

### 1.2. Tries.

**Definition 1.** Let  $X$  be a finite set of infinite words produced by the same source. A *trie*  $\text{Tr}(X)$  is a structure defined by the following rules:

( $R_0$ ) If  $X = \emptyset$  (the empty set),  $\text{Tr}(X)$  is the empty tree.

( $R_1$ ) If  $X = \{x\}$ ,  $\text{Tr}(X)$  consists of a single leaf node represented by  $\square$  that contains  $x$ .

( $R_2$ ) If  $X$  is of cardinality greater than or equal to 2,  $\text{Tr}(X)$  is an *internal node* represented by  $\bullet$  to which are attached  $r$  subtrees:

$$\text{Tr}(X) = \left\langle \bullet, \text{Tr}(\underline{T}_{[a_1]}X), \text{Tr}(\underline{T}_{[a_2]}X), \dots, \text{Tr}(\underline{T}_{[a_r]}X) \right\rangle.$$

The edge that attaches the subtree  $\text{Tr}(\underline{T}_{[a_j]}X)$  is labelled by the symbol  $a_j$ . Notice a little abuse in ( $R_2$ ): if there is no word in  $X$  beginning with  $a_j$ , then  $\underline{T}_{[a_j]}X$  is not defined, and we consider that is equal to the empty set. Hence  $\text{Tr}(\underline{T}_{[a_j]}X)$  is the empty tree, and it is as though there were no subtree corresponding to  $a_j$  (see Figure 1).

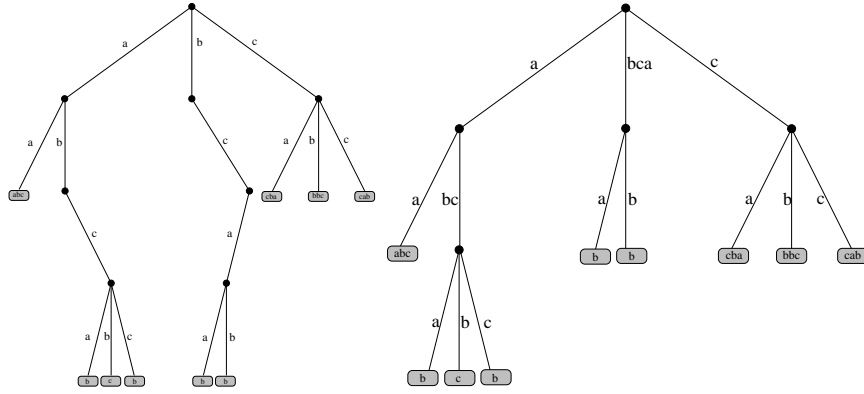


FIGURE 1. Standard trie and corresponding Patricia trie.

**1.3. Patricia Tries.** A Patricia trie is a trie from which all unary nodes are eliminated. Hence with any finite set  $X$  of infinite words produced by the same source, we associate a Patricia trie  $\text{PaTr}(X)$ . The first two rules are the same, but the last rule ( $R'_2$ ) is more sophisticated:

( $R'_2$ ) If  $X$  is of cardinality greater than or equal to 2, we have two cases:

( $R'_{2,1}$ ) if  $\underline{\sigma}(X)$  consists of a single symbol, then  $\text{PaTr}(X)$  equals  $\text{PaTr}(\underline{T}X)$ .

( $R'_{2,2}$ ) if  $\underline{\sigma}(X)$  has at least two distinct symbols,  $\text{PaTr}(X)$  is an *internal node* generically represented by  $\bullet$  to which are attached  $r$  subtrees,

$$\text{PaTr}(X) = \left\langle \bullet, \text{PaTr}(\underline{T}_{[a_1]}X), \text{PaTr}(\underline{T}_{[a_2]}X), \dots, \text{PaTr}(\underline{T}_{[a_r]}X) \right\rangle.$$

The edges of the Patricia trie are labelled by words. These words are obtained from the associated trie by concatenating all the labels of the collapsed edges.

**1.4. Additive parameters.** The *depth* of a node in a tree is the number of edges of the path that connects it to the root. The *size* of a tree is the number of its internal nodes. The *path length* of a tree is the sum of the depths of all (nonempty) external nodes.

**1.5. Algebraic analysis of additive parameters.** In a standard trie built on the set  $X = \{x_1, \dots, x_n\}$ , the structure of a node labelled by a prefix  $w$  is a finite string called a *slice* given by

$$\underline{\sigma} \underline{T}_{[w]}X := \left( \underline{\sigma} \underline{T}_{[w]}(x_1), \dots, \underline{\sigma} \underline{T}_{[w]}(x_n) \right).$$

An additive parameter  $\gamma$  on  $X$  is defined by a toll parameter  $\delta$  defined on finite strings and the recursive rule:

$$\gamma[X] = \begin{cases} 0, & \text{if } |X| \leq 1, \\ \delta[\underline{\sigma}(X)] + \sum_{m \in \Sigma} \gamma[\underline{T}_{[m]}X], & \text{if } |X| \geq 2, \end{cases}$$

Let  $|s|$  and  $\#(s)$  denote the number of symbols of the string  $s$  and the number of distinct symbols of  $s$ , respectively. The parameters of interest are the size on tries and Patricia tries,

$$\delta_S(s) = \begin{cases} 1 & \text{if } |s| \geq 2, \\ 0 & \text{otherwise,} \end{cases} \quad \delta_{PS}(s) = \begin{cases} 1 & \text{if } \#(s) \geq 2, \\ 0 & \text{otherwise,} \end{cases}$$

and the internal path length on tries and Patricia tries

$$\delta_L(s) = \begin{cases} |s| & \text{if } |s| \geq 2, \\ 0 & \text{otherwise,} \end{cases} \quad \delta_{PL}(s) = \begin{cases} |s| & \text{if } \#(s) \geq 2, \\ 0 & \text{otherwise.} \end{cases}$$

Size of Tr	$\widehat{S}(n) = \sum_{w \in \Sigma^*} (1 - (1 + (n-1)p_w)(1-p_w)^{n-1})$
Path Length of Tr	$\widehat{L}(n) = \sum_{w \in \Sigma^*} np_w(1 - (1-p_w)^{n-1})$
Size of PaTr	$\widehat{S}_P(n) = \sum_{w \in \Sigma^*} \left( 1 - (1-p_w)^n - \sum_{i \in \Sigma} \left( (1-p_w(1-p_{[i w]}))^n - (1-p_w)^n \right) \right)$
Path Length of PaTr	$\widehat{L}_P(n) = \sum_{w \in \Sigma^*} np_w \left( 1 - (1-p_w)^{n-1} - \sum_{i \in \Sigma} p_{[i w]} (1-p_w(1-p_{[i w]}))^{n-1} \right)$

TABLE 1. Expectations of size and path length for tries (Tr) and Patricia tries (PaTr).

**1.6. Expectation of parameters.** Let  $(\mathcal{P}_z, \mathcal{S})$  denote the Poisson model of rate  $z$  relative to the source  $\mathcal{S}$ , and  $p_w$  the probability that a given infinite word begins with the prefix  $w$ . If the cardinality of  $X$  is a random Poisson variable of rate  $z$ , the length of the slice  $\sigma \underline{T}_{[w]} X$  is also a random Poisson variable of rate  $zp_w$ . Hence the expectation of parameter  $\gamma$  is a sum of expectations of parameter  $\delta$ ,  $\mathbf{E}[\gamma; \mathcal{P}_z, \mathcal{S}] = \sum_{w \in \Sigma^*} \mathbf{E}[\delta; \mathcal{P}_{zp_w}, B_w]$ .

The expectation of the parameter is given by  $\mathbf{E}[\delta; \mathcal{P}_z, B] = e^{-z} \frac{\partial}{\partial u} F_\delta(z, u, p_1, \dots, p_r) \Big|_{u=1}$ , where  $F_\delta(z, u, x_1, \dots, x_r) = \sum_{s \in \Sigma^*} \frac{z^{|s|}}{|s|!} u^{\delta(s)} x_1^{|s|_1} \dots x_r^{|s|_r}$ .

Using algebraic depoissonization [3], based on the equalities  $\mathbf{E}[Y; \mathcal{P}_z] = e^{-z} \sum_{n \geq 0} \mathbf{E}[Y; \mathcal{B}_n] \frac{z^n}{n!}$  and thus  $\mathbf{E}[Y; \mathcal{B}_n] = n! [z^n] e^z \mathbf{E}[Y; \mathcal{P}_z] \frac{z^n}{n!}$ , one can return to the Bernoulli model. Finally, the expectations of interest are given in Table 1.

## 2. Tools for the Asymptotics of the Expectations

**2.1. Mellin analysis and Dirichlet series.** To get asymptotics for the expressions found previously, we first note that they belong to the paradigm of harmonic sums. Their Mellin transforms are given in Table 2, where  $\Lambda(s) = \sum_{w \in \mathcal{M}^*} p_w^s$  and

$$\begin{aligned}
 \Lambda_S(s) &= - \sum_{w \in \Sigma^*} p_w^s - \sum_{w \in \Sigma^*} p_w^s \sum_{i \in \Sigma} [(1-p_{[i|w]})^s - 1] \\
 (1) \quad &= (s-1)\Lambda(s) - s \sum_{k \geq 2} \frac{(-1)^k}{k!} \left( \prod_{i=2}^{k-1} (s-i) \right) [(s-1)\Lambda^{[k]}],
 \end{aligned}$$

$$\begin{aligned}
 \Lambda_L(s) &= \sum_{w \in \Sigma^*} p_w^s \sum_{i \in \Sigma} [(1-p_{[i|w]})^{s-1} - 1] \\
 (2) \quad &= \sum_{k \geq 2} \frac{(-1)^k}{(k-1)!} \left( \prod_{i=2}^{k-1} (s-i) \right) [(s-1)\Lambda^{[k]}],
 \end{aligned}$$

with  $\Lambda^{[k]}(s) = \sum_{w \in \Sigma^*} p_w^s \sum_{i \in \Sigma} p_{[i|w]}^k$ , for  $k \geq 1$ ,

**2.2. Dynamical sources.** We have to restrict ourselves to a class of dynamical sources  $\mathcal{S}$  (see [4] for more details and [2] for its use in a study of standard tries),

- (a) a finite or denumerable alphabet  $\Sigma$ ,
- (b) a topological partition of  $\mathcal{I} := (0, 1)$  with disjoint open intervals  $\mathcal{I}_a$ , for  $a \in \Sigma$ ,
- (c) an encoding mapping  $\sigma$  which is constant and equal to  $a$  on each  $\mathcal{I}_a$ ,

Size of Tr	$S^*(s) = -\Lambda(-s)(s+1)\Gamma(s)$
Path Length of Tr	$L^*(s) = -\Lambda(-s)\Gamma(s+1)$
Size of PaTr	$S_P^*(s) = \Gamma(s)\Lambda_S(-s)$
Path Length of PaTr	$L_P^*(s) = -\Gamma(s+1)(\Lambda(-s) + \Lambda_L(-s))$

TABLE 2. Mellin transforms of expectations.

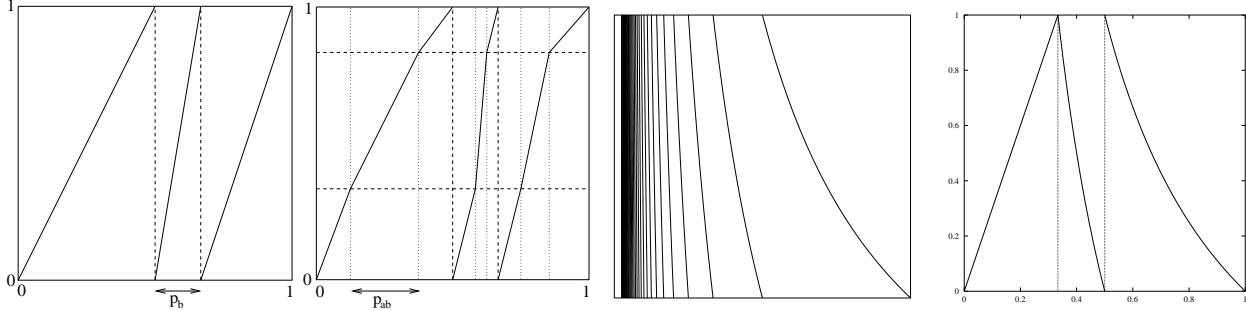


FIGURE 2. Memoryless source, Markov chain of order 1, continued fraction source, heteroclinical source.

(d) a shift mapping  $\mathcal{T}$  whose restriction to  $\mathcal{I}_a$  is a real analytic bijection from  $\mathcal{I}_a$  to  $\mathcal{I}$ .

Besides,  $\mathcal{T}$  has to satisfy more precise properties. If we let  $h_a$  be the local inverse of  $T$  restricted to  $\mathcal{I}_a$  and  $\mathcal{H}$  be the set  $\mathcal{H} = \{h_a \mid a \in \Sigma\}$ , then we add properties on bounds of the first derivatives, among which Rényi's condition which plays an important rôle in the study of conditional probabilities. This condition states that, if  $h_a$  are the local inverse of  $T$ , supposed to be locally holomorphic, restricted to  $\mathcal{I}_a$ , then there exists a constant  $K$  that bounds the ratio  $|h_a''(x)/h_a'(x)|$  for all branch  $h_a$  and all  $x \in [0, 1]$ . With each  $h_a$ , that are only defined on  $\mathcal{I}_a$ , we associate its analytical extension  $\tilde{h}_a$  to the whole set  $\mathcal{I}$ .

If  $M$  maps  $x \in [0, 1]$  to  $(\sigma(x), \sigma T(x), \sigma T^2(x), \dots) \in \Sigma^\infty$ ,  $T$ , and  $\sigma$  are linked with the previously defined  $\underline{T}$  and  $\underline{\sigma}$  by  $\underline{\sigma}M \equiv \sigma$  and  $\underline{T}M \equiv MT$ .

Figure 2 displays several types of dynamical sources:

*Memoryless sources.* We have affine branches of slope  $1/p_a$  on intervals  $\mathcal{I}_a := (q_a, q_{a+1})$ , where  $q_a = \sum_{i < a} p_i$ .

*Markov chains.* Each  $\mathcal{I}_a$  of a memoryless source is divided in  $r$  intervals  $\mathcal{I}_{a,b}$ ,  $b \in \Sigma$ , of length  $p_{ab} = p_{[b|a]} \cdot p_a$  on which  $T : \mathcal{I}_{a,b} \rightarrow \mathcal{I}_b$  has slope  $\frac{p_a}{p_{ab}} = \frac{p_b}{p_{[b|a]}} \cdot \frac{1}{p_a}$ . Notice that when the order  $d$  of the Markov chain goes to infinity in a certain sense, one obtains at the limit a source with unbounded memory.

*Continued fractions.* With  $\Sigma = \mathbb{N}$ ,  $\mathcal{I}_a := \left(\frac{1}{a+1}, \frac{1}{a}\right)$ ,  $T(x) = \frac{1}{x} - \lfloor \frac{1}{x} \rfloor$ , and  $\sigma(x) = \lfloor \frac{1}{x} \rfloor$ , corresponding to a continued fraction source, we obtain a source with unbounded memory.

*Heteroclinical sources.* A source for which derivatives in different intervals can be of different signs is called *heteroclinical*. Otherwise the source is *homoclinical*, like the sources presented before.



**2.3. Ruelle operators, multi-secants and prefix probabilities.** In the context of dynamical systems, with transformations  $T$  of local inverses  $h_a$  are associated a transfer operator,

$$\mathcal{G}[f](x) := \sum_{a \in \Sigma} |h'_a(x)| f \circ h_a(x),$$

whose interest lies in the following property: if  $X$  is a random variable with density function  $f$ , then the density of  $T(X)$  is  $\mathcal{G}[f]$ . The Ruelle operator generalizes it by introducing a complex parameter  $s$ , interpreted in statistical physics as the temperature:

$$\mathcal{G}_s[f](x) := \sum_{a \in \Sigma} \tilde{h}_a(x)^s f \circ h_a(x).$$

To deal with probabilities of prefixes of words  $p_w$  and hence with fundamental intervals, we have to replace tangents with secants  $H[h](x, y) := \left| \frac{h(x)-h(y)}{x-y} \right|$ , leading to a first generalization  $\mathbf{G}_s$  of the Ruelle operator, acting on functions  $L$  of two complex variables:

$$\mathbf{G}_s[L](x) := \sum_{a \in \Sigma} \tilde{H}_a^s[h_a](x, y) L(h_a(x), h_a(y)).$$

To deal with conditional probabilities, we have to resort to a further generalization  $\mathfrak{G}_s$  of the Ruelle operator involving multiseccants instead of secants:

$$\mathfrak{G}_s^{[m]}[L] := \sum_{a \in \Sigma} \mathfrak{H}_s^{[m]}[h_a] L \circ V[h_a],$$

where the multiseccants are defined by  $\mathfrak{H}_s^{[m]}[h](x, y, z, t) = H[h]^{s-m}(x, y)H[h]^m(z, t)$ , and  $V$  by  $V[h](x, y, z, t) = (h(x), h(y), h(z), h(t))$ .

Let  $F$  be the distribution associated with the initial density  $f$  of a source  $(\mathcal{S}, f)$ . The probability  $p_w$  that a word begins with some prefix  $w$  is  $\left| F(h_w(0)) - F(h_w(1)) \right|$ . For the special case  $F = \text{Id}$ , it will be denoted  $p_w^*$ . Let  $Q := H[F]$  be the secant of the initial distribution. Then the quasi-inverses of  $\mathbf{G}_s$  and  $\mathfrak{G}_s^{[k]}$  are related to Dirichlet series in the following way:

$$\Lambda(s) = \sum_{w \in \mathcal{M}^*} p_w^s = (\text{Id} - \mathbf{G}_s)^{-1}[Q^s](0, 1); \quad \Lambda^{[k]}(s) = \sum_{i \in \Sigma} \left( \text{Id} - \mathfrak{G}_s^{[k]} \right)^{-1} \left[ \mathfrak{H}_s^{[k]}[F] \right] (0, 1, h_i(0), h_i(1)).$$

Thanks to a theorem similar to the Perron–Frobenius theorem, we have the decomposition

$$(\text{Id} - \mathbf{G}_s)^{-1} = \frac{\lambda(s)}{1 - \lambda(s)} \mathbf{P}_s + (\text{Id} - \mathbf{N}_s)^{-1},$$

and a similar decomposition for the multi-secant operator. We deduce the asymptotics:

$$\lim_{s \rightarrow 1} (s - 1)(\text{Id} - \mathfrak{G}_s)^{-1}[L](x) = \frac{-1}{\lambda'(1)} \Psi_1(x) \int_0^1 \ell(t) dt,$$

where  $\Psi_1(x)$  is an eigenfunction associated with the dominant eigenvalue and chosen according to a proper normalization, and  $\ell$  is the diagonal mapping of  $L$ . We get similar results for the  $\Lambda^{[m]}$  that also have 1 as pole of order 1, and their respective residues  $r_m$  are related to the dominant eigenfunctions  $\Psi_1^{[m]}$  of the operators  $\mathfrak{G}_1^{[m]}$ , which allows us to find the singular expansion

$$\Lambda(s) = \Lambda^{[1]}(s) \asymp \frac{-1}{\lambda'(1)(s - 1)} + C(\mathcal{S}),$$

where  $C(\mathcal{S})$  is a constant depending on the source  $\mathcal{S}$  and the initial density  $f$ . Using the equalities (1) and (2) we can then get asymptotics for  $\Lambda_S(1)$  and  $\Lambda_L(1)$ .

Size of Tr	$S(n) \approx \frac{1}{h(\mathcal{S})}n$
Path Length of Tr	$L(n) \sim \frac{1}{h(\mathcal{S})}n \log n + \left(C(\mathcal{S}) - \frac{\gamma}{h(\mathcal{S})}\right)n$
Size of PaTr	$S_P(n) \approx \frac{1}{h(\mathcal{S})}(1 - C_1(\mathcal{S}))n$
Path Length of PaTr	$L(n) \sim \frac{1}{h(\mathcal{S})}n \log n + \left(C(\mathcal{S}) - \frac{\gamma + C_2(\mathcal{S})}{h(\mathcal{S})}\right)n$

TABLE 3. Asymptotics of expectations.

### 3. Results: Asymptotics

**3.1. General expressions.** Let  $h(\mathcal{S}) = -\lambda'(1) = \lim_{\ell \rightarrow \infty} \sum_{w \in \mathcal{M}^\ell} p_w^* |\log p_w^*|$  be the entropy of fundamental intervals and, besides  $C(\mathcal{S})$  encountered before, define the constants

$$C_1(\mathcal{S}) = 1 - \sum_{k \geq 2} \frac{1}{k(k-1)} K^{[k]}(\mathcal{S}) = 1 - \lim_{\ell \rightarrow \infty} \sum_{w \in \mathcal{M}^\ell} p_w^* \sum_{w \in \mathcal{M}^\ell} (1 - p_{[i|w]}^*) \left| \log(1 - p_{[i|w]}^*) \right|,$$

$$C_2(\mathcal{S}) = \sum_{k \geq 1} \frac{1}{k} K^{[k+1]}(\mathcal{S}) = \lim_{\ell \rightarrow \infty} \sum_{w \in \mathcal{M}^\ell} p_w^* \sum_{w \in \mathcal{M}^\ell} p_{[i|w]}^* \left| \log(1 - p_{[i|w]}^*) \right|.$$

For random tries built from  $n$  words emitted by a source  $\mathcal{S}$ , asymptotics of expectations are given in Table 3.

**3.2. Example.** For a memoryless source with probabilities  $\{p_i\}$ :

$$h(\mathcal{S}) = \sum_{i \in \mathcal{M}} p_i |\log p_i|, \quad C(\mathcal{S}) = \frac{\sum_{i \in \mathcal{M}} p_i \log^2 p_i}{\left(\sum_{i \in \mathcal{M}} p_i \log p_i\right)^2},$$

$$C_1(\mathcal{S}) = 1 - \sum_{i \in \mathcal{M}} (1 - p_i) |\log(1 - p_i)|, \quad C_2(\mathcal{S}) = \sum_{i \in \mathcal{M}} p_i |\log(1 - p_i)|.$$

Similar formulae are available for Markov chains and continued fraction sources. Simulations are in agreement with theory.

## 4. Conclusion and Open Questions

For the average value of the size, a Patricia trie turns out to be better than a trie, and Rényi's condition is not necessary. For the average value of the path length, there is only a correcting term  $C_2$  of order 2, and our proofs made use of Rényi's condition. An open question (see [1] for details) would be to know whether this correcting term remains valid for sources for which Rényi's condition does not hold, although all the natural sources we are aware of do satisfy that condition.

### Bibliography

- [1] Bourdon (Jérémie). – Size and path length of Patricia tries: dynamical sources context. *Random Structures & Algorithms*, vol. 19, n° 3-4, 2001, pp. 289–315. – Special issue “Analysis of Algorithms” dedicated to Don Knuth.
- [2] Clément (J.), Flajolet (P.), and Vallée (B.). – Dynamical sources in information theory: a general analysis of trie structures. *Algorithmica*, vol. 29, n° 1-2, 2001, pp. 307–369.
- [3] Jacquet (Philippe) and Szpankowski (Wojciech). – Analytical de-Poissonization and its applications. *Theoretical Computer Science*, vol. 201, n° 1-2, 1998, pp. 1–62.
- [4] Vallée (Brigitte). – Dynamical sources in information theory: fundamental intervals and word prefixes. *Algorithmica*, vol. 29, n° 1-2, 2001, pp. 262–306.

# New and Old Problems in Pattern Matching

Wojciech Szpankowski

Computer Science Department, Purdue University (USA)

June 25, 2001

Summary by Mireille Régnier

## Abstract

This talk presents three problems in pattern matching and their analysis. Different methods are used, that rely on complex analysis and probability theory.

## 1. Statement of the Problems

Some pattern  $H$  (or a set  $\mathcal{H}$  of patterns) is searched in a text  $T$ . The text  $T$  is generated by a random probabilistic source that is either a Bernoulli source or a Markov source or a mixing source. In the string matching and the subsequence matching problems,  $H$  is given: the model is deterministic. In the *repetitive patterns problem*, in Section 4,  $H$  is a string of  $T$  repeated elsewhere.

## 2. String Matching

One counts the number of occurrences of a given word  $H$  or a given finite set of words,  $\mathcal{H}$ , in a text of size  $n$ . This number is denoted  $O_n(H)$  or  $O_n(\mathcal{H})$ . This counting relies on the decomposition of the text  $T$  onto languages, the so-called *initial*, *minimal*, and *tail languages*.

**Definition 1.** Given two strings  $H$  and  $F$ , the *overlap set* is the set of suffixes of  $H$  that are also prefixes of  $F$ . The suffixes of  $F$  in the associated factorizations of  $F$  form the *correlation set*  $\mathcal{A}_{H,F}$ . In the Bernoulli model, one defines the *correlation polynomial* of  $H$  and  $F$  as

$$A_{H,F}(z) = \sum_{w \in \mathcal{A}_{H,F}} P(w)z^{|w|}.$$

When  $H$  is equal to  $F$ ,  $\mathcal{A}_{H,H}$  is named the *autocorrelation set* and denoted  $\mathcal{A}_{H,H}$ ; the *autocorrelation polynomial* is defined as

$$A_H(z) = \sum_{w \in \mathcal{A}_{H,H}} P(w)z^{|w|}.$$

For example, let  $H = 11011$  and  $F = 1110$ . Then the overlap set of  $H$  and  $F$  is  $\{11, 1\}$  and the correlation set is  $\mathcal{A}_{H,F} = \{10, 110\}$ . Similarly,  $\mathcal{A}_{F,H} = \{11\}$ . It is worth noticing that  $\mathcal{A}_{F,H} \neq \mathcal{A}_{H,F}$ . Intuitively, the concatenation of a word in  $\mathcal{A}_{H,F}$  to  $H$  creates an (overlapping) occurrence of  $F$ .

**Definition 2.** Let  $H$  be a given word.

- (i) The *initial language*  $\mathcal{R}$  is the set of words containing only one occurrence of  $H$ , located at the right end.

- (ii) The *tail language*  $\mathcal{U}$  is defined as the set of words  $u$  such that  $Hu$  has exactly one occurrence of  $H$ , which occurs at the left end.
- (iii) The *minimal language*  $\mathcal{M}$  is the set of words  $w$  such that  $Hw$  has exactly two occurrences of  $H$ , located at its left and right ends.

With these notations, any text that contains exactly  $k$  occurrences of  $H$ ,  $k \geq 1$ , rewrites unambiguously as

$$rm_1 \dots m_{k-1}u$$

where  $r \in \mathcal{R}$ ,  $m_i \in \mathcal{M}$ , and  $u \in \mathcal{U}$ . In other words, this set  $\mathcal{T}_k$  of words satisfies  $\mathcal{T}_k = \mathcal{R}\mathcal{M}^{k-1}\mathcal{U}$ . The power of this approach comes from the equations that can be written on these languages, that translate into equations on their generating functions in the Bernoulli model *and* the Markov model. Moreover, it turns out that these generating functions—hence the whole counting problem—only depend on the probability of  $H$ , denoted  $P(H)$ , and the so-called correlation set.

**Theorem 1.** *Let  $H$  be a given pattern of size  $m$ , and  $T$  be a random text generated by a Bernoulli model. The generating function of the set  $\mathcal{T}_k$  satisfies*

$$T_k(z) = z^m P(H) \frac{(D_H(z) + 1 - z)^{k-1}}{D_H(z)^{k+1}}, \quad k \geq 1,$$

$$T_0(z) = \frac{A_H(z)}{D_H(z)}$$

where

$$D_H(z) = (1 - z)A_H(z) + z^m P(H).$$

Moreover, the bivariate generating function satisfies

$$T(z, u) = \sum_k T_k(z) u^k = \frac{u}{1 - u \frac{D_H(z) + 1 - z}{D_H(z)}} \frac{z^m P(H)}{D_H(z)^2}$$

These results extend to the Markovian model and to the case of multiple pattern matching [3].

### 3. Subsequence Matching

A pattern  $W = w_1 \dots w_m$  is hidden in a text  $T$  if there exist indices  $1 \leq i_1 < \dots < i_m \leq n$  such that  $t_{i_1} = w_1, \dots, t_{i_m} = w_m$ . For example, *date* is hidden 4 times in the text *hidden pattern* but it is not a substring. We focus on cases where the sequence of indices satisfies additional constraints  $i_{j+1} - i_j \leq d_j$ , where  $d_j$  is either an integer or  $\infty$ . Such a sequence is called an *occurrence*. One denotes  $(d_1, \dots, d_{m-1})$  by  $\mathcal{D}$ . For example, when  $\mathcal{D} = (3, 2, \infty, 1, \infty, \infty, 4, \infty)$  the set  $I = (5, 7, 9, 18, 19, 22, 30, 33, 50)$ , satisfies the constraints.

The number of occurrences,  $\Omega_n$ , is asymptotically Gaussian. This is proved in [1] by the moments method: all moments of the properly normalized random variable converge to the corresponding moments of the Gaussian law. For any sequence  $I$  that satisfies the constraints, one denotes  $X_I$  the random variable that is 1 if  $t_{i_1} = w_1, \dots, t_{i_m} = w_m$ . Then,

$$\Omega_n = \sum_I X_I.$$

The computation of the moments relies on a generalization of correlation sets. Let

$$\mathcal{U} = \{u_1, \dots, u_{b-1}\}$$

be the subset of indices  $j$  for which  $d_j = \infty$ . Any occurrence  $I$  satisfying the constraints can be divided into  $b$  blocks:

$$[i_1, i_{u_1}], [i_{u_1+1}, i_{u_2}], \dots, [i_{u_{b-1}+1}, i_m].$$

The collection of these blocks is called the *aggregate* of  $I$  and denoted  $\alpha(I)$ . In the example above, the aggregate  $\alpha(I)$  is

$$\alpha(I) = [5, 9], [18, 19], [22], [30, 33], [50].$$

*Deriving the mean.* The collection of occurrences of  $W$  can be described as

$$\mathcal{A}^* \times \{w_1\} \times \mathcal{A}^{\leq d_1} \times \{w_2\} \times \dots \times \mathcal{A}^{\leq d_{m-1}} \times \{w_m\} \times \mathcal{A}^*,$$

where  $\mathcal{A}$  is the alphabet and  $\mathcal{A}^{\leq d_j}$  is the collection of words of size less than or equal to  $d_j$ . It follows that the generating function of expectations is

$$\sum_n \mathbf{E}(\Omega_n) z^n = \frac{1}{(1-z)^{b-1}} \times \prod_{i=1}^m p_{w_i} z \times \prod_{i \notin \mathcal{U}} \frac{1-z^{d_i}}{1-z},$$

where  $p_{w-i}$  is the probability of character  $w_i$ . Hence, the expectation satisfies

$$(1) \quad \mathbf{E}(\Omega_n) = \frac{n^b}{b!} \prod_{i \notin \mathcal{U}} d_i \prod_{i=1}^m p_{w_i} \left( 1 + O\left(\frac{1}{n}\right) \right)$$

*Deriving the variance and higher moments.* The variance rewrites

$$\mathbf{Var}(\Omega_n) = \sum_{I, J} \mathbf{E}(X_I X_J) - \mathbf{E}(X_I) \mathbf{E}(X_J).$$

In the Bernoulli model, the two random variables  $X_I$  and  $X_J$  are independent whenever the blocks of  $I$  and  $J$  do not overlap. Hence, the contribution to the variance is zero. If  $\alpha(I)$  and  $\alpha(J)$  overlap, one defines the aggregate  $\alpha(I, J)$  as the set of blocks obtained by merging the blocks of  $\alpha(I)$  and  $\alpha(J)$  that overlap. The number of blocks in  $\alpha(I, J)$ , denoted  $\beta(I, J)$ , is upper bounded by  $2b - 1$ . For such a pair  $(I, J)$ , the text can be rewritten as an element of the language

$$\mathcal{A}^* \times \mathcal{B}_1 \times \mathcal{A}^* \times \dots \times \mathcal{B}_{\beta(I, J)} \times \mathcal{A}^*$$

and the generating function of the covariance rewrites

$$\sum_n \mathbf{Var}(\Omega_n) z^n = \sum_{p \geq 1} \sum_{\beta(I, J) = 2b-p} \frac{1}{(1-z)^{2b-p}} P_p(z),$$

where  $P_p$  are polynomials of the variable  $z$  that generalize the correlation polynomials defined in [2] (see Definition 1). The asymptotic order of each term is  $n^{2b-p}$ . Hence, the dominating contribution is due to the intersecting pairs such that  $\beta(I, J) = 2b - 1$ , and

$$\mathbf{Var}(\Omega_n) \sim n^{2b-1} \sigma^2$$

where the variance coefficient  $\sigma$  can be easily evaluated for any given pattern by dynamic programming.

The proof is similar for higher moments.

#### 4. Repetitive Pattern Matching

Given a pattern  $H$  found in a text  $T$ , one searches for a second *approximate* occurrence of  $H$ . A word  $F$  is a  $D$ -approximate occurrence of a word  $H$  if the Hamming distance between  $F$  and  $H$  is smaller than  $D$ . Recall that the Hamming distance between two words of size  $m$ , say  $H = H_1 \dots H_m$  and  $F = F_1 \dots F_m$  is

$$d_H(H, F) = \sum_{i=1}^m 1_{H_i \neq F_i}.$$

The usual parameters on trees, such as the *depth of insertion*, *height*, *fill-up*,  $\dots$ , are extended in the approximate case. Notably:

**Definition 3.** The *depth*  $L_n$  is the largest integer  $K$  such that

$$\min \left\{ d(T_i^{i-K+1}, T_n^{n+K}) \mid 1 \leq i \leq n - K + 1 \right\} \leq D.$$

Rényi's entropy is generalized. Given a word  $H$ , the  $D$ -ball with center  $H$ , denoted  $B_D(H)$ , is the set of words that are within distance  $D$ .

**Definition 4.** Given a text  $T$ , Rényi's entropy of order 0 is

$$r_0(D) = \lim_{k \rightarrow \infty} \frac{-\mathbf{E} \left[ \log \mathbf{P} (B_D(T_1^k)) \right]}{k},$$

when this limit exists.

Asymptotic properties are proved for the depth, the height and the fill-up, that depend on Rényi's entropy. Notably, the convergence in probability of the depth of insertion in a trie extends for this approximate scheme:

$$\frac{L_n}{\log n} \rightarrow \frac{1}{r_0(D)}, \quad n \rightarrow \infty.$$

The proof relies on the subadditive ergodic theorem and asymptotic equipartition property.

#### Bibliography

- [1] Flajolet (P.), Guivarc'h (Y.), Vallée (V.), and Szpankowski (W.). – Hidden patterns statistics. In *ICALP'01*. – 2001. Proceedings of the 28th ICALP Conference, Crete, Greece, July 2001.
- [2] Guibas (L. J.) and Odlyzko (A. M.). – String overlaps, pattern matching, and nontransitive games. *Journal of Combinatorial Theory. Series A*, vol. 30, n° 2, 1981, pp. 183–208.
- [3] Régnier (M.) and Szpankowski (W.). – On pattern frequency occurrences in a Markovian sequence. *Algorithmica*, vol. 22, n° 4, 1998, pp. 631–649. – Average-case analysis of algorithms.

## Genome Analysis and Sequences with Random Letter Distribution

*Michel Termier*

Institut de Génétique et Microbiologie, Université de Paris XI (France)

April 2, 2001

*Summary by Mathias Vandenberg*

### Abstract

The information content of genomes of different organisms reflects their mode of physical organisation. For the last decades the wet lab biologist's research interests has been to decipher this information content, with the purpose of extracting useful biological features. The reliability of the information extraction process, mainly based on the textual nature of the underlying messages, was hard to achieve. Therefore, an approach based on the comparison of naturally occurring sequences and randomly generated sequences, is used for discerning the artefacts in sequences and for improving the power of our genome models.

### Introduction

The building plan for vegetative life is based on the assembly and catalytic function of proteins and active RNAs. The complete set of instructions that is needed to generate the building blocks of the reproductory system is called a "genome." Any production of living tissue from these building blocks will give rise to an accumulation of secondary metabolites, which are of adverse influence for the survival of the species. The secondary effects of metabolite production are at the basis for the requirement of the genome to be able to respond to the induced environmental changes. To counter this problem, a cell of an organism will only bring to expression those genes that are required at some specific moment in the cell's life cycle. For this purpose, a genome disposes of regulatory systems in the generation processes of building blocks. These systems can be compared to logical gates that are situated in upstream sequences of most information that needs to be processed. This permits a modulation in the usage of information. The genomic information is stocked in a linear fashion, which facilitates the tracking of information. At the time the sequencing of the human genomic sequence is being accomplished, several tasks remain to be addressed:

- the decomposition of the genomic sequence into streams of messages;
- the distinction of these "messages" in contrast to the "non-coding bulk information";
- assignment of biologically significant functions to the messages.

Our bioinformatics team is mainly interested in providing an answer to basically two questions:

1. How can messages be extracted from genomic sequences in order to perform the function assignment task?
2. What is the nature of the message contained within any linear macromolecular structure?

## 1. First Task: Message Extraction and Function Assignment

The approach consists in observing the known words in the vocabulary of the genome. These known words have been indexed through many years of genetic experiments, with the use of techniques handled in molecular biology wet labs. Through this biology-related knowledge accumulation, the following facts are at the basis for the study of genomic sequences:

- the start and end points (the START and STOP signals) of a nucleic acid sequence correspond to the beginning and to the end of a diffusible product (= protein);
- the information content of a nucleic acid sequence is translated in a unidirectional fashion to the corresponding protein through some basic transcription rules:

$$\text{DNA} \rightarrow \text{messenger (mRNA)} \rightarrow \text{protein};$$

- for the yeast organism, experiments have demonstrated that at least 99 triplets are required between the START and STOP signals, which leads to the coding sequence expression [5]:

$$\text{START}(n^3 \setminus \text{STOP})^{99}(n^3 \setminus \text{STOP})^* \text{STOP},$$

with  $n = \{a, c, g, t\}$ ,  $\text{START} = atg$ ,  $\text{STOP} = \{taa, tag, tga\}$ ;

- by replacing the T-based nucleotides with U, this expression proves to be universally true for the genes describing the intermediate messenger molecules (mRNA) in the steps between DNA and protein;
- for the genomic sequences of higher eucaryotes, the protein-describing sequences are interspersed with non-coding intronic sequences (introns, non-coding bulk information);
- a multitude of other signals exists, regulating the expression of specific coding regions, and responsible for the organism's physiological response in precise environmental conditions.

**1.1. Mechanisms for processing signals in messages.** There exist mechanisms for processing complex signals, both within eucaryotes as well as within viral species. The *eucaryotic* mechanism is described as *alternative splicing*: a protein-encoding sequence can generate different proteins at the time mRNA is being spliced, according to different translational systems. Sample mechanisms for this group of organisms are read-through (the transcription machinery is reading through and beyond the STOP codon), and hopping (the transcription machinery is skipping the STOP codon and the codons surrounding it). The *retro-viral* mechanism is called *re-encoding*, which implies that different proteins can be obtained at the time the mRNA is being translated. Sample mechanisms for this group are frameshift (the reading frame for translation is changed, which induces an alteration of the encoded amino acids), read-through and hopping. Several features can be conferred to some sequences that are responsible for a frameshift:

1. Slipping sequences (structure X XXY YYZ).
2. A badly positioned classical STOP signal: the ribosome loses his grip on the sequence and gets positioned again in phase  $-1$ .
3. A ribosome-blocking structure.

Regulatory sequences that are responsible for the modulation of DNA transcription in a less error-prone fashion are:

1. *Inhibitor signals*. Their role is to bind proteins so that the RNA polymerase can no longer bind to the sequence to initiate transcription.
2. *Activator signals*. There exists a multitude of signals per protein-encoding sequence, according to the specific function of the protein to be generated.

Usually, these regulatory sequences are short sequences, whose observed frequency is higher (hence unexpected) in comparison to a random word composed of the same letters.



**1.2. Modelling a genomic sequence.** A Markov model is frequently used for modelling a genomic sequence. The number of sequences that can be generated by this model, increases with the order of the Markov model, and reaches a plateau.

For a Bernoulli-type distribution of the nucleotides, the actual sequence follows a Gaussian distribution. Additionally, when [A+T] increases, the amount of START and STOP signals increases. This implies that the certainty of finding a gene increases.

Regulatory signals are words with biased composition, with respect to the global word distribution of the sequence. These signals have been selected for their properties in the course of evolution. They have been generated according to mechanisms which include random events [2, 3].

**1.3. The importance of codon usage biases.** In the context of genetic expression, the codon usage bias is correlated with the level of tRNAs available, and with the abundance of protein generated. The level of protein-encoding sequences that are significantly biased is of the order of 20% of the total amount of sequences. Within this respect, several observations have been made:

- the biased structure helps in regulating the transcription turnover [6];
- there is a positional codon bias according to the strand on which the gene is situated [4];
- there is a codon usage bias according to the life cycle of the organism and the cellular location of the metabolic activity [1];
- there is a bias in relation with mRNA stability problems [9];
- some horizontal transfers can have effects on the codon usage [8].

The codon usage bias determining the level of codons corresponding to the amino acids of proteins has a direct effect in the genomic sequence composition of the organism. This bias, which is the result of an interaction of horizontal transfer and metabolic constraints, is at the basis of the selection of efficient proteins. The codon usage bias reveals information about the nucleotide triplet usage of the encoded protein and about the eventual external origin of the sequence in the organism. The significance of the codon usage bias can be evaluated by using weighted linguistics approaches. This consists in heuristically weighting the codons used to encode the amino acids, instead of using an average weight for every amino acid that is encoded by several triplets. This prevents from having resulting frequencies that diverge from the observed values.

Nevertheless, the probability of finding reasonable codon compositions through linguistic methods is fairly low, because:

- global linguistics are calculated on a larger set of oligonucleotides than the number of oligos that determine the proteins;
- the number of codons in a gene equals one third of the number of possible triplets;
- the different genes are built up from codons of different composition, and this is increasing the background noise accordingly.

## 2. Second Task: Determining the Nature of the Message

Life on any other planet besides Earth can only be detectable for us if it is based on our carbon chemistry. Any sequential organic macromolecule contains constitutional information, if textual organization can be detected within it.

Different approaches exist for the detection of organized information:

1. *Complexity analysis of sequences.* The complexity of sequences is difficult to compute. Ed Trifonov introduced in 1990 the notion of linguistic complexity [7] that reflects the linguistic wealth of a sequence. This complexity is easily computable as  $C = \prod_{i=1}^{n-1} u_i$ , with  $u_i$  the ratio of the words found in a sliding window at position  $i$  in a sequence, versus the total number of different words that could possibly be found. Computations are made along

windows, by multiplying the  $u$  ratios of words of all possible lengths in the window. This implies that all redundancies are eliminated. The value of  $C$  varies from 0 to 1.

2. *Shannon's entropy measure*  $H(X) = -\sum_i P(x_i) \cdot \log(P(x_i))$ . The entropy  $H(X)$  is maximal in the case of a random equiprobable sequence. A reduction in entropy corresponds to a generation of information. This implies that the measurement of the amount of information can be done by:

$$I(X) = H(\text{without message}) - H(\text{with message}).$$

This way, the amount of information can be quantified by comparing a randomly generated Markovian sequence (sequence without message) with a naturally occurring sequence. This measure is related to global information content, but does not give any idea on the distribution of the coding zones of the sequence. It is a common observation in information-bearing texts that coding zones are separated from each other by areas that are more or less deprived of information. If the hypothesis of a non-terrestrial genome makes sense, then its linguistics must respond to the following criterions:

- it must be based on a restricted alphabet;
- it bears coding subsequences that are separated from each other in a way that is recognizable by certain molecules;
- the coding subsequences are likely to share some common characteristics;
- these sequences are constructed using linguistics that can vary from one “genome” to another;
- the reading direction of the sequences is oriented (this should facilitate their regulation);
- the method used to copy the message determines the ordered relation between the coding sequences.

### Bibliography

- [1] Chiapello (H.), Ollivier (E.), Landès-Devauchelle (C.), Nitschké (P.), and Risler (J.-L.). – Codon usage as a tool to predict the cellular location of eukaryotic ribosomal proteins and aminoacyl-tRNA synthetases. *Nucleic Acids Research*, vol. 27, n° 14, 1999, pp. 2848–2851.
- [2] Grantham (R.). – Workings of the genetic code. *Trends in Biochemical Sciences*, n° 5, 1980, pp. 327–331.
- [3] Grantham (R.), Gautier (C.), Gouy (M.), Jacobzone (M.), and Mercier (R.). – Codon catalog usage is a genome strategy modulated for gene expressivity. *Nucleic Acids Research*, n° 9, 1981, pp. 43–74.
- [4] Lafay (B.), Lloyd (A.T.), McLean (M.J.), Devine (K.M.), Sharp (P.M.), and Wolfe (K.H.). – Proteome composition and codon usage in spirochaetes: species-specific and DNA strand-specific mutational biases. *Nucleic Acids Research*, n° 27, 1999, pp. 1642–1649.
- [5] Oliver (S.G.), van der Aart (Q.J.), Agostoni-Carbone (M.L.), Aigle (M.), Alberghina (L.), Alexandraki (D.), Antoine (G.), Anwar (R.), Ballesta (J.P.), and Benit (P.). – The complete DNA sequence of yeast chromosome III. *Nature*, n° 357, 1992, pp. 38–46.
- [6] Olivier (E.), Delorme (M.O.), and Henaut (A.). – Dos DNA occurs along yeast chromosomes, regardless of functional significance of the sequence. *Comptes rendus de l'Académie des sciences Paris*, n° 318, 1995, pp. 599–608.
- [7] Popov (O.), Segal (D. M.), and Trifonov (E. N.). – Linguistic complexity of protein sequences as compared to texts of human languages. *BioSystems*, n° 38, 1996, pp. 65–74.
- [8] Rocha (E.P.C.), Viari (A.), and Danchin (A.). – Oligonucleotide bias in bacillus subtilis: general trends and taxonomic comparisons. *Journal of Applied Probability*, n° 36, 1998, pp. 179–193.
- [9] Seffens (W.) and Digby (D.). – mRNAs have greater negative folding free energies than shuffled or codon choice randomized sequences. *Nucleic Acids Research*, n° 27, 1999, pp. 1578–1584.

## Random Sequences and Genomic Analysis

*Alain Denise*

IGM & LRI, Université de Paris XI (France)

April 2, 2001

### **Abstract**

A crucial problem in genomic analysis is to distinguish “biologically significant” signals in sequences from those that are part of the ground noise. To this end, biological sequences are compared with those expected to be met “by chance.” Models of random sequences frequently used in this perspective will be briefly described, as will be analytical methods (developped notably in the Algorithms Project at Inria!) and experimental methods (random sequence generation) used to solve these problems. Then, recent works on random sequence generation according to a model that is more constrained than those studied so far will be presented, together with a framework in which it applies to the study of genomic sequences.



# The Primal-Dual Schema for Approximation Algorithms: Where Does It Stand, and Where Can It Go?

Vijay Vazirani

Georgia Institute of Technology (USA)

December 11, 2000

Summary by Claire Kenyon

## Introduction

NP-hard problems cannot be solved exactly *and* efficiently at the same time. Can they be approximated in polynomial time? When doing so, we want a guarantee: for every instance, the solution must be within some factor of the optimal solution. Such questions are discussed systematically in Vijay Vazirani's book [6] on which the present lecture is based.

Linear programming duality theory provides many efficient algorithms with a good *approximation* factor. Designing *exact* algorithms is a main topic of the paper by Grötschel, Lovász, and Schrijver in 1981; see [3]. As we shall see, the primal-dual scheme provides the broad outline of an algorithm; working out the details for each individual problem then often provides a specific approximate solution with good complexity characteristics.

### 1. The Vertex Cover Problem

Given a graph, a subset of its vertices is a vertex cover if and only if every edge has at least one vertex in the subset. Each vertex has a cost—the cover having cost equal to the sum of the costs of its vertices—and we wish to obtain the cover of minimum cost. This problem is NP-hard (as proved by Karp in 1971, see [5]). We need to compare the cost of an approximate solution constructed by an algorithm to the cost of the optimal solution (OPT), but we do not know the cost of OPT; so we need a good lower bound on the cost of OPT. This is a key first step in the design of approximation algorithms.

**1.1. Linear programming approximation.** To the end of obtaining bounds on OPT for vertex cover, we start with an integer programming formulation of it. There is one variable  $x_v$  for each vertex  $v$ , and it is equal to 0 or 1; there is one constraint for each edge  $\{u, v\}$ , i.e.,  $x_u + x_v \geq 1$ , which expresses that the sum of its two endpoint variables is at least 1; it is then required to minimize a linear combination of vertex variables times vertex costs, i.e.  $\sum_v \text{cost}(v) \times x_v$ .

We then do a relaxation of the problem by allowing the variables to be real numbers between 0 and 1 (instead of being integers). Each feasible solution provides a *fractional* vertex cover whose cost is necessarily a lower bound to OPT. We know since the works of Khachian and Karmarkar around 1980 that linear programming is polynomial-time solvable, both theoretically and effectively. The best fractional solution is thus polynomial-time computable, which gives us our lower bound. The relaxation algorithm is then as follows:

**Linear Programming Algorithm.** First find the optimal fractional solution, then put in the cover all vertices  $v$  such that  $x_v \geq 1/2$ . It is easy to see that this is a vertex cover, and

the cost is at most 2 times the lower bound, hence at most 2 times OPT. This algorithm has the defect of requiring to solve a linear program, a polynomial-time but expensive step.

**1.2. A combinatorial algorithm.** The principle of a combinatorial algorithm that has an approximation factor of 2 is as follows. Initially the cover  $C$  is empty. While  $C$  is not a vertex cover, pick an uncovered edge  $\{u, v\}$ , look at the smaller of the two current costs of  $u$  and  $v$ , subtract this smaller current cost from the costs of  $u$  and of  $v$ , put the corresponding vertex in  $C$ , and charge its cost to the edge. What we charge to the edges turns out (by induction) to be a lower bound on OPT. The cost of the cover is obviously at most twice the amount charged to the edge. Hence this technique gives rise to a combinatorial algorithm with an approximation factor of 2; the outcome is in fact a very fast linear-time algorithm.

This alternative algorithm is actually related to the LP-based algorithm seen previously. There is currently no approximation algorithm known which beats this factor of 2.

## 2. LP Relaxation and Dual LP

An original linear programming (LP) problem (the “primal”) always admits a “dual” formulation.

**Primal linear program (LP).** Determine  $\min \sum_v \text{cost}(v) \times x_v$  subject to  $\forall e \ x_u + x_v \geq 1$  and  $\forall v \ x_v \geq 0$ .

One can prove an upper bound on the OPT solution to the primal LP by exhibiting a particular solution  $(x_v)$  which satisfies all the constraints. One can prove a lower bound by exhibiting a particular linear combination of the constraints which equals the objective function. This corresponds to a dual LP solution.

**Dual linear program.** Determine  $\max \sum_e y_e$  subject to  $\forall v \ \sum_{e|v \in e} y_e \leq \text{cost}(v)$  and  $\forall e \ y_e \geq 0$ .

Equality of the optimal solutions of the primal and dual programs constitutes the *strong duality theorem*. The idea of a primal-dual algorithm is precisely to use a feasible solution of the dual LP as a *lower* bound on OPT. (Note that duality exchanges ‘min’ and ‘max’.)

How to design the primal-dual algorithm? We need the complementary slackness theorem, which says that if  $x$  is a feasible solution to the primal LP and  $y$  a feasible solution to the dual LP, then both are optimal if and only if for every  $v$  either  $x_v = 0$  or  $\sum_{e|v \in e} y_e = \text{cost}(v)$ , and for every edge  $e$  either  $y_e = 0$  or  $x_u + x_v = 1$ . Thus if  $(x, y)$  are not both optimal, we can find a slack and decrease the corresponding  $x_v$  or increase the corresponding  $y_e$ . To design an *approximation* algorithm, we change the equality relative to  $\text{cost}(v)$  into an inequality.

**Primal-dual algorithm for vertex cover.** Initially  $x$  and  $y$  are set to 0. Let  $C$  be the set of “tight” vertices. While  $C$  is not a cover, do: pick an uncovered edge  $e$ , pick  $y_e$  and raise it until one of its two endpoints is tight. Iteratively improve the primal and dual solutions until a primal feasible solution is obtained; compare the primal and dual solutions to establish the approximation guarantee.

The set cover problem can be solved in the same fashion. In this problem, one has a set  $U$  of elements and a collection of subsets  $U_I$ , each with a positive cost, and one wishes to construct a minimum collection of subsets whose union is  $U$ . (Exercise: Let the frequency of element  $e$  be the number of subsets containing  $e$ , and let  $f$  be the maximal frequency of an element. Design a primal-dual approximation algorithm with an approximation factor of  $f$ .) By design, the best approximation factor we can get by these methods is the integrality gap, i.e., the ratio between the OPT solution to the integer linear program and the OPT solution to the relaxed linear program.

*History.* This paradigm started in 1955 (Kuhn) in the context of weighted bipartite matching. The primal-dual terminology is due to Dantzig, Ford, and Fulkerson in 1956. It was used to design *exact* algorithms for many polynomial-time algorithms much before linear programming was recognized to be polynomial-time solvable. Examples of this technique include matching, network flow, shortest paths, minimum spanning trees, branchings, and so on.

These exact primal-dual algorithms all use the fact that the polyhedron defined by the LP has integral vertices, and so the LP has integral optimal solutions. It is the relaxation of the complementary slackness solutions that essentially leads to approximation algorithms.

In 1981 Bar-Yehuda and Even [2] gave an approximation algorithm with a factor of 2 for vertex cover. In retrospect, their work can be reframed in the setting of primal-dual algorithms so that it can be regarded as the first primal-dual approximation algorithm.

### 3. Other Problems

Many other problems can be solved approximately using the primal-dual approach. We give a short list below and refer to the book [6] for details.

*Steiner tree problem.* Given a graph and a set of red vertices in the graph, find a tree which connects all the red vertices (possibly using the other graph vertices in the tree) and has minimal total cost. Gauß also had a version on the plane (given a set of vertices in the plane, connect them into a tree, possibly branching out at other points in the plane).

*Steiner network problem.* Design a network with a prescribed number of edge-disjoint paths between pairs of vertices. There are numerous applications of this problem in networks.

*Steiner forest problem.* The connectivity requirement is 0 or 1 between pairs of vertices. In 1991 factor-of-2 algorithms were designed by Agrawal, Klein, and Ravi [1] on the one hand, Goemans, Williamson on the other hand. These authors use the idea of *simultaneously* raising the violated minimal constraints. In 1992 Williamson, Goemans, Vazirani, and Mihail [7] found a  $2k$  approximation algorithm for the extended Steiner network problem when the maximum connectivity requirement is  $k$ ; their algorithm has been implemented at Bellcore.

*Facility location problem.* What is given is a set of locations for installing proxy servers and a set of clients; the goal is to minimize the sum of server installation cost plus the sum of client's connection costs. For this problem, in the late 1990s, several primal-dual approximation algorithms using LP rounding were designed; they are nice but not so practical. Recently Jain and Vazirani [4] got an approximation algorithm with a factor of 3 based on a practical combinatorial solution, which stems from the primal-dual scheme.

*The  $k$ -median problem.* This problem is like the facility location problem, except that facilities are free, one is constrained to open at most  $k$  facilities; what is required is to minimize the connection cost. This has applications to data mining *inter alia*. In 1998 there was an  $O(1)$ -approximation primal-dual algorithm based on LP-rounding, but that again had the disadvantage of requiring to solve a linear program. In 1999 Jain and Vazirani designed a combinatorial algorithm that is more complicated and relies on randomized rounding. This last algorithm can then be derandomized using the method of conditional expectations.

The techniques discussed in this talk are very robust in the sense that once you solve one problem, you can get solutions to many closely related problems as well.

#### 4. Open Problems

Our approximation algorithms always deal with dual variables in a greedy fashion, whereas exact primal-dual algorithms are much more sophisticated: there is a long way to go to bring the two approaches closer!

Some of the main open problems are: get a factor better than 2 for vertex cover, and better than  $3/2$  for the traveling salesman path; get a factor of 2 for the Steiner network; design a bidirected cut relaxation for Steiner trees.

#### Bibliography

- [1] Agrawal (Ajit), Klein (Philip), and Ravi (R.). – When trees collide: an approximation algorithm for the generalized Steiner problem on networks. *SIAM Journal on Computing*, vol. 24, n° 3, 1995, pp. 440–456.
- [2] Bar-Yehuda (R.) and Even (S.). – A linear-time approximation algorithm for the weighted vertex cover problem. *Journal of Algorithms*, vol. 2, n° 2, 1981, pp. 198–203.
- [3] Grötschel (M.), Lovász (L.), and Schrijver (A.). – The ellipsoid method and its consequences in combinatorial optimization. *Combinatorica*, vol. 1, n° 2, 1981, pp. 169–197.
- [4] Jain (Kamal) and Vazirani (Vijay V.). – Approximation algorithms for metric facility location and  $k$ -median problems using the primal-dual schema and Lagrangian relaxation. *Journal of the ACM*, vol. 48, n° 2, 2001, pp. 274–296.
- [5] Karp (Richard M.). – Reducibility among combinatorial problems. In *Complexity of computer computations (Proc. Sympos., IBM Thomas J. Watson Res. Center, Yorktown Heights, N. Y., 1972)*, pp. 85–103. – Plenum, New York, 1972.
- [6] Vazirani (Vijay V.). – *Approximation algorithms*. – Springer-Verlag, Berlin, 2001, xx+378p.
- [7] Williamson (David P.), Goemans (Michel X.), Mihail (Milena), and Vazirani (Vijay V.). – A primal-dual approximation algorithm for generalized Steiner network problems. *Combinatorica*, vol. 15, n° 3, 1995, pp. 435–454.



## Distributed Decision Making: The Case of No Communication

*Paul Spirakis*

Computer Technology Institute, Patras University (Greece)

November 20, 2000

### **Abstract**

We examine the case of  $n$  agents trying to achieve a global goal without any communication. Our analysis for the bottleneck probability of scheduling loads in common finite buffers also includes the first exact expressions for the density of a general sum of uniform random variables, this being obtained via a new polyhedral combinatorial approach.



## **Part III**

# **Computer Algebra and Applications**



## Thirty Years of Integer Factorization

*François Morain*

LIX, École polytechnique (France)

February 5, 2001

*Summary by Marianne Durand*

### Abstract

Factoring integers is quite an old challenge. Thirty years ago, two researchers factored the mythic number  $F_7 = 2^{2^7} + 1$ . A few years later public-key cryptography was born, and with it the famous RSA algorithm. Even if the security of RSA is not equivalent to integer factorization, factoring the RSA key is the simplest way to decode everything, so a lot of people tried to factor. In 1990,  $F_9 = 2^{2^9} + 1$ , the ninth Fermat number was factored, with the help of hundreds of computers. In august 1999, it was the turn of the first ordinary 512-bit integer. What follows is a survey of thirty years of factorization, describing the different methods used and the technical problems met.

### 1. Introduction

Factoring is of great interest since it allows to use the properties of prime number in arithmetic. It is the keystone of the RSA algorithm, the mostly used encryption algorithm. RSA is an asymmetric public key algorithm that is based on the fact that the product of two very large prime numbers can not be easily factored, whereas to check if a number is prime can be done quickly. The complexity class of testing the primality of an integer is  $NP \cap co-NP$ . Factoring a number is in  $NP$ , but can be done in polynomial time on a quantum computer!

Method	Complexity
sieve	$p$
$\rho$	$\sqrt{p}$
elliptic curve method	$L_p[1, 1/2]$
quadratic sieve (QS)	$L_N[1/2, c]$
number field sieve (NFS)	$L_N[1/3, c]$

TABLE 1. Complexity of factorization methods ( $N$  is the integer to be factored,  $p$  its smallest factor)

A lot of different methods exist to factor a number, starting from the linear sieve up to the algebraic sieve, including methods based on elliptic curves. Their complexity can be expressed in terms of the function

$$L_x[\alpha, c] = e^{c \log^\alpha x (\log \log x)^{1-\alpha}}.$$

Some complexities are given in Table 1. The smallest factor  $p$  of  $N$  is usually of order  $\sqrt{N}$ . The letter  $c$  stands for a constant and is not specified as it depends on the algorithm and its implementation. These methods are detailed in the next section.

## 2. Combination of Congruences

The method of combination of congruences is an extension of Kraitchik's method. The latter aims at finding an integer  $x$  such that  $x^2 \equiv 1 \pmod{N}$  and  $x \not\equiv \pm 1 \pmod{N}$ , then at testing if  $\text{pgcd}(x-1, N)$  is non-trivial. If so, it is a factor of  $N$ . The quadratic congruence approach refines the way the square root of 1 is found. The first step consists in finding pairs of integers  $(u_i, v_i)_{i \in I}$  such that  $u_i^2 \equiv v_i \pmod{N}$  and  $u_i^2 \not\equiv \pm v_i$ . The second step is to find a subset  $J \subset I$  such that  $\prod_{j \in J} v_j$  is a square, noted  $V_J^2$ . This step is detailed later. If we note  $\prod_{j \in J} u_j = U_J$  then step 2 implies  $U_J^2 \equiv V_J^2 \pmod{N}$ . As we also assume that  $V_J$  and  $N$  are together prime (otherwise we have a factor of  $N$ ) then  $x = U_J/V_J \pmod{N}$  is well defined and is a square root of 1. There is a probability greater than 1/2 that it gives a non trivial factorization of  $N$ . This extension is interesting because in order to find the pairs  $(u_i, v_i)$ , we can use an algorithm that eventually rejects or ignore some valid pairs, to go faster. One solution for this is Dixon's method. The idea is to restrict the search to integers  $v_i$  that can be factored on a small set of given small prime integers  $P_k = (p_1, \dots, p_k)$ . To find pairs  $(u_i, v_i)$  according to Dixon's method, we choose an integer  $u_i$ , and try to factor  $u_i^2$  on the set  $P_k$ . If we succeed, then we keep the pair  $(u_i, u_i^2)$ . The integer  $u_i$  has to be greater than  $\sqrt{N}$ , so as to give a non-trivial pair.

Once the pairs  $(u_i, v_i)$  are found, the second step is to find a subspace  $J$  such that  $\prod_{j \in J} v_j$  is a square. As the factorization of each  $v_i$  is already known, this can be seen as a linear algebra problem. Assume that there are  $k+1$  valid pairs available. Consider the matrix  $M$  of size  $(k, k+1)$  with coefficients 0 and 1 viewed in the field  $\mathbb{Z}/2\mathbb{Z}$  and such that  $M[i, j]$  is equal to the exponent of  $p_i$  in the factorization of  $v_j$ . This matrix has a rank smaller than  $k$ , so there exists a linear combination of the columns equals to 0. The subset  $J$  corresponds to the non-zero coefficients in the linear combination, and we can check that  $\prod_{j \in J} v_j$  is a square, because all its factors are of even degree. To exhibit a concrete linear combination equal to zero is made easier by the sparsity of the matrix  $M$ . As a matter of fact, the techniques of Wiedemann or of Lanczos have complexity  $O(k^{2+\epsilon})$  on sparse matrices, whereas the Gauss pivot has complexity  $O(k^3)$ . Then we have the expression of  $V_J$  easily, and a square root of 1 that may give a factorization of  $N$ . This algorithm has a complexity  $L_N[1/2, c]$ , where  $c$  is a constant that depends on the algorithm.

## 3. Sieves

A sieve algorithm searches a lot of candidates satisfying a certain property. Then it makes some tests systematically on all candidates, and at the end keeps the ones that have passed all the tests successfully. One of the first sieves concerning primality and factorization is the Erastothene sieve. The sieve technique is useful in factorization for the search of the set of pairs  $(u, v)$  such that  $u^2 \equiv v \pmod{N}$ .

The basic quadratic sieve, found by Pomerance in 1981 is an extension of the combination of congruence, with a specific choice algorithm for the pairs  $(u_i, v_i)$ . The idea is to choose  $u_i = i + \lfloor \sqrt{N} \rfloor$ , which implies

$$(1) \quad v_i = \left( i + \lfloor \sqrt{N} \rfloor \right)^2 - N.$$

The advantage is that  $v_i$  is close to  $2i\sqrt{N}$ , and thus  $v_i \ll N$ , this increases the probability that the prime factors of  $v_i$  are small. To check that these factors are in the prime number basis  $P_k$  we use a sieve algorithm. This sieve algorithm can be described as follows. First fill an array  $S$  such that  $S[i] = v_i$  for  $i$  from 1 to a bound  $L$ , then for every  $p$  in the prime number basis  $P_k$ , for the two roots of the equation  $(i + \lfloor \sqrt{N} \rfloor)^2 \equiv N \pmod{p}$  noted  $i_{\pm}(p)$ , do  $i \leftarrow i_{\pm}(p)$ , and while  $i < L$  do  $S[i] \leftarrow S[i]/p$  and  $i \leftarrow i + p$ . This algorithm is justified by the equivalence  $p|v_i \iff (i + \lfloor \sqrt{N} \rfloor)^2 \equiv N \pmod{p}$ . Then at the end of the loops, for every  $i$  such that  $S[i] = 1$ ,  $v_i$  is factored on  $P_k$ . The complexity of this algorithm is  $L_N[1/2, 3/\sqrt{8}]$ , and the cost in memory space is  $L_N[1/2, 1/\sqrt{8}]$ . The algorithm can be optimized in many ways, for example the large prime or double large prime variation that we are going to detail in the next paragraph.

The large prime variation owes its name to the use of large primes, not in the prime factor basis, and smaller than the square of the largest prime in the basis  $P_k$ . The sieving stage of the algorithm can easily be modified to find new relations  $v_i = q \prod p^{\alpha_p}$ , where  $q$  is a large prime. Now we can combine two relations using the same large prime  $q$ , namely  $v_1 = q \prod p^{\alpha_p}$  and  $v_2 = q \prod p^{\beta_p}$ , and see that  $v_1 v_2 / q^2$  is factored on  $P_k$ . This large prime technique allows us to search for more “good” pairs  $(u_i, v_i)$  and so to get more candidates to factor  $N$ . In practice it means a speed-up by a factor of approximately 2.5 [5]. The double large prime variation is quite similar, the difference is that two large primes are allowed in the factorization of the integers  $v_i$ . For example if  $v_1 = q_1 q_2 \prod p_j^*$ ,  $v_2 = q_2 q_3 \prod p_j^*$ , and  $v_3 = q_1 q_3 \prod p_j^*$  ( $p^*$  stands for any power of  $p$ ), then  $v_1 v_2 v_3 / (q_1 q_2 q_3)^2$  is factored on the prime basis. The choice of  $v_i$ ,  $v_j$  and  $v_k$  such that their product can be factored upon the prime basis  $P_k$  modulo squares of large primes can be modelled by a graph problem. Let  $G$  be the graph with vertex  $q_i$  and multiple edges  $q_i, q_j$  labelled by the multiples  $v_k$  of  $q_i q_j$ . A useful relation corresponds to a cycle in the graph  $G$ . This technique was used for the sieving step of a 138-digit number in 1990, as the non-optimized sieve was too big to be handled [5] (see also [4]).

The algebraic sieve [2] or number field sieve (NFS) algorithm is based on the factorization in a number field. Given a polynomial  $P \in \mathbb{Z}[X]$  irreducible over  $\mathbb{Q}$ , we will work in the number field  $\mathbb{Q}[X]/(P(X)) = \mathbb{Q}(\theta)$  where  $\theta$  is a root of  $P$ . In the ring  $\mathbb{Z}[\theta]$  we can talk about the primality or the prime decomposition of an element, and the norm of the number  $a - b\theta$  is  $\prod (a - b\theta_i)$  where  $\theta_i$  are all the roots of the polynomial  $P$ . In particular the norm does not depend on the particular choice of  $\theta$ . The description of the algorithm requires the following notation. First let  $m$  be an integer such that  $P(m) \equiv 0 \pmod{N}$ , then consider the ring homomorphism  $\phi$  that maps  $\mathbb{Z}[\theta]$  onto  $\mathbb{Z}/N\mathbb{Z}$  and that satisfies  $\phi(\theta) = m$ . We are now looking for a set  $\mathcal{A}$  of pairs  $(a, b)$  such that  $\prod_{\mathcal{A}} (a - b\theta) = (A - B\theta)^2$  and  $\prod_{\mathcal{A}} (a - bm) = Z^2$ . These properties give  $\phi((A - B\theta)^2) \equiv (A - Bm)^2 \equiv Z^2 \pmod{N}$ . Then  $(A - Bm)/Z$  is a square root of 1, that provides a candidate to factor  $N$ . The choice of the polynomial  $P$  plays a large part in the efficiency of the algorithm [6]. If the degree of  $P$  is  $O((\log N)^{1/3} (\log \log N)^{2/3})$  then the complexity is  $L_N[1/3, c]$ , where  $c$  is a constant.

The way the factorization is done in  $\mathbb{Z}[\theta]$  needs to be explained as it is a non trivial part of the algorithm. The idea is to factor first the norm of  $a - b\theta$ ,  $\text{Norm}(a - b\theta) = \pm \prod p^{\alpha_p(a,b)}$ . This helps because the factorization of  $a - b\theta$  follows the factorization of its norm. If  $p$  is a factor of  $N(a - b\theta)$ , and  $p$  does not divide  $b$  (this being a pathological case), then there exists an integer  $r$  such that  $a - br \equiv 0 \pmod{p}$  and  $P(r) \equiv 0 \pmod{p}$ . We denote by  $[p, r]$  the ideal of  $\mathbb{Z}[\theta]$  such that any element  $x - y\theta$  of  $[p, r]$  satisfies  $\text{Norm}(x - \theta y) \equiv 0 \pmod{p}$  and  $x - yr \equiv 0 \pmod{p}$ . This family of ideals is very interesting because  $(a - b\theta) \in \prod [p, r]^{\alpha_p(a,b)}$ , where  $(a - b\theta)$  is the ideal generated by  $a - b\theta$ .

Now that we know how to factor a number in  $\mathbb{Z}[\theta]$ , we apply the sieve algorithm over the pairs  $(a, b)$ . The factorization algorithm can be optimized by a good choice of the polynomial  $P$  [1]. The variant SNFS, Special Number Field Sieve, targets the numbers  $b^n \pm 1$  by the choice of  $P$ . The

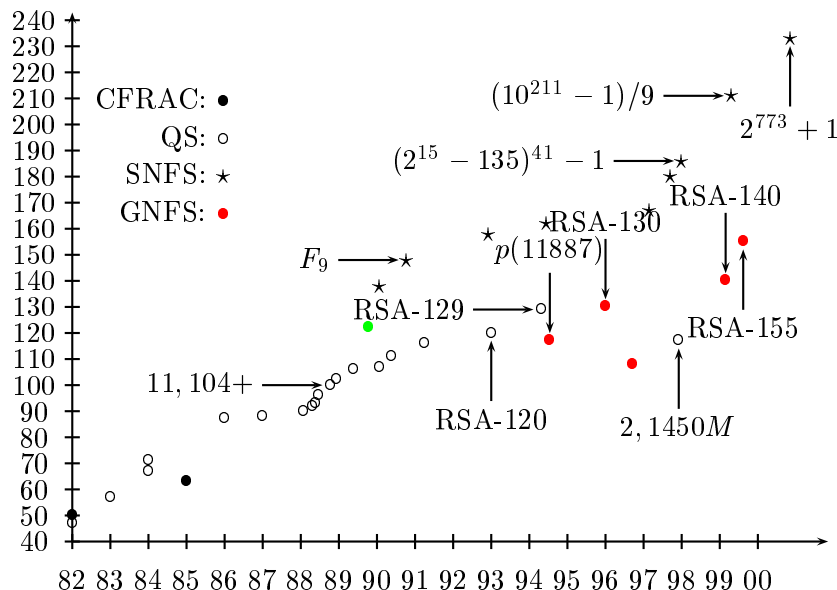


FIGURE 1. Size in bits of the factored numbers depending on the year.

general NFS algorithm becomes better than the quadratic sieve with large primes optimizations for numbers of size around 130 digits.

#### 4. Records and Conclusion

Figure 1 shows the evolution of the factorization records. For each specific algorithm, the progress follows Moore's law that states that the speed of computers double every 18 months. Then for each change of algorithm, there is a jump. Remark that the SNFS algorithm factors specific numbers, that are thus larger than for GNFS that factors general numbers [3]. The linear algebra is often the limiting factor, and unless there is a new idea on the subject, RSA can still be used for some times if used with a key big enough.

#### Bibliography

- [1] Bernstein (Daniel J.) and Lenstra (A. K.). – A general number field sieve implementation. In Lenstra (A.) and Lenstra (H.) (editors), *The development of the number field sieve*, pp. 103–126. – Springer, Berlin, 1993.
- [2] Buhler (J. P.), Lenstra, Jr. (H. W.), and Pomerance (Carl). – Factoring integers with the number field sieve. In Lenstra (A.) and Lenstra (H.) (editors), *The development of the number field sieve*, pp. 50–94. – Springer, Berlin, 1993.
- [3] Cavallar (Stefania), Dodson (Bruce), Lenstra (Arjen K.), Lioen (Walter M.), Montgomery (Peter L.), Murphy (Brian), te Riele (Herman), Aardal (Karen), Gilchrist (Jeff), Guillerm (Gerard), Leyland (Paul C.), Marchand (Joël), Morain (François), Muffett (Alec), Putnam (Chris), Putnam (Craig), and Zimmermann (Paul). – Factorization of a 512-bit RSA modulus. In Preneel (B.) (editor), *Advances in cryptology—EUROCRYPT'00 (Bruges, 2000)*, pp. 1–18. – Springer, Berlin, 2000.
- [4] Lenstra (A. K.) and Manasse (M. S.). – Factoring with two large primes. *Mathematics of Computation*, vol. 63, n° 208, 1994, pp. 785–798.
- [5] Lenstra (Arjen K.) and Manasse (Mark S.). – Factoring with two large primes (extended abstract). In Damgård (I. B.) (editor), *Advances in cryptology—EUROCRYPT '90 (Aarhus, 1990)*, pp. 72–82. – Springer, Berlin, 1991.
- [6] Murphy (Brian). – Modelling the yield of number field sieve polynomials. In Buhler (J. P.) (editor), *Algorithmic number theory (Portland, OR, 1998)*, pp. 137–150. – Springer, Berlin, 1998. Proceedings of the Third International Symposium ANTS-III.



## Variations on Computing Reciprocals of Power Series

*Arnold Schönhage*

Institut für Informatik, Universität Bonn (Germany)

February 5, 2001

*Summary by Ludovic Meunier*

### Abstract

Fast algorithms for polynomial division with remainder are key tools in computer algebra. The power series domain defines a suitable framework where such algorithms can be efficiently constructed. While revisiting Kung's article [5], Arnold Schönhage discusses algebraic complexity bounds for the computation of reciprocals of power series and describes a new algorithm for this task involving Graeffe's root squaring steps.

### 1. Introduction

By means of Newton's iteration, reciprocals of power series modulo  $x^{n+1}$  can be computed with complexity  $O(M(n))$ , where  $M(n)$  denotes the complexity of multiplication (see, e.g., [6] for a survey). However, the Bachmann–Landau  $O$ -notation hides a multiplicative constant, which needs to be investigated, for instance in order to determine cross-over points when a collection of algorithms is available.

Section 2 sets the required background by recalling a few definitions from algebraic complexity. Section 3 presents an algorithm for computing reciprocals of power series, while discussing complexity bounds. Section 4 describes a new algorithm and its implementation over  $\mathbb{Z}$ .

### 2. Algebraic Complexity

Let  $F$  be a field and let  $A(x) = \sum_{i \geq 0} a_i x^i \in F[[x]]$  denote a formal power series of the indeterminate  $x$ . Here, *formal* means that convergence matters are out of concern. Let  $D = F(a_0, a_1, \dots)$  define a domain where  $a_i$ 's are regarded as indeterminates. If  $D$  is endowed with the four arithmetic operations ( $+$ ,  $-$ ,  $*$ ,  $/$ ) and a scalar multiplication, then an algorithm that inputs the power series  $A(x)$  consists of a finite sequence of operations in  $D$ . Counting these operations defines the algebraic complexity, which is an intuitive way of reflecting performances of the algorithm. Two models of complexity are worth considering. The *arithmetic* complexity, denoted by  $L$ ,<sup>1</sup> charges one unit of cost for each operation in  $D$ , while the *nonscalar* complexity, denoted by  $C$ , only counts nonscalar multiplications and divisions.

### 3. Kung's Algorithm Revisited

The underlying algorithm used for the accurate cost calculation is based on Newton's iteration for reciprocals, as discussed by Kung in [5].

---

<sup>1</sup>For notational convenience, arithmetic complexity is also denoted by  $M$  (resp.  $\Lambda$ ) for multiplication (resp. fast Fourier transform).

**3.1. Kung's algorithm.** Let  $R(x)$  be the reciprocal of the unit  $A(x)$  with respect to the field  $D[[x]]$ . Define the function  $f$  from the subdomain of  $D[[x]]$  whose elements have nonzero constant term to  $D[[x]]$  by  $f(s) = s^{-1} - A(x)$ . Thus  $R(x)$  is just the zero of  $f$ .

Newton's iteration is a second-order iteration<sup>2</sup> and consists of a linear approximation of  $f$ . Newton's iteration function  $\mathcal{N}$  is given by:

$$(1) \quad \mathcal{N}(s) = s - \frac{f(s)}{f'(s)} = s(2 - A(x)s),$$

where  $f'$  denotes the derivative of  $f$ , which is defined algebraically (see [8]). Let  $n$  be a power of two and

$$(2) \quad A_{2n}(x) = A(x) \bmod x^{2n+1},$$

$$(3) \quad R_n(x) = 1/A(x) \bmod x^{n+1}.$$

Newton's iteration features a quadratic convergence (see [3, Chap. 4]): the number of accurate terms doubles at each iteration. This may be expressed by

$$(4) \quad R_{2n}(x) = \mathcal{N}(R_n(x)) \bmod x^{2n+1}.$$

From (2) and (3), there exists a polynomial  $P$  of degree at most  $n - 1$  such that

$$(5) \quad R_n(x)A_{2n}(x) = 1 + x^{n+1}P(x) \bmod x^{2n+1}.$$

Combining (1), (5) and the expansion (4) leads to a recursive formula that computes the reciprocal of  $A(x)$  modulo  $x^{2n+1}$ :

$$(6) \quad \frac{1}{A(x)} = R_{2n}(x) = R_n(x)(1 - x^{n+1}P(x)) \bmod x^{2n+1}.$$

Equations (5) and (6) both charge  $M(n) + O(n)$  units of cost. Therefore, the overall arithmetic complexity of Kung's algorithm is bounded by

$$(7) \quad L(2n) \leq L(n) + 2M(n) + O(n).$$

Unfolding this recurrence leads to  $L(n) = O(M(n))$  for all known multiplications.

The derivation of the exact arithmetic complexity from (7) depends on a specific algorithm for multiplication of polynomials. The next section describes a multiplication algorithm involving fast Fourier transform (FFT). Originally, Kung derived (7) for nonscalar complexity, where  $M(n) = 2n + 1$ , and found  $C(n) < 4n$ . Actually, the lowest upper bound presently known for the nonscalar complexity is  $C(n) < 3.75n$ . Kalorkoti derived this latter result from Kung's third-order iteration [4] and taking advantage that squaring modulo  $x^{n+1}$  is less expensive than multiplying modulo  $x^{n+1}$  (see [2, Chap. 2]).

**3.2. FFT and fast multiplication.** The  $N$ -point FFT defines a ring isomorphism from the quotient  $F[[x]]/(x^N)$  to  $F^N$ . It is an evaluation-interpolation map where the evaluation points, also called Fourier points, are the  $N$ th roots of unity. Actually, the FFT is the evaluation-interpolation map whose implementation yields the lowest known complexity. Indeed, the symmetry properties of the  $N$ th roots of unity allow a divide-and-conquer implementation [3, Chap. 4]. The arithmetic complexity of  $N$ -point FFT is bounded by  $\Lambda(N) \leq 3/2N \log N - N + 1$  (see [2, Chap. 2]).

The FFT performs fast back and forth conversions from an evaluated form to its interpolated form. Thus, low complexity algorithms can be achieved by taking advantage of each representation. In particular, fast multiplication consists in converting both operands into their evaluation forms with two FFTs, performing a coefficient-wise multiplication, and delivering the result with one

---

<sup>2</sup>Third-order iteration is mentioned later and consists of a parabolic approximation.

backward FFT. Schönhage shows that multiplication of polynomials of degree  $n$  (some restrictions on  $n$  are needed and discussed later) according to this method has algebraic complexity

$$(8) \quad M(n) = (9 + o(1)) n \log n.$$

**3.3. Kung's algorithm revisited.** Direct substitution of (8) into (7) leads to

$$L(n) \leq (18 + o(1)) n \log n.$$

However, Schönhage obtains a lower multiplicative constant by deferring the last backward FFT.  $R_n$  and  $A_{2n}$  are first converted into their evaluation forms, requiring two direct  $N$ -point FFTs, which cost  $2\Lambda(N)$ . Then, steps (5) and (6) compute the evaluation form of  $R_n P$ , involving two coefficient-wise multiplications and two subtractions, which add  $4N$  units of cost. One ultimate backward  $N$ -point FFT interpolates  $R_n P$  with  $\Lambda(N)$  operations. Therefore, (7) becomes

$$L(2n) \leq L(n) + 3\Lambda(N) + 4N.$$

A typical value for  $N$  is the lowest power of two that is greater than  $d = \deg(R_n(x)A_{2n}(x)) = 3n$ . However, a significant overhead is expected when  $d$  is slightly greater than the nearest power of two. In this case, the arithmetic complexity for the  $N$ -point FFT is  $\Lambda(N) < 3d \log(2d)$ . Thus, Schönhage suggests for  $N$  a scaled power of two of the form  $N = c 2^\nu$ , where  $\nu = \lceil \log(d) \rceil - \lfloor \log \log(d+1) \rfloor$  and  $c = \lceil d/2^\nu \rceil$ . This latter choice for  $N$  yields a lower bound

$$\Lambda(N) \leq d(3/2 \log(d) + 13/5 \log \log(d+1) + O(1)).$$

This precise count yields the arithmetic complexity for reciprocals

$$L(n) \leq (27/2 + o(1)) n \log n.$$

Surprisingly, Newton's third-order iteration does not yield a better bound for arithmetic complexity, as opposed to the case of nonscalar complexity (see Section 3.1).

#### 4. A New Algorithm over $\mathbb{Z}$

Algorithms for division of polynomials reduce the division task to multiplications. However, while featuring an attractive asymptotic complexity, such reductions may involve detours and tricks whose implementations lead to tremendous multiplicative constants. Indeed, earlier algorithms for division of polynomials shared this drawback. Therefore, Schönhage suggests a new fast algorithm by means of Graeffe's root squaring with a low constant and ready for an immediate implementation due to its extreme simplicity.

**4.1. Graeffe's root squaring method.** Graeffe's squaring method originates in numerical analysis for solving polynomial equations [1]. This method proceeds from any polynomial  $A(x)$  in  $F[x]$  to the even polynomial  $G(x^2) = A(x)A(-x)$ .

In  $F[[x]]$  the reciprocal of  $A(x)$  modulo  $x^{n+1}$  may be written as

$$(9) \quad \frac{1}{A(x)} = \frac{A(-x)}{A(-x)A(x)} \bmod x^{n+1}.$$

In equation (9), the denominator of the right hand-side contains at most  $n+1$  terms, but only half of them are significant when computing modulo  $x^{n+1}$ . Therefore, Graeffe's rule reduces the task of inverting  $n+1$  terms to a half-sized problem. Thus, the corresponding algorithm works recursively as follows (notations are those of (2) and (3)). With  $k = \lfloor n/2 \rfloor$ , Graeffe's step computes

$$G_k(x^2) = A_n(x)A_n(-x) \bmod x^{n+1},$$

charging at most  $n + 1$  nonscalar units of cost. Indeed, typically, nonscalar complexity for such a multiplication is  $C(n) = 2n + 1$  (see [2, Chap. 2]). However, the polynomial  $A_n$  may be rewritten as

$$A_n(x) = A_n^{(\text{even})}(x^2) + xA_n^{(\text{odd})}(x^2),$$

which shows that both  $A_n(x_0)$  and  $A_n(-x_0)$ , for any  $x_0$  lying in the ground field, can be computed together as follows

$$A_n(\pm x_0) = A_n^{(\text{even})}(x_0^2) \pm x_0 A_n^{(\text{odd})}(x_0^2).$$

Therefore, Graeffe's step requires at most  $n + 1$  essential multiplications, by evaluation of  $A_n$  for  $n + 1$  distinct squares. The reciprocal of  $G_k(x)$  modulo  $x^{k+1}$ , denoted by  $H_k(x)$ , is determined by recursive calls. An ultimate multiplication

$$R_n(x) = A_n(-x)H_k(x^2) \bmod x^{n+1}$$

delivers the result, charging extra  $n + 2k + 1$  units of nonscalar cost. Then, the nonscalar complexity is bounded by  $C(n) \leq 6n + 2 \log(n/2)$ , which is slightly weaker than Kalorkoti's (see Section 3.1) but the implementation of Graeffe's approach is straightforward.

**4.2. Application to reciprocals over  $\mathbb{Z}$ .** This section deals with units of the ring  $\mathbb{Z}[[x]]$  of the form  $A(x) = 1 + \sum_{i>0} a_i x^i$ . This form naturally arises with divisions by monic polynomials computed via the substitution  $x \mapsto 1/x$ .

Basically, the implementation of Graeffe's method consists in mapping polynomials to integers expressed in some radix  $r_0$  notation, so that multiplication of integers can be used. This idea is based on Kronecker's trick of encoding polynomials with bounded coefficients in a single integer. Let  $\phi_{r_0}$  be a ring morphism from  $\mathbb{Z}_n[x]$  (i.e., polynomials of  $\mathbb{Z}[x]$  of degree less than  $n$ ) to  $\mathbb{Z}$  that evaluates polynomials at  $r_0 \in \mathbb{N}$ . If there exists a constant  $\beta$  such that  $|a_i| < \beta^i$  holds for each  $i > 0$ , then the bit size of the coefficients of  $R$  and  $G$  can be bounded. Thus, under this assumption,  $r_0 \in \mathbb{N}$  can be chosen such that the evaluation map  $\phi_{r_0}$  is a bijection and  $N$  can be optimally determined. The arithmetic complexity can easily be derived

$$L(n) = 6M(\tau n^2),$$

where  $\tau = \log(3\beta)$  and where the Schönhage–Strassen algorithm for multiplication of integers, which features the lowest known complexity  $M(m) = O(m \log(m) \log \log(m))$  [7], is likely to be used.

### Bibliography

- [1] Bareiss (Erwin H.). – Resultant procedure and the mechanization of the Graeffe process. *Journal of the ACM*, vol. 7, 1960, pp. 346–386.
- [2] Bürgisser (Peter), Clausen (Michael), and Shokrollahi (M. Amin). – *Algebraic complexity theory*. – Springer-Verlag, Berlin, 1997, xxiv+618p. With the collaboration of Thomas Lickteig.
- [3] Geddes (K. O.), Czapor (S. R.), and Labahn (G.). – *Algorithms for computer algebra*. – Kluwer Academic Publishers, Boston, MA, 1992, xxii+585p.
- [4] Kalorkoti (K.). – Inverting polynomials and formal power series. *SIAM Journal on Computing*, vol. 22, n° 3, 1993, pp. 552–559.
- [5] Kung (H. T.). – On computing reciprocals of power series. *Numerische Mathematik*, vol. 22, 1974, pp. 341–348.
- [6] Salvy (Bruno). – *Asymptotique automatique*. – Research Report n° 3707, Institut National de Recherche en Informatique et en Automatique, 1999. 20 pages.
- [7] Schönhage (A.) and Strassen (V.). – Schnelle Multiplikation grosser Zahlen. *Computing (Arch. Elektron. Rechnen)*, vol. 7, 1971, pp. 281–292.
- [8] van der Waerden (B. L.). – *Modern algebra*. – Frederick Ungar Publishing Co., New York, N. Y., 1949, vol. I, xii+264p.

# Fast Multivariate Power Series Multiplication in Characteristic Zero

Grégoire Lecerf

GAGE, École polytechnique (France)

June 11, 2001

Summary by Ludovic Meunier

## Abstract

Let  $S$  be a multivariate power series ring over a field of characteristic zero. The article [5] presents an asymptotically fast algorithm for multiplying two elements of  $S$  truncated according to *total* degree. Up to logarithmic factors, the complexity of the algorithm is optimal, in the sense that it is linear in the size of the output.

## 1. Introduction

Let  $k$  be a field of characteristic zero. We write  $S = k[[x_1, \dots, x_n]]$  for the multivariate power series ring in the  $n$  variables  $x_1, \dots, x_n$ . Let  $\mathfrak{J}$  be any ideal of  $S$ . By computing at *precision*  $\mathfrak{J}$  in  $S$ , we understand computing modulo the ideal  $\mathfrak{J}$  in  $S$ . In other words, power series in  $S$  are regarded as vectors in the  $k$ -algebra  $S/\mathfrak{J}$ . We denote by  $\mathfrak{m}$  the maximal ideal  $(x_1, \dots, x_n)$  in  $S$  and by  $d$  any positive integer. The paper [5] sets the problem of a fast algorithm for multiplying two power series in  $S$  truncated in *total* degree  $d$ , that is computed at precision  $\mathfrak{m}^{d+1}$ .

The general question of a fast algorithm for multivariate multiplication in  $S$  modulo *any* ideal remains an open problem and has received very little attention in the literature. Previous works (e.g., [2]) investigated computation modulo the ideal  $(x_1^{d+1}, \dots, x_n^{d+1})$ , that is truncation according to *partial* degree with respect to each variable  $x_i$ . The method used is called Kronecker's substitution and is briefly discussed in Section 3.

The need for multiplication routines modulo  $\mathfrak{m}^{d+1}$  arises in various fields, such as polynomial system solving [7] and treatment of systems of partial differential equations.

The efficiency of the algorithm is measured with respect to the model of *nonscalar* complexity. By nonscalar complexity, we understand the number of primitive operations in the field  $k$  needed to complete the algorithm, independently of the sizes of the numbers involved (see [3]). We now introduce some notation. We denote by  $D = \deg(\mathfrak{m}^{d+1})$  the *degree* of the ideal  $\mathfrak{m}^{d+1}$ .  $D$  is the number of monomials in  $S$  which are not in  $\mathfrak{m}^{d+1}$ , that is the dimension of the  $k$ -algebra  $S/\mathfrak{m}^{d+1}$ . Simple combinatorial considerations give

$$D = \deg(\mathfrak{m}^{d+1}) = \binom{d+n}{n}.$$

We set  $C := \deg(\mathfrak{m}^d)$  and denote by  $\mathcal{M}_u(\delta)$  the complexity of the multiplication of two univariate polynomials of degree  $\delta$  in  $k[t]$ .

The next section presents the algorithm; its complexity belongs to

$$(1) \quad \mathcal{O}(D \log^3 D \log \log D).$$

Since  $D$  is the size of the output, the algorithm is optimal, up to the logarithmic factors.

## 2. The Algorithm

**2.1. Description.** The first step of the algorithm consists in translating the multivariate problem into a univariate one. This is motivated by the fact that fast algorithms for univariate power series multiplication are known (e.g., [6]).

Let  $t$  be a new variable. We consider the substitution

$$\begin{aligned} \tilde{\mathcal{R}}_t : S/\mathfrak{m}^{d+1} &\longrightarrow k[x_1, \dots, x_n][[t]]/(t^{d+1}) \\ f(x_1, \dots, x_n) &\longmapsto f(x_1t, \dots, x_nt). \end{aligned}$$

If  $f$  is an element of  $S/\mathfrak{m}^{d+1}$ ,  $\tilde{\mathcal{R}}_t(f)$  is a univariate power series in the *single* variable  $t$  truncated at degree  $d$ . It can then be written  $\tilde{\mathcal{R}}_t(f) = \tilde{f}_0 + \tilde{f}_1t + \dots + \tilde{f}_dt^d$ , where each coefficient  $\tilde{f}_i$  is a homogeneous multivariate polynomial in the variables  $x_1, \dots, x_n$  of *total* degree  $i$ . This remark on the degree suggests that:

1. the substitution  $\tilde{\mathcal{R}}_t$  is optimal, in the sense that it provides us with a representation of  $f$  that retains exactly the monomials that form a basis of  $S/\mathfrak{m}^{d+1}$ . In particular, the algorithm does not suffer from any overhead caused by unnecessary terms (see Section 3);
2. in view of the homogeneity of the  $\tilde{f}_i$ , keeping all of the variables  $x_i$  is redundant. The substitution defined by

$$\begin{aligned} \mathcal{R}_t : S/\mathfrak{m}^{d+1} &\longrightarrow k[x_2, \dots, x_n][[t]]/(t^{d+1}) = (k[[t]]/(t^{d+1})) [x_2, \dots, x_n] \\ f(x_1, \dots, x_n) &\longmapsto f(t, x_2t, \dots, x_nt) \end{aligned}$$

reduces the complexity in the step of *evaluation-interpolation* (see below):  $n - 1$  variables, instead of  $n$  variables, are actually needed.

The second step of the algorithm performs the multiplication. Let  $f$  and  $g$  be two power series in  $S/\mathfrak{m}^{d+1}$  and  $h$  be the product  $fg$  in  $S/\mathfrak{m}^{d+1}$ . The equality  $h = fg$  turns into

$$(2) \quad \mathcal{R}_t(h) = \mathcal{R}_t(f)\mathcal{R}_t(g).$$

Consequently, we concentrate on a fast way to compute  $\mathcal{R}_t(h)$ . We use an evaluation-interpolation scheme. We first consider the *evaluation* map at the point  $P = (p_2, \dots, p_n)$  in  $k^{n-1}$  defined by

$$\begin{aligned} \mathcal{E}_P : (k[[t]]/(t^{d+1})) [x_2, \dots, x_n] &\longrightarrow k[[t]]/(t^{d+1}) \\ f(x_2, \dots, x_n) &\longmapsto f(P). \end{aligned}$$

We then apply  $\mathcal{E}_P$  to equation (2), which yields

$$(3) \quad \mathcal{E}_P(\mathcal{R}_t(h)) = \mathcal{E}_P(\mathcal{R}_t(f))\mathcal{E}_P(\mathcal{R}_t(g)) \pmod{t^{d+1}}.$$

Equation (3) holds for any point  $P$  and computes the product  $\mathcal{R}_t(h)$  at  $P$  by using a *univariate* power series multiplication algorithm. Such an algorithm is described in [6].

The last step of the algorithm consists in reconstructing  $h$  from a set of values of  $\mathcal{R}_t(h)$ . We regard  $\mathcal{R}_t(h)$  as a multivariate polynomial in the variables  $x_2, \dots, x_n$ . There exists an *interpolation* map

$$\begin{aligned} \mathcal{I} : (k[[t]]/(t^{d+1}))^C &\longrightarrow (k[[t]]/(t^{d+1})) [x_2, \dots, x_n] \\ (f(P_1), \dots, f(P_C)) &\longmapsto f(x_2, \dots, x_n), \end{aligned}$$

which recovers  $\mathcal{R}_t(h)$  from a set of  $C$  pairwise distinct values  $\{\mathcal{E}_{P_1}(\mathcal{R}_t(h)), \dots, \mathcal{E}_{P_C}(\mathcal{R}_t(h))\}$ . The evaluation points  $P_i$ , for  $i$  in  $1, \dots, C$ , are chosen to be powers of distinct prime numbers, namely  $P_i = (p_2^i, \dots, p_n^i)$ , where  $p_j$  are distinct prime numbers. Note the key point is that the characteristic of the ground field  $k$  is zero, so that all  $\mathcal{E}_{P_i}(\mathcal{R}_t(h))$  have pairwise distinct values. An implementation

of both maps  $\mathcal{E}_P$  and  $\mathcal{I}$  is described by J. Canny, E. Kaltofen, and Y. Lakshman in [4]. Their method relies on fast univariate multipoint evaluation and interpolation (e.g., [1]).

Finally, we reconstruct  $h$  from  $\mathcal{R}_t(h)$ . If  $\mathcal{R}_t(h) = h_0 + h_1t + \cdots + h_d t^d$  is given,  $h$  is obtained by homogenizing each  $h_i$  in degree  $i$  with respect to the variable  $x_1$  and then evaluating at  $t = 1$ .

We are now ready to unfold the algorithm.

```

MultivariatePS_Mult := proc(f,g)
  (1)   $F \leftarrow \mathcal{R}_t(f); G \leftarrow \mathcal{R}_t(g);$            // new representation
  (2)  for  $i$  in  $(P_1, \dots, P_C)$  do                     // evaluation
         $F_{P_i} \leftarrow \mathcal{E}_{P_i}(F); G_{P_i} \leftarrow \mathcal{E}_{P_i}(G);$ 
  (3)  for  $i$  to  $C$  do                                       // univariate multiplication
         $H_{P_i} \leftarrow F_{P_i} G_{P_i};$ 
  (4)   $\mathcal{R}_t(h) \leftarrow \mathcal{I}(H_{P_1}, \dots, H_{P_C});$            // interpolation
  (5)   $h \leftarrow$  homogenization in degree with respect to  $x_1$  // reconstruction
        in  $\mathcal{R}_t(h);$ 
  return  $h;$ 

```

The next section derives the complexity result claimed by (1).

**2.2. Complexity.** Steps 1 and 5 can be performed in  $\mathcal{O}(C)$  operations. We examine the cost of Steps 2, 3, and 4 separately:

- Step 2 evaluates the  $d$  coefficients of  $F$  and  $G$  at  $C$  points. The  $C$  points  $P_i$  are chosen to be powers of the  $n - 1$  distinct prime numbers  $(p_2, \dots, p_n)$ , namely  $P_i = (p_2^i, \dots, p_n^i)$ . Each coefficient can be computed in  $\mathcal{O}(\mathcal{M}_u(C) \log C)$  operations, according to the algorithm for fast multipoint evaluation given in [4]. This yields an overall complexity of  $\mathcal{O}(d\mathcal{M}_u(C) \log C)$  for Step 2.
- Step 3 performs  $C$  univariate power series products. Each multiplication requires  $\mathcal{O}(\mathcal{M}_u(d))$  operations. Complexity of Step 3 is then  $\mathcal{O}(C\mathcal{M}_u(d))$ .
- Step 4 interpolates the  $d$  coefficients of  $H$ . Each interpolation requires  $\mathcal{O}(\mathcal{M}_u(C) \log C)$  operations, also using the algorithm presented in [4]. Step 4 then requires  $\mathcal{O}(d\mathcal{M}_u(C) \log C)$  operations.

The overall complexity of the algorithm is then derived by replacing  $\mathcal{M}_u(C)$  by its estimate  $\mathcal{O}(C \log C \log \log C)$  obtained in [6] and noting that  $C < D \log(D)/d$ . This yields

$$\mathcal{O}(D \log^3 D \log \log D).$$

**2.3. Generalization.** We mention that van der Hoeven generalized the algorithm to the case when

$$\mathfrak{J} = (x_1^{d_1} \dots x_n^{d_n}, \quad \text{for } \alpha_1 d_1 + \cdots + \alpha_n d_n > d),$$

where the  $\alpha_i$  are positive integers, by using the substitution defined by

$$\begin{aligned} \mathcal{V}_t: S/\mathfrak{J} &\longrightarrow k[x_2, \dots, x_n][[t]]/(t^{d+1}) \\ f(x_1, \dots, x_n) &\longmapsto f(t^{\alpha_1}, x_2 t^{\alpha_2}, \dots, x_n t^{\alpha_n}) \end{aligned}$$

instead of  $\mathcal{R}_t$ . The rest of the algorithm remains unaltered.

### 3. Appendix: Kronecker's Substitution

Kronecker's substitution is defined by the map

$$\begin{aligned} \mathcal{K}_t : S/\mathfrak{J} &\longrightarrow k[[t]]/t^{(2d+1)^n} \\ f(x_1, \dots, x_n) &\longmapsto f(t, t^{2d+1}, \dots, t^{(2d+1)^{n-1}}), \end{aligned}$$

where  $\mathfrak{J} = (x_1^{d+1}, \dots, x_n^{d+1})$ . This substitution truncates power series in *partial* degree  $d$  with respect to each variable  $x_i$ . Let  $f$  be a power series in  $S/\mathfrak{J}$ , one recovers the coefficient of  $x_1^{e_1} \dots x_n^{e_n}$  in  $f$  by simply reading off the coefficient of  $t^{e_1 + (2d+1)e_2 + \dots + (2d+1)^{n-1}e_n}$  in  $\mathcal{K}_t(f)$ . The cost of this algorithm is the cost of the multiplication of two univariate polynomials of degree  $(2d)^n$ , that is  $\mathcal{O}(\mathcal{M}_u((2d)^n))$ . This is the lowest known complexity for multivariate power series multiplication modulo the ideal  $(x_1^{d+1}, \dots, x_n^{d+1})$ . In particular, when addressed in this context, the algorithm presented above requires precision  $\mathfrak{m}^{nd+1}$  and yields a similar complexity.

Kronecker's substitution may be used to compute modulo  $\mathfrak{m}^{d+1}$  as well. However, it results in a significant overhead of  $\mathcal{O}(2^n n!)$ , for fixed  $n$  and  $d \gg n$ , with respect to the size of the power series.

#### Bibliography

- [1] Aho (Alfred V.), Hopcroft (John E.), and Ullman (Jeffrey D.). – *The design and analysis of computer algorithms*. – Addison-Wesley Publishing Co., Reading, Mass.-London-Amsterdam, 1975, x+470p. Second printing, Addison-Wesley Series in Computer Science and Information Processing.
- [2] Brent (R. P.) and Kung (H. T.). – Fast algorithms for composition and reversion of multivariate power series (preliminary version). In *Proceedings of a Conference on Theoretical Computer Science Department of Computer Science, University of Waterloo, Waterloo, Ontario (August 1977)*, pp. 149–158. – 1977.
- [3] Bürgisser (Peter), Clausen (Michael), and Shokrollahi (M. Amin). – *Algebraic complexity theory*. – Springer-Verlag, Berlin, 1997, xxiv+618p. With the collaboration of Thomas Lickteig.
- [4] Canny (John F.), Kaltofen (Erich), and Lakshman Yagati. – Solving systems of nonlinear polynomial equations faster. In Gonnet (Gaston) (editor), *Symbolic and Algebraic Computation (International Symposium ISSAC'89, Portland, Oregon, USA, July 17-19, 1989)*, pp. 121–128. – ACM Press, 1989. Conference proceedings.
- [5] Lecerf (Grégoire) and Schost (Éric). – Fast multivariate power series multiplication in characteristic zero. – Available from <http://www.medicis.polytechnique.fr/~schost/>, 2001.
- [6] Schönhage (A.). – Schnelle Multiplikation von Polynomen über Körpern der Charakteristik 2. *Acta Informatica*, vol. 7, n° 4, 1976/77, pp. 395–398.
- [7] Schost (Éric). – *Sur la résolution des systèmes polynomiaux à paramètres*. – PhD thesis, École polytechnique, Palaiseau, France, January 2001. Defended on December 7, 2000.



## A Tutorial on Closed Difference Forms

Burkhard Zimmermann

RISC, Linz (Austria)

January 15, 2001

Summary by Frédéric Chyzak

### Abstract

Zeilberger’s theory of closed difference forms provides with a deeper understanding of the creative telescoping method used to prove many ( $q$ -)hypergeometric (multi-)sum identities, and of “companion” or “dual” identities. By introducing new types of summation domains, the closed form approach allows to discover new identities of the form “sum equals sum,” including new summatory representations of  $\zeta(3)$ . A transform similar to a pullback (change of variables) of differential forms is introduced, and permits to find more new identities. This summary is freely inspired by [1, 2, 4, 5] and the talk.

### 1. Comparison Between Differential and Difference Calculi

By mimicking differential calculus [2], Zeilberger has developed a complete *difference calculus* [4]. This theory, which we recall here, culminates with a discrete analogue to Stokes’s theorem.

Given a  $\mathbb{C}$ -vector space  $V$ , which will take the role of a tangent space momentarily, an *alternate multilinear  $p$ -form* on  $V$  is just a multilinear map  $\phi : V^p \rightarrow \mathbb{C}$  that satisfies the rule

$$\phi(v_1, \dots, v_{i+1}, v_i, \dots, v_p) = -\phi(v_1, \dots, v_p).$$

This represents a  $p$ -volume measure, in the sense that it assigns an (oriented) volume to the parallelepipedic polyhedron determined by the vectors  $v_i$ . By a natural convention, 0-forms are just constants. To a  $p$ -form  $\phi$  and a  $q$ -form  $\psi$ , one associates a  $(p+q)$ -form, i.e., a  $(p+q)$ -volume measure, by means of the *exterior product*  $\phi \wedge \psi$ :

$$(\phi \wedge \psi)(v_1, \dots, v_{p+q}) = \sum_{\sigma \in S_{p,q}} \epsilon(\sigma) \phi(v_{\sigma(1)}, \dots, v_{\sigma(p)}) \psi(v_{\sigma(p+1)}, \dots, v_{\sigma(p+q)})$$

where  $S_{p,q}$  denotes the set of permutations of  $\{1, \dots, p+q\}$  with  $\sigma(1) < \dots < \sigma(p)$  and  $\sigma(p+1) < \dots < \sigma(p+q)$ , and where  $\epsilon(\sigma)$  denotes the signature of the permutation  $\sigma$ . Consider the direct sum  $\mathcal{A}(V) = \bigoplus_{p \geq 0} \mathcal{A}_p(V)$  of the vector spaces  $\mathcal{A}_p(V)$  of alternate  $p$ -forms. By extending the exterior product by linearity, we obtain an associative multiplication on  $\mathcal{A}(V)$ , which becomes a graded algebra with the product rule  $\psi \wedge \phi = (-1)^{pq} \phi \wedge \psi$  for a  $p$ -form  $\phi$  and a  $q$ -form  $\psi$ .

Next, an *alternate difference  $p$ -form*, or for short a *difference  $p$ -form*, is a map  $\omega$  which to each element  $\xi$  of a real manifold  $M$  associates a multilinear  $p$ -form  $\omega(\xi)$  on the tangent space  $V = T_\xi M$ . Exterior products of difference forms are defined pointwise. At this point, difference forms and differential forms share the same definition. In the following however, we focus to the case when  $M$  is a submanifold of  $\mathbb{R}^d$ : each  $\omega(\xi)$  is then an alternate form on  $V = \mathbb{R}^d$ . By imposing the additional property  $\omega(\xi_1, \dots, \xi_d) = \omega([\xi_1], \dots, [\xi_d])$ , we obtain forms that are piecewise constant, as well as their coefficients. (Compare this situation with the theory in the differential setting,

where one insists in having  $C^\infty$  forms and  $C^\infty$  coefficients.) The possible variations of forms with  $\xi$  is at the origin of the notions of *exterior differential* and *exterior difference* introduced below.

In the differential setting, a kind of a derivation is defined on differential forms in the following way. One starts with the usual derivative  $\omega'$ , which satisfies the asymptotic relation  $\omega(\xi + v) = \omega(\xi) + \omega'(\xi)(v) + o(v)$  as  $v \rightarrow 0$ . Each  $\omega'(\xi)$  is a linear map from  $V = \mathbb{R}^d$  to the vector space  $\mathcal{A}_p(V)$ , and can be viewed as a multilinear map from  $V^{p+1}$  to  $\mathbb{C}$  that is not alternate, but alternate in its last  $p$  variables only. Making it alternate by an averaging technique, we obtain the *exterior differential*  $d\omega$  given by

$$(d\omega)(\xi)(v_0, \dots, v_p) = \sum_{i=0}^p (-1)^i (\omega'(\xi)(v_i))(v_0, \dots, \hat{v}_i, \dots, v_p).$$

In the difference case, we start with another linearization instead of the derivative  $\omega'$  to define the exterior difference of  $\omega$ , namely by secants instead of tangents. Let  $\omega^\Delta(\xi)$  be the linear map on  $V$  defined by  $\omega(\xi + v) = \omega(\xi) + \omega^\Delta(\xi)(v) + R(v)$  and  $R(v)$  is zero for each element  $v = e_i$  of the canonical basis of  $V = \mathbb{R}^d$ . Again,  $(v_0, \dots, v_p) \mapsto \omega^\Delta(\xi)(v_0)(v_1, \dots, v_p)$  is alternate in its last  $p$  variables only, but the full alternate nature is recovered by the *exterior difference*  $d\omega$  defined by

$$(d\omega)(\xi)(v_0, \dots, v_p) = \sum_{i=0}^p (-1)^i (\omega^\Delta(\xi)(v_i))(v_0, \dots, \hat{v}_i, \dots, v_p).$$

As opposed to the classical exterior differential, exterior difference heavily depends on the choice of a basis on  $V$ ; but like it, it satisfies  $d \circ d = 0$ .

Denote  $(n_1, \dots, n_d)$  the dual basis of the canonical basis of the manifold  $\mathbb{R}^d$  that contains  $M$ . As in the differential setting, the exterior difference  $dn_i$  of the restriction of  $n_i$  to  $M$  (i.e., or the  $i$ th coordinate function on  $M$ ) plays a special role: the  $dn_i$  form a basis for the ring of difference form, and the  $dn_{i_1} \wedge \dots \wedge dn_{i_p}$  for  $i_1 < \dots < i_p$  span the vector space (respectively, free module) of  $p$ -forms. Exterior differential and exterior difference share a formally simple, easy-to-memorize formulation on the canonical basis  $(dn_1, \dots, dn_d)$ : for  $\omega = f dn_{i_1} \wedge \dots \wedge dn_{i_r}$ , we get

$$d\omega = df \wedge dn_{i_1} \wedge \dots \wedge dn_{i_r}$$

where the exterior differential is  $df = \sum_{i=1}^d \frac{\partial f}{\partial \xi_i} dn_i$ , and the exterior difference  $df = \sum_{i=1}^d (\Delta_i f) dn_i$ , where  $\Delta_i$  is the finite difference operator defined by  $(\Delta_i f)(\xi_1, \dots, \xi_d) = f(\xi_1, \dots, \xi_i + 1, \dots, \xi_d) - f(\xi_1, \dots, \xi_d)$ .

In order to make the link between difference forms and summation, we restrict to hypercubic manifolds given by setting some of the coordinates  $\xi_i$  to 0 and letting all others vary freely in  $[0, 1)$ , and to the manifolds obtained after translating the latter by vectors with integer entries. Note that all those elementary manifolds (in various dimensions) have volume 1, and that we have restricted difference forms to be constant on such sets. As a consequence, the integral of a form  $f dn_1 \wedge \dots \wedge dn_d$  on  $[0, 1)^d$  is just  $f(0, \dots, 0)$ , as is for  $i_1 < \dots < i_r$  the integral of  $f dn_{i_1} \wedge \dots \wedge dn_{i_r}$  on the hypercube defined by  $0 \leq \xi_j < 1$  for each  $j = i_k$  and  $\xi_j = 0$  for all other  $j$ . By integration over a union of elementary manifolds, we are naturally led to integral representing sums; for example:

$$\int_{\mathbb{R}^d} f dn_1 \wedge \dots \wedge dn_d = \sum_{(n_1, \dots, n_d) \in \mathbb{Z}^d} f(n_1, \dots, n_d).$$

We are now ready to derive a difference variant of Stokes's theorem: consider the oriented hypercube  $\Omega = [0, 1)^d$  and its boundary  $\partial\Omega$  defined as usual as a formal linear combination of  $2d$  faces,

$$\partial\Omega = F(\xi_1 = 0) - F(\xi_2 = 0) + \dots + (-1)^{d+1} F(\xi_d = 0) - F(\xi_1 = 1) + F(\xi_2 = 1) + \dots + (-1)^d F(\xi_d = 1),$$

where  $F(\xi_i = a)$  is the (oriented) face  $\Omega \cap \{\xi \mid \xi_i = a\}$ . Boundaries of other elementary manifolds are obtained by translating  $\partial\Omega$ , keeping the same coefficients. In this way, we can define the integral of a form over a linear combination of manifolds to be the very same linear combination of integrals of the same form over the manifolds. For

$$(1) \quad \omega = \sum_{i=1}^d f_i dn_1 \wedge \cdots \wedge \hat{dn}_i \wedge \cdots \wedge dn_d$$

we get

$$\begin{aligned} \int_{\partial\Omega} \omega &= \sum_{i=1}^d (-1)^i \int_{F(\xi_i=1)-F(\xi_i=0)} f_i dn_1 \wedge \cdots \wedge \hat{dn}_i \wedge \cdots \wedge dn_d \\ &= \left( \sum_{i=1}^d (-1)^i f_i(0, \dots, 1, \dots, 0) - \sum_{i=1}^d (-1)^i f_i(0, \dots, 0) \right) dn_1 \wedge \cdots \wedge dn_d \\ &= \sum_{i=1}^d (-1)^i (\Delta_i f_i)(0, \dots, 0) dn_1 \wedge \cdots \wedge dn_d = \int_{\Omega} d\omega. \end{aligned}$$

We could have as well considered forms  $\omega$  defined on the integer lattice  $\mathbb{Z}^d$ , and defined their sums  $\sum_{\Omega} \omega$  on a manifold  $\Omega$  by the integrals  $\int_{\Omega} \omega$  of the form  $\omega$  extended to  $\mathbb{R}^d$  by  $\omega(\xi_1, \dots, \xi_d) = \omega([\xi_1], \dots, [\xi_d])$ . We shall adopt this equivalent viewpoint from the next section on. By linearity with respect to manifolds, we obtain the following discrete variant of Stokes's formula [4].

**Theorem 1** (Zeilberger–Stokes formula). *For any difference  $p$ -form  $\omega$  such that  $\omega(\xi_1, \dots, \xi_d) = \omega([\xi_1], \dots, [\xi_d])$  on any manifold  $\Omega$  that is a linear combination of elementary hypercubic manifolds, we have  $\sum_{\partial\Omega} \omega = \sum_{\Omega} d\omega$ .*

## 2. Closed Form Identities (Pun Intended!)

An interesting situation is that of a *closed (difference) form*, which by definition is a difference form  $\omega$  such that  $d\omega = 0$ . In this case, the sum  $\sum_{\partial\Omega} \omega = 0$  for any manifold  $\Omega$  on all of which  $\Omega$  is defined, owing to Theorem 1 above. If more specifically  $\omega$  is given by (1), we obtain

$$\sum_{i=1}^d \sum_{\partial\Omega} f_i dn_1 \wedge \cdots \wedge \hat{dn}_i \wedge \cdots \wedge dn_d = 0,$$

in other words a relation between a priori infinite sums! Using the leeway available in the choice of  $\Omega$  yields several kinds of identities: sum equals constant, sum equals sum, etc. In the following, we detail this situation in the special case  $r = 2$ . Let us denote  $dn$  and  $dk$  for  $dn_1$  and  $dn_2$ , respectively, and consider a closed 1-form  $\omega = g dn + f dk$ , so that  $\Delta_n f = \Delta_k g$ .

**2.1. Stripe-shaped manifolds.** Consider  $\Omega = \mathbb{R}^+ \times [0, n] = \{(x, y) \mid x \geq 0 \text{ and } 0 \leq y \leq n\}$  and the closed form  $\omega$  obtained for

$$f(n, k) = \binom{m}{k} \binom{n}{k} \binom{p+n+m-k}{n+m} \quad \text{and} \quad g(n, k) = \frac{mk - p(n+1)}{(n+m+1)(n+1-k)} f(n, k).$$

Stokes's theorem on  $\Omega$  then yields (after elementary manipulations of binomial sums)

$$\sum_{k=0}^n \binom{m}{k} \binom{n}{k} \binom{p+n+m-k}{n+m} = \sum_{k \in \mathbb{N}} f(n, k) = \sum_{k \in \mathbb{N}} f(0, k) + \sum_{l=0}^n g(l, 0) = \binom{m+p}{m} \binom{n+p}{n}.$$

More generally, many *closed-form* identities like the one above, where “closed form” now means that both the summand and the sum are hypergeometric sequences, correspond to a “closed form” that involves the summand as one of its coefficients. Hence Zeilberger’s “pun intended.”

But some magic takes place here: changing  $\Omega$  to  $[0, k] \times \mathbb{R}^+$  and summing with respect to  $n$  instead of  $k$ , the same method sometimes yields a *companion identity*. Moreover, the more variables there are, the more amplified this phenomenon is: for  $r$  variables and in lucky cases where all summations make sense, a single closed difference  $(r - 1)$ -form with hypergeometric coefficients can be viewed as a simultaneous encoding of  $r$  closed form summation identities [4].

**2.2. Triangular-shaped manifolds.** Zeilberger observed that for a closed form  $\omega_1 = g_1 dn + f_1 dk$ , the functions  $f_s(n, k) = f_1(sn, k)$  and  $g_s(n, k) = g_1(sn, k) + g_1(sn + 1, k) + \cdots + g_1(sn + s - 1, k)$  provide for each  $s > 1$  with another closed form  $\omega_s = g_s dn + f_s dk$ . Basing on this, Amdeberhan and Zeilberger [1] derived the following representations for  $\zeta(3)$ :

$$\begin{aligned} \zeta(3) &= \frac{5}{2} \sum_{n=1}^{\infty} \frac{(-1)^{n-1}}{\binom{2n}{n} n^2} = \frac{1}{4} \sum_{n=1}^{\infty} \frac{(-1)^{n-1} (56n^2 - 32n + 5)}{(2n-1)^2 \binom{3n}{n} \binom{2n}{n} n^3} \\ &= \frac{1}{72} \sum_{n=0}^{\infty} \frac{(-1)^n (5265n^4 + 13878n^3 + 13761n^2 + 6120n + 1040)}{(4n+3)(4n+1)(3n+2)^2(3n+1)^2(n+1) \binom{4n}{n} \binom{3n}{n}}. \end{aligned}$$

Specifically, they considered  $\Omega = \{(x, y) \mid y \geq \lfloor x + 1 \rfloor\}$  and the functions

$$f_1(n, k) = (-1)^k \frac{k!^2(n-k-1)!}{(n+k+1)!(k+1)} \quad \text{and} \quad g_1(n, k) = 2(-1)^k \frac{k!^2(n-k)!}{(n+k+1)!(n+1)^2}.$$

The representations above have respectively been obtained for  $s = 1, 2$ , and  $3$ ; their general terms decrease like  $O(n^{-3/2}4^{-n})$ ,  $O(n^{-2}27^{-n})$ ,  $O(n^{-2}64^{-n})$ , respectively—at the cost of more and more operations for each term, though! Changing  $\Omega$  to  $\Omega_s = \{(x, y) \mid y \geq s\lfloor x + 1 \rfloor\}$  leads to other representations [1], like, for  $s = 2$ ,

$$\zeta(3) = \sum_{n=0}^{\infty} \frac{(-1)^n P(n)}{80(5n+4)(5n+3)(5n+2)(5n+1)(4n+3)^2(4n+1)^2(2n+1)^2(n+1) \binom{5n}{n} \binom{4n}{n}}$$

where  $P = 1613824n^8 + 7638016n^7 + 15700096n^6 + 18317312n^5 + 13278552n^4 + 6131676n^3 + 1763967n^2 + 289515n + 20782$ . The general term is now  $O(n^{-2}(27/3125)^{-n})$ , with  $27/3125 \approx 115.74$ .

To sketch the proof, we apply Stokes’s theorem to  $\omega_s$  on  $\Omega_s$ , and obtain:

$$\sum_{n=0}^{\infty} g_s(n, 0) + \sum_{k=0}^{\infty} f_s(sk + s, k) + \sum_{k=0}^{\infty} (g_s(sk, k) + \cdots + g_s(sk + s - 1, k)) = 0.$$

Next, noting that  $g_1(n, 0) = 2/(n+1)^3$  and grouping the sums over  $k$  yields the announced identity.

**2.3. Finite triangular-shaped and rectangular-shaped manifolds.** Other identities like

$$\frac{\Gamma(x+n)\Gamma(y+n)}{\Gamma(n)\Gamma(x+y+n)} {}_3F_2\left(\begin{matrix} x, y, v+n-1 \\ v, x+y+n \end{matrix} \middle| 1\right) = \frac{\Gamma(x+k)\Gamma(y+k)}{\Gamma(k)\Gamma(x+y+k)} {}_3F_2\left(\begin{matrix} x, y, v+k-1 \\ v, x+y+k \end{matrix} \middle| 1\right)$$

and  $\sum_{n+m=s} \binom{2n}{n} \binom{2m}{m} = 4^s$  are based on other choices for  $\Omega$ , like a rectangle  $[0, k] \times [0, n]$  or a “triangle”  $\{(x, y) \mid \lfloor x \rfloor + \lfloor y \rfloor \leq s\}$  for  $\Omega$  [5].

### 3. Closed Forms with Holonomic Coefficients

Consider a closed form  $\omega = g \, dn + f \, dk$  with hypergeometric coefficients. Since  $f$  is hypergeometric in  $n$ , one can find some rational function  $R$  of  $(n, k)$  such that  $\Delta_k g = \Delta_n f = Rf$ . It is also well-known that if a hypergeometric sequences  $h$  has a hypergeometric anti-difference  $H$ , there has to be some rational function  $S$  such that  $H = Sh$ . Here we get  $g = S\Delta_n f = SRf$ . This situation extends to more variables, which legitimates Zeilberger's focus to closed forms whose coefficients are all multiples of the same hypergeometric sequence  $f$  by polynomials in the variables; he called such forms *WZ forms* [4]. Here we extend this situation to forms whose coefficients are rational multiples of the same holonomic sequence, and make the link between closed forms and creative telescoping explicit.

Let a summation identity  $\sum_{k=a}^b f_{n,k} = F_n$  be given, where both  $f$  and  $F$  are holonomic  $\partial$ -finite sequences. In view of verifying it, knowing  $F$  allows to compute a non-zero operator  $P_0(n, S_n)$  such that  $P_0 \cdot F = 0$ . Proving the identity thus reduces to proving  $\sum_{k=a}^b (P_0 \cdot f)(n, k) = 0$ . By restricting to holonomic hypergeometric summands and right-hand sides, Zeilberger's presentation essentially only dealt with the case  $P_0 = S_n - 1$ :  $F$  can always be assumed to be 1, otherwise we replace  $f(n, k)$  with  $f(n, k)/F(n)$ . In this spirit, we now require that  $P_0$  be a right multiple of  $S_n - 1$  and write  $P_0 = (S_n - 1)R$  this factorization.

The holonomy of  $f$  ensures that there exists a pair  $(P, Q)$  with non-zero  $P$  such that

$$(2) \quad (P + (S_k - 1)Q) \cdot f = 0.$$

Provided that there exists such a pair for  $P = P_0$ , the operator  $Q$  can be computed by Chyzak's  $\partial$ -finite extension of Gosper's algorithm [3]. Let  $\mathcal{A}$  be the algebra of difference operators with respect to  $n$  and  $k$  with coefficients that are rational functions in  $n$  and  $k$ , and introduce the module  $\mathfrak{M} = \mathcal{A} \cdot f$ . The form

$$(3) \quad \omega = (R \cdot f) \, dk - (Q \cdot f) \, dn,$$

whose coefficients all lie in  $\mathfrak{M}$  is closed:

$$d\omega = ((S_n - 1)R \cdot f) \, dn \wedge dk - ((S_k - 1)Q \cdot f) \, dk \wedge dn = \left( (P + (S_k - 1)Q) \cdot f \right) \, dn \wedge dk = 0.$$

Conversely, assume that there exists a closed form  $\omega$  (with coefficients in  $\mathfrak{M}$ ) given by (3). By closedness, we have  $((S_n - 1)R + (S_k - 1)Q) \cdot f = 0$ , whence after summation over  $k$ , and provided that  $R$  involves neither  $k$  nor  $S_k$ ,

$$(S_n - 1)R \cdot \sum_{k=a}^b f(n, k) = 0.$$

More generally, if the  $r$ -form  $f \, dk_1 \wedge \dots \wedge dk_r + \sum_{i=1}^r (P_i \cdot f) \, dn \wedge dk_1 \wedge \dots \wedge \hat{dk}_i \wedge \dots \wedge dk_r$  is closed,

$$\text{i.e., } (S_n - 1) \cdot f + (S_{k_1} - 1)P_1 \cdot f + \dots + (S_{k_r} - 1)P_r \cdot f = 0,$$

the  $r$ -fold summation  $\sum_{k_1, \dots, k_r} f$  yields a constant with respect to  $n$ .

### 4. Extended WZ Cohomology

Is it easily shown that any 1-form with coefficients defined on  $\mathbb{Z}^r$  is exact. Even more is true: any 1-form with holonomic coefficients derives from a holonomic sequence. More specifically, a 1-form  $\omega$  given by (3) is exact if and only if there exists a function  $\phi(n, k)$  such that  $\omega = d\phi$ , or more explicitly

$$-(Q \cdot f) = (S_n - 1) \cdot \phi \quad \text{and} \quad R \cdot f = (S_k - 1) \cdot \phi.$$

This always holds if we look for unconstrained  $\phi$ : simply define  $\phi$  by

$$\phi(n, k) = \sum_{i=0}^{k-1} (R \cdot f)(0, i) - \sum_{j=0}^{n-1} (Q \cdot f)(j, k).$$

The non-trivial problem is to impose  $\phi \in \mathfrak{M}$ . (For example, when  $f$  is hypergeometric, all coefficients of  $\omega$  as well as  $\phi$  have to be rational multiples of  $f$ .) Then, not all 1-forms  $\omega$  remain exact. Viewing closed forms modulo exact forms we are led to a cohomology that Zeilberger named *WZ cohomology* in [4] in the case of hypergeometric  $f$ , and that we call *extended WZ cohomology* in the more general case of holonomic  $\partial$ -finite  $f$ . Following Zeilberger [4], we suggest the following extended research problem: characterize those holonomic  $\partial$ -finite sequences  $f$  for which there exists a non-exact closed form with coefficients in  $\mathfrak{M} = \mathcal{A} \cdot f$  and compute the corresponding cohomology.

## 5. Pullbacks

In the differential case, the notion of *pullback* propagates a change of variables in functions to the level of differential forms, thus permitting change of variables in integrals: for a differentiable map  $\phi$  from a manifold  $N$  to another manifold  $M$ , one gets a mapping  $\phi^*$  that transforms a  $p$ -form  $\omega$  on  $M$  to a  $p$ -form on  $N$  while preserving closedness of forms by simply requiring

$$(4) \quad (\phi^* \omega)(\xi)(v_1, \dots, v_p) = \omega(\phi(\xi))(\phi'(\xi)(v_1), \dots, \phi'(\xi)(v_p)).$$

In the difference case, a simple example of a pullback has already been given in Section 2.2: the closed form  $\omega_s$  is the pullback of the closed form  $\omega_1$  under the map given by  $\phi(n, k) = (sn, k)$ . However, no simple definition of a pullback seems possible: the obvious guess that mimicks (4), substituting  $\phi^\Delta$  for  $\phi'$ , unfortunately does not preserve closedness (taking finite differences is not a local operation). Zimmermann [5] and Gessel independently gave a definition for the case of a linear mapping  $\phi$  that maps integer points to integer points.

The key observation is that for a linear transform  $l = \phi(n)$ , defined by  $l_i = \sum_j a_{i,j} n_j$ , shifting by 1 with respect to  $n_j$  after performing the substitution induced by  $\phi$  is equivalent to doing shifts with respect to each  $l_i$  before substituting, as detailed by the formula  $S_{l_j} \phi^* = \phi^* S_{n_1}^{a_{1,j}} \dots S_{n_r}^{a_{r,j}}$ . It then follows from a technical but easy calculation that  $\Delta_{l_j} \phi^* = \phi^* \sum_i P_{i,j} \Delta_{n_i}$  for some operators  $P_{i,j}$ . Imposing the natural relations  $\phi^*(f) = f \circ \phi$  and  $\phi^*(df) = d(\phi^* f)$  for 0-forms  $f$  leads to

$$\sum_i \phi^*((\Delta_{n_i} f) dn_i) = \sum_j (\Delta_{l_j} (\phi^* f)) dl_j = \sum_{i,j} \phi^*(P_{i,j} \Delta_{n_i} f) dl_j.$$

Choosing  $f$  such that  $df = (\Delta_{n_i} f) dn_i$ , we get  $\phi^*(g dn_i) = \sum_j \phi^*(P_{i,j} g) dl_j$ , a definition that proves to preserve closedness.

## Bibliography

- [1] Amdeberhan (Tewodros) and Zeilberger (Doron). – Hypergeometric series acceleration via the WZ method. *Electronic Journal of Combinatorics*, vol. 4, n° 2, 1997, p. Research Paper 3. 4 pages. – The Wilf Festschrift (Philadelphia, PA, 1996).
- [2] Cartan (Henri). – *Cours de calcul différentiel*. – Hermann, Paris, 1977, *Collection Méthodes*.
- [3] Chyzak (Frédéric). – An extension of Zeilberger's fast algorithm to general holonomic functions. *Discrete Mathematics*, vol. 217, n° 1-3, 2000, pp. 115–134. – Formal power series and algebraic combinatorics (Vienna, 1997).
- [4] Zeilberger (Doron). – Closed form (pun intended!). In *A tribute to Emil Grosswald: number theory and related analysis*, pp. 579–607. – American Mathematical Society, Providence, RI, 1993.
- [5] Zimmermann (Burkhard). – *Difference forms and hypergeometric summation*. – Master's thesis, Universität Wien, Vienna, Austria, February 2000.

# Transformations Exhibiting the Rank for Skew Laurent Polynomial Matrices

*Manuel Bronstein*

Projet CAFÉ, INRIA (France)

June 11, 2001

*Summary by Alin Bostan*

## Abstract

This talk presents an algorithm to perform transformations exhibiting the rank (TER) on a large class of matrices with entries in skew polynomial rings. This algorithm only uses elementary linear algebra operations and has various applications in solving very general linear functional systems.

## 1. Motivation

The question of finding polynomial solutions for linear functional systems is of particular interest in treating various problems in differential and difference algebra, as well as in combinatorics. It appears as a basic subtask in algorithms for finding all rational solutions of differential and ( $q$ -)difference equations, for computing liouvillian solutions of differential equations and ( $q$ -)hypergeometric solutions of ( $q$ -)difference equations. It also applies in factoring linear differential and difference operators, or in designing effective Gröbner basis algorithms in multivariate Ore algebras, which in turn are used in generalization of Gosper’s algorithm for indefinite hypergeometric summation and Zeilberger’s “creative telescoping” algorithm for definite summation and integration.

The traditional computer algebra approach to solving functional systems is via an elimination method like the cyclic-vector method, which converts the system to scalar equations (this procedure is called *uncoupling*). The major problem of this approach is the increase in size of the coefficients of equations.

The algorithm described in the next section offers a direct alternative for transforming a linear system of recurrences into an equivalent one of a *simpler form*, well-suited for the purpose of computing solutions with finite support of such a system. This gives a useful tool for constructing polynomial solutions of very general linear functional systems; see Sections 4.1 and 4.2 below.

The main advantage of this approach is that it does not require preliminary uncoupling of linear systems, but only performs elementary linear algebra operations on the original matrix.

## 2. Description of the Algorithm

The existence of canonical forms for matrices over various types of rings, such as principal ideal domains, has been known since the middle of the last century; their computation has important applications in both theoretical and practical areas of mathematics, science, and engineering.

Suppose that we consider matrices over a ring for which the notion of *rank* makes sense. A method for obtaining canonical forms of a matrix is performing *elementary operations* on its rows. Here, by elementary operation we mean permuting two rows, adding a multiple of a row to another row, and multiplying a row by a nonzero element of the base ring. Such a finite sequence of elementary row

operations on a matrix  $A$  can be represented by a matrix  $E$ . It will be called a *TER* (*transformation exhibiting the rank*) if it has the additional property that the rank of  $A$  equals the number of nonzero rows of the matrix  $EA$ .

In the commutative case, Gaussian elimination is the classical example of a TER, but it is a very greedy one, because of the exponential growth of the intermediate expressions; see [5]. The *Popov form* from linear control theory [8, 9] and the *reduced matrix form* [10, 11] are two other examples.

In [6] Mulders and Storjohann gave a simple algorithm that computes a simplified, non-canonical version of the Popov form, called the *weak Popov form* of a polynomial matrix. The algorithm performs only *delicate* elementary transformations which avoid intermediate expression swell. As a by-product, fast algorithms are obtained for computing the rank, the determinant, the Hermite form, the triangular factorization, and also the Popov form.

In the following, we describe an algorithm that computes a TER in a non-commutative setting.

Let  $R$  be an integral domain and  $\sigma$  an automorphism of  $R$ . Localizing the skew polynomial ring  $R[X; \sigma]$  at the set of powers of  $X$ , we obtain the *skew Laurent polynomial ring*

$$S = R[X, X^{-1}; \sigma],$$

with the commutation rules  $X \cdot r = \sigma(r) \cdot X$ , for all  $r$  in  $R$  (and therefore  $X^{-1} \cdot r = \sigma^{-1}(r) \cdot X^{-1}$ ). It is a left Ore domain, in the sense that any nonzero elements of  $S$  have a nonzero common left multiple in  $S$ . This implies that for any  $S$ -module  $M$ , the rank of  $M$ , denoted by  $\text{rk}(M)$  is a well-defined notion; see [4]. If  $A$  is a matrix with entries in  $S$ , we will call the *rank of  $A$*  the rank of the  $S$ -module generated by the rows of the matrix  $A$ .

We detail an algorithm which computes a TER of a  $n \times m$  matrix  $A$  with entries in the skew Laurent polynomial ring  $S = R[X, X^{-1}; \sigma]$ . If we write

$$A = A_t X^t + A_{t-1} X^{t-1} + \cdots + A_{s+1} X^{s+1} + A_s X^s,$$

where  $s \leq t$  are integers,  $A_i$  are matrices with entries in  $R$ , the leading matrix  $A_t$  and the trailing matrix  $A_s$  are nonzero, we are interested in finding a TER  $E$  such that the trailing matrix (respectively the leading matrix) of  $EA$  be nonsingular.

Remark that a straightforward application of the algorithm given in [6] does not do the job, even in the commutative case. The algorithm hereafter is essentially the algorithm proposed in [2] for the particular case of recurrence polynomials and improves the EG-elimination method [1].

The algorithm consists in iterating the following two basic steps, as long as the first operation can be performed:

1. look for a nonzero  $v \in R^n$  in the left kernel of the trailing (respectively leading) matrix of  $A$ , i.e., such that  $v^T A_s = 0$  [respectively  $v^T A_t = 0$ ] and such that  $v_i$  is zero whenever the  $i$ th row of  $A$  is zero;
2. choose  $i_0$  in the set of indices  $i$  such that the maximal degree in  $X$  of the polynomials of the  $i$ th row of  $A$  be maximal [respectively, its valuation be minimal] and replace this row by  $X^{-1} v^T A$  [respectively by  $X v^T A$ ].

Remark that  $\sum_i \deg({}_i A)$  decreases after each iteration, where  ${}_i A$  denotes the  $i$ th row of  $A$ , so the above algorithm terminates after at most  $n(t - s + 1)$  iterations.

Let  $N$  denote the number of iterations necessary for the previous algorithm to terminate and  $A^{(p)}$  the matrix obtained from  $A = A^{(0)}$  after  $p$  iterations. Then it can easily be seen that the number  $r$  of nonzero rows in the matrix  $A^{(N)}$  equals its rank, as any linear nontrivial dependency over  $S$  of these nonzero rows would imply a linear nontrivial dependency over  $R$  of the corresponding rows of its trailing matrix.



On the other hand, the ranks of the matrices  $A^{(p)}$  do not change all along the algorithm. This is implied by the formula  $\text{rank}A^{(p)} = \text{rank}A^{(p+1)} + \text{rank}(\mathcal{M}^{(p)}/\mathcal{M}^{(p+1)})$ , where  $\mathcal{M}^{(p)}$  denotes the  $S$ -module generated by the rows of the matrix  $A^{(p)}$ , and by the fact that  $\mathcal{M}^{(p)}/\mathcal{M}^{(p+1)}$  is a torsion module, therefore of rank zero.

This shows that the previous algorithm provides a TER for  $A$ .

### 3. Complexity

The previous algorithm only needs to compute nonzero elements of the kernels of matrices with entries in  $R$ . When  $R$  is a polynomial ring over some field  $K$  of characteristic 0, which is the case for differential and ( $q$ -)difference equations, one can use modular and probabilistic methods (like [7]) to find elements of the kernel. Their worst-case complexity is  $O(n^3 d^2)$  operations in  $K$ , where  $d$  is a bound on the degrees of the entries of  $A$ . Since the algorithm loops at most  $n(t-s+1)$  times, its complexity is  $O((t-s)n^4 d^2)$ . Refinements are possible; see [2].

### 4. Applications

**4.1. Desingularisation of recurrences.** As mentioned in the first section, linear systems of recurrences with variable coefficients are of interest in combinatorics and numeric computation. In addition, as shown in [3], they give a useful tool for constructing solutions of very general linear functional equations.

Consider the system  $A_t(n)Y_{n+t} + \dots + A_{s+1}(n)Y_{n+s+1} + A_s(n)Y_{n+s} = 0$ , where  $A_i$  are  $m \times m$  matrices with entries in the polynomial ring  $K[n]$ . This system is equivalent to  $AY = 0$ , where  $A = A_t E^t + \dots + A_s E^s$  is now viewed as a matrix with entries in  $K[n][E, E^{-1}; \sigma]$ ,  $\sigma$  being the shift automorphism of  $K[n]$ .

If either the leading matrix  $A_t$  or the trailing matrix  $A_s$  is nonsingular, its determinant is a nonzero polynomial in  $K[n]$  and the finite set of its integer roots gives the singularities of the recurrence and the possible degrees of polynomial solutions of the initial system. If the matrices  $A_s$  and  $A_t$  are singular, one faces the necessity to transform such a recurrence system into an equivalent one, with nonsingular leading (or trailing matrix). The following method is taken from [2]. If  $\text{rank}A = m > \text{rank}A_t$ , then applying the previous algorithm to the matrix  $A$  yields a new matrix

$$A^* = A_t^* E^t + \dots + A_{s'}^* E^{s'}$$

such that  $\text{rank}A_t^* = m$ .

**4.2. Solutions with finite support.** As already mentioned, the question of finding polynomial solutions of linear functional systems may be reduced to the problem of finding solutions with finite support  $(Y_0, Y_1, \dots, Y_N, 0, \dots)$  of the previous recurrence system; see [3]. In [2] a similar method to that of Section 4.1 was given, in order to find constraints on the set of the possible values of the bound  $N$  for the support of such a solution.

If  $\text{rank}A = m = \text{rank}A_s$  then we can find a finite set of candidates for  $N$ , given by the relation  $\delta(N-s) = 0$  for  $\delta(n) = \det A_s$ . If  $\text{rank}A = m > \text{rank}A_s$ , then applying the previous TER to the matrix  $A$  gives a matrix  $A^* = A_{t'}^* E^{t'} + \dots + A_s^* E^s$  where  $\text{rank}A_s^* = m$  and  $(\det A_s^*)(N-s) = 0$ .

**4.3. Hensel lifting for singular linear systems.** Let  $A$  be a nonsingular matrix with entries in  $K[X]$ , where  $K$  is a field. We consider the problem of recovering a  $v \in K(X)$  such that  $Av = b$ , or determine that no such  $v$  exists.

$X$ -adic lifting works by computing a vector series  $w = w_0 + w_1 X + w_2 X^2 + \dots$ , with each  $w_i \in K^n$  and such that

$$A(w_0 + w_1 X + w_2 X^2 + \dots) = b.$$

A rational solution  $v$  of the system  $Av = b$  is then reconstructed from the truncated series solution  $w \pmod{X^l}$  using Padé approximation. In general, we can compute the series solution  $w$ , by undetermined coefficients method, only when  $A$  is nonsingular modulo  $X$ .

In the case  $A(0)$  is singular, one can manage by applying the previous TER to the extended matrix  $[A \mid b]$  to transform the system  $AY = b$  into an equivalent one  $A^*Y = b^*$ , with  $A^*(0)$  nonsingular. A similar idea already appeared in [7].

**4.4. Solving linear differential systems.** We now consider the problem of solving a linear differential system  $Y' = B(x)Y$  where  $B$  is a  $m \times m$  matrix with entries in  $K[x]$ . By solving such a system we mean finding its formal power solutions. The system may be written in the compressed form  $AY = 0$ , where  $A$  is a matrix with entries in  $K[X][D; d/dx]$ .

Using the isomorphism of  $K$ -algebras:

$$\mathcal{R} : K[x, x^{-1}][D; d/dx] \longrightarrow K[n][E, E^{-1}; \sigma]$$

given by  $\mathcal{R}x = E^{-1}$  and  $\mathcal{R}D = (n+1)E$ , we remark that there is a bijective correspondence between formal power solutions  $Y = \sum_{n \geq 0} Y_n x^n$  of the linear differential system  $AY = 0$  and sequences  $Y = (Y_n)_{n \geq 0}$ , solutions of the recurrence system  $\mathcal{R}(A)(Y) = 0$ . This reduces the problem of finding (polynomial) solutions of the differential system  $AY = 0$  to finding solutions (with finite support) of the recurrence system  $\mathcal{R}(A)(Y) = 0$ .

### Bibliography

- [1] Abramov (Sergei A.). – EG-eliminations. *Journal of Difference Equations and Applications*, vol. 5, n° 4-5, 1999, pp. 393–433.
- [2] Abramov (Sergei A.) and Bronstein (Manuel). – On solutions of linear functional systems. In Mourrain (Bernard) (editor), *ISSAC'01 (July 22-25, 2001. London, Ontario, Canada)*. pp. 1–6. – ACM Press, 2001. Proceedings of the 2001 International Symposium on Symbolic and Algebraic Computation.
- [3] Abramov (Sergei A.), Bronstein (Manuel), and Petkovšek (Marko). – On polynomial solutions of linear operator equations. In Levelt (A. H. M.) (editor), *Symbolic and Algebraic Computation*. pp. 290–296. – ACM Press, New York, 1995. Proceedings of ISSAC'95, July 1995, Montreal, Canada.
- [4] Cohn (P. M.). – *Free rings and their relations*. – Academic Press Inc., London, 1985, second edition, *London Mathematical Society Monograph*, vol. 19, xxii+588p.
- [5] Fang (X. G.) and Havas (G.). – On the worst-case complexity of integer Gaussian elimination. In Küchlin (Wolfgang W.) (editor), *ISSAC'97 (July 21-23, 1997. Maui, Hawaii, USA)*. pp. 28–31. – ACM Press, New York, 1997. Conference proceedings.
- [6] Mulders (Thom) and Storjohann (Arne). – *On lattice reduction for polynomial matrices*. – Technical Report n° 356, ETH Zürich, Institute of Scientific Computing, December 2000. 26 pages.
- [7] Mulders (Thom) and Storjohann (Arne). – Rational solutions of singular linear systems. In Traverso (Carlo) (editor), *ISSAC'00 (August 6-9, 2000. St Andrews, Scotland)*. pp. 242–249. – ACM Press, 2000. Proceedings of the 2000 International Symposium on Symbolic and Algebraic Computation.
- [8] Popov (V. M.). – Invariant description of linear, time-invariant controllable systems. *SIAM J. Control*, vol. 10, 1972, pp. 252–264.
- [9] Villard (G.). – Computing Popov form and Hermite forms of polynomial matrices. In Lakshman (Y. N.) (editor), *ISSAC'96 (July 24-26, 1996. Zurich, Switzerland)*. pp. 250–258. – ACM Press, New York, 1996. Conference proceedings.
- [10] von zur Gathen (Joachim). – Hensel and Newton methods in valuation rings. *Mathematics of Computation*, vol. 42, n° 166, 1984, pp. 637–661.
- [11] von zur Gathen (Joachim) and Gerhard (Jürgen). – *Modern computer algebra*. – Cambridge University Press, New York, 1999, xiv+753p.

## A Criterion for Non-Complete Integrability of Hamiltonian Systems

*Delphine Boucher*

Université de Limoges (France)

January 15, 2001

*Summary by Philippe Dumas and Bruno Salvy*

### Abstract

Finding polynomial solutions of linear differential equations is a building block implemented in several algorithms of computer algebra systems. In particular, this is a necessary sub-step when looking for rational, algebraic or Liouvillian solutions of linear differential equations. When there are no parameters, several algorithms are available, but the general case with parameters is undecidable. However, special families can be handled by *ad hoc* methods. Such methods were developed by Boucher who applied them to the nice example of integrability of the 3-body problem. The key idea there is to rely on a recent result of Morales-Ruiz and Ramis who relate complete integrability and differential Galois group. It turns out that special properties of this group can be related to computable properties of an appropriate linear differential equation, which leads Boucher to a “simple” sufficient condition for non-complete integrability.

### 1. Polynomial Solutions of Linear Differential Equations

The classical method to find polynomial solutions of linear differential equations over  $\mathbb{K}(x)$ , where  $\mathbb{K}$  is a field, starts by determining a bound on the degree of potential solutions. This is a bound on the integer solutions of the *indicial equation* at infinity.

Once a bound on the degree has been found, one uses an indeterminate coefficients method. The linear system on these coefficients has a band-matrix structure which can be exploited to accelerate the computation [1]. This linear system is rectangular, with more equations than unknown coefficients, thus existence of solution is related to the vanishing of a determinant.

When parameters occur in the equation ( $\mathbb{K}$  is a field of rational functions), there are two difficulties: the size of the matrix may depend on the parameters and even when it does not, the determinant which must vanish is a polynomial in the parameters. Using Matijasevich’s result on the undecidability of Hilbert’s 10th problem (*Is there a finite process which determines if a polynomial equation is solvable in integers?*), it is possible to show that this problem itself is undecidable. More precisely, Jacques-Arthur Weil observes that the equation

$$y'(x) - \left( \frac{a_1}{x-1} + \dots + \frac{a_m}{x-m} + P(a_1, \dots, a_m) \right) y(x) = 0$$

has rational solutions if and only if  $P(a_1, \dots, a_m) = 0$  has integral solutions.

There are still cases where all polynomial solutions can be found: this happens when either the size of the matrix is bounded independently of the parameters and the vanishing of the required determinant can be determined or when the structure of the matrix is sufficiently regular to make the decision possible. Examples of both cases are given in [4].

## 2. Complete Integrability

**2.1. Hamiltonian Mechanics.** In the Hamiltonian approach to classical mechanics, the state of a system is characterized by  $2n$  variables,  $q_i$  (positions) and  $p_i$  (momenta),  $i = 1, \dots, n$ , living in an open subset  $U$  of  $\mathbb{R}^{2n}$  (the phase space). More generally, the phase space of a system is the cotangent fibre bundle  $T^*M$  of an  $n$ -dimensional real manifold  $M$ . The formulae we give below are expressions in a chart of useful quantities. The state variables satisfy

$$(1) \quad \dot{p}_i = \frac{\partial H}{\partial q_i}, \quad \dot{q}_i = -\frac{\partial H}{\partial p_i},$$

where a dot denotes a derivative with respect to time and  $H(p, q, t)$  is the *Hamiltonian*. Physically, the Hamiltonian often represents the energy of the system. The system (1) governs the evolution of the system (in the phase space  $U$ ). Solutions  $\gamma(t)$  of (1) are the trajectories of the system.

In a more abstract setting,  $\mathbb{R}^{2n}$  is endowed with a non-degenerate 2-form

$$\omega = \sum_{i=1}^n dp_i \wedge dq_i,$$

known as Liouville's symplectic 2-form. Since  $\omega$  is non-degenerate, it induces an isomorphism between  $\mathbb{R}^{2n}$  and its dual under which  $-dH$  is the image of a vector field  $X_H$ . In this language, the Hamiltonian system (1) reduces to

$$\dot{\gamma} = X_H(\gamma).$$

First integrals are functions  $F(p, q)$  that are constant along the solutions  $\gamma(t)$ . A necessary and sufficient condition is

$$\{F, H\} := \sum_i \frac{\partial F}{\partial p_i} \frac{\partial H}{\partial q_i} - \frac{\partial H}{\partial p_i} \frac{\partial F}{\partial q_i} = 0,$$

where  $\{F, H\}$  is known as the *Poisson bracket* of  $F$  and  $H$ . In particular, the Hamiltonian itself is a first integral.

Two first integrals are *in involution* if their Poisson bracket vanishes. A Hamiltonian system is *completely integrable* when it possesses a set of  $n$  first integrals in involution that are independent (i.e., their Jacobian matrix is regular in the open set  $U$ ).

Informally, a completely integrable system can be “solved” in terms of its first integrals. Indeed, given a first integral, a process known as *symplectic reduction* makes it possible to reduce the number of degrees of freedom by 1, i.e., the dimension by 2 [2, p. 91].

**2.2. Many-Body Problem.** In the many-body problem,  $n$  particles obeying Newton's law are governed by the following Hamiltonian:

$$H(p, q) = \frac{1}{2} \sum_i \frac{\|p_i\|^2}{m_i} - \sum_{i \neq j} \frac{m_i m_j}{\|q_j - q_i\|}.$$

Note that here each  $p_i$  and  $q_i$  has coordinates in  $\mathbb{R}^3$ , thus the phase space has dimension  $6n$ .

Apart from the Hamiltonian itself, known first integrals for this system are the momentum of the centre of mass and the angular momentum  $\sum q_i \wedge p_i$ . Thus, the number of degrees of freedom can be reduced from  $3n$  to  $3n - 6$  (or from  $2n$  to  $2n - 4$  in the planar case).

For the 3-body problem, Poincaré proved that there are no other *complex analytic* first integrals. Bruns proved a similar result for *complex algebraic* first integrals.

**2.3. Theorem of Morales-Ruiz and Ramis.** We now present a simple version of a result of Morales-Ruiz and Ramis in [10, 11, 12, 13] (see also [3]) on non-complete integrability in terms of *meromorphic* first integrals. The Hamiltonian is analytic over an open set of  $\mathbb{C}^{2n}$  and  $t$  (time) is a *complex* variable. Given a non-stationary trajectory  $\Gamma(t)$ , following an idea of Poincaré, one considers the *linear* differential equation that must satisfy a “small” variation  $\eta$ , such that  $\Gamma(t) + \eta(t)$  is solution of the Hamiltonian system. This equation

$$(2) \quad \dot{\eta} = X'_H(\gamma) \cdot \eta$$

is called the *variational equation* along  $\Gamma$ . A theorem of Morales-Ruiz and Ramis relates complete integrability and Galois group of this equation. (For an introduction to differential Galois theory, see [14] or the summary of Ulmer’s talk in this seminar in 1994.) However, since the Galois group is often very difficult to compute, it is useful to consider a differential equation of lower order. This is achieved by the following result.

**Theorem 1** (Morales-Ruiz and Ramis). *If the system possesses  $n$  meromorphic first integrals in the neighbourhood of  $\Gamma$ , independent and in involution, then the connected component of identity in the differential Galois group of the normal variational equation along  $\Gamma$  is abelian.*

Similar earlier results of Ziglin based on the monodromy group and of Churchill, Singer *et alii* based on the Galois group did not extend to the case where the variational equation has an irregular singular point. In this theorem, the *normal* variational equation is an equation obtained from the variational equation through symplectic reduction. Indeed,  $dH(\Gamma(t)) \cdot \eta$  is a first integral of the variational equation, as can be seen by a first-order expansion.

### 3. Boucher’s Criterion and its Application

It is not necessary to compute the Galois group of a linear differential equation in order to detect that it is not abelian. Thanks to a sufficient criterion [5, 6], Boucher has proved that the planar 3-body problem is not completely integrable in terms of meromorphic first integrals. Unfortunately, the formulae involved in this derivation are much too large to be reproduced here. Thus we content ourselves with a sketch of the steps and a description of the tools used in the calculations.

#### 3.1. Criterion.

**Theorem 2.** *Assume that the linear differential operator  $L$  can be factored as  $KM$ , with  $M = \text{lcm}(L_1, \dots, L_m)$  where the  $L_i$ ,  $i = 1, \dots, m$ , are irreducible (and  $\text{lcm}$  denotes the least common left multiple). Assume moreover that  $M(y) = 0$  has a formal solution with a logarithm. Then the connected component of the differential Galois group of  $L(y) = 0$  is not abelian.*

Given a linear differential equation, this theorem reduces the task to factoring and finding formal solutions. Factoring can be done by an algorithm of van Hoeij [19, 20], and formal solutions can be computed at any singularity, including infinity [15, 20].

**3.2. Application to the 3-Body Problem.** Tsygintsev and Boucher have proved independently that the planar 3-body problem is not completely integrable in terms of meromorphic first integrals. Their approaches [5, 17] follow the same initial steps till the normal variational equation. Then [17] uses Ziglin’s result. We now outline Boucher’s approach.

*Reduced Hamiltonian.* Using the first integrals obtained in Section 2.2, the problem is reduced to a Hamiltonian with three degrees of freedom, given in [17]. The parameters in this equation are the three masses  $m_1, m_2, m_3$  and the value  $c$  of the angular momentum (which reduces to a scalar in this dimension). By homogeneity, we can freely assume  $m_3 = 1$ . (Note that these transformations make the resulting expressions asymmetric with respect to the bodies.)

In order to apply Theorem 1, we need a particular solution of the system. This is provided by the celebrated *Lagrange solutions*. In these solutions, the three particles have orbits on similar conics with a common focus located at their centre of mass (see [8, p. 400]). Since any particular solution can be chosen, Tsygvintsev and Boucher concentrate on the parabolic orbit (for angular momentum  $c \neq 0$ ).

*Variational Equation.* The variational equation (2) is a linear system of order  $n = 6$ . The normal variational equation is obtained via a linear change of symplectic basis as follows. We observe that  $X_H$  itself is a solution of the variational equation. It will be the first vector  $e_1$  of the new basis. Next, we compute a basis  $(e_1 = X, e_2, \dots, e_n, e_{n+2}, \dots, e_{2n})$  of the kernel of  $dH(\Gamma(t))$  satisfying  $\omega(e_i, e_{n+i}) = 1$  for  $1 < i \leq n$  and  $\omega(e_i, e_j) = 0$  otherwise. Finally, we compute a vector  $e_{n+1} = Y$  such that  $\omega(e_i, Y) = 0$  for  $i \neq 1$  and  $\omega(X, Y) = 1$ . In the new basis  $(e_1, \dots, e_{2n})$ , the first column of the matrix of the variational equation is 0, since  $X_H$  is a solution. Now, for any vector field  $\eta$ ,  $\omega(X, \eta) = -dH(\Gamma(t)) \cdot \eta$ , therefore for any solution  $\eta$  of the variational equation, the value of this first integral is the coordinate of  $\eta$  on the vector  $Y$  in the new basis. The normal variational equation is obtained by setting this coordinate to 0 and considering the induced matrix  $A$  on the subspace with basis  $(e_2, \dots, e_n, e_{n+2}, \dots, e_{2n})$ .

*Cyclic Vector.* The criterion of Theorem 1 applies to equations rather than systems. A classical method to convert a system of order  $m$  into an equation  $L(u) = 0$  is to start from a random vector  $u$  and find a linear dependency between the  $m+1$  vectors  $u, u', \dots, u^{(m)}$  where the derivatives are computed using the matrix  $A$ . Unfortunately, this process generically introduces spurious singularities that are roots of the determinant of the change of basis  $(u, u', \dots, u^{(m-1)})$ . Boucher therefore selects cyclic vectors in such a way that no new singularity occurs and this requires distinguishing two cases depending on the value of the mass  $m_1$ .

*Right Factors.* In the simplest case of Boucher's criterion, the operator  $L$  has an irreducible right factor  $M$  whose formal solutions exhibit logarithms. This requires  $M$  to have order at least 2. Factors of order  $k$  are found by constructing an auxiliary equation  $L^{\wedge k}$  of order  $\binom{m}{k}$  whose solutions are Wronskians of  $k$  independent solutions of  $L$  [7]. (Note that this can be computed directly from  $L$ .) Indeed, a monic right factor of order  $k$  has for coefficient of order  $k-1$  the logarithmic derivative  $w'/w$  of some particular Wronskian of its solutions. Finding right factors then amounts to looking for so-called *exponential solutions* of  $L^{\wedge k}$  (i.e., those with logarithmic derivative that is rational). From a basis of such solutions, corresponding to linear combinations of Wronskians, Plücker's relations help select those that are indeed Wronskians [16]. From there, the complete factor can be reconstructed. Exponential solutions are found by looking at formal solutions at all singularities of the equation [19]. This requires a discussion in the parametric case. If a factor is found, the next step is to check whether this factor is irreducible, or to find conditions on the parameters that make it irreducible. This is done again by searching for factors of the factor. It turns out that in this application, in all cases an irreducible right factor of order 2 is found.

*Logarithms.* Logarithms in formal solutions occur when the indicial equation at a singularity has roots that differ by an integer. A necessary and sufficient condition has been given by Frobenius [9, p. 404–406]. Again, in all generality nothing can be said when parameters are present but Boucher manages to show that logarithms are present in all cases for this application.

#### 4. Conclusion

This application is a very good showcase for many of the algorithms that have been developed in computer algebra for linear differential equations: formal solutions, factorization, polynomial solutions, . . .

What Boucher has shown is that, even in the presence of parameters, these algorithms can be exploited to provide useful information by concentrating on those points where specific quantities such as the indicial equation or its solutions do not depend “too much” on the parameters.

A recent trend in computer algebra is to revisit all these algorithms that have been designed for equations and extend them to deal with systems, without using the cyclic vector. It would be a natural step to try and adapt Boucher’s criterion so that the symplectic structure is not lost. (Work on this has been started by Boucher and Weil.)

*Remark.* A new result of Tsygvintsev [18] shows the stronger result that there is no additional meromorphic first integral. Also, Theorem 2 has been extended to the case when  $L$  is a product of irreducible factors one of which has a solution with logarithms.

#### Bibliography

- [1] Abramov (Sergei A.), Bronstein (Manuel), and Petkovšek (Marko). – On polynomial solutions of linear operator equations. In Levelt (A. H. M.) (editor), *Symbolic and Algebraic Computation*. pp. 290–296. – ACM Press, New York, 1995. Proceedings of ISSAC’95, July 1995, Montreal, Canada.
- [2] Arnol’d (V. I.), Kozlov (V. V.), and Neishtadt (A. I.). – *Dynamical systems. III* – Springer-Verlag, Berlin, 1988, xiv+291p. Mathematical Aspects of Classical and Celestial Mechanics. Translated from the Russian.
- [3] Audin (Michèle). – Intégrabilité et non-intégrabilité de systèmes hamiltoniens [d’après S. Ziglin, J. Morales-Ruiz, J.-P. Ramis, . . .]. – To appear in *Astérisque*. Bourbaki seminar 884, March 2001. Preliminary version available at <http://irma.u-strasbg.fr/~maudin/publications.html>.
- [4] Boucher (Delphine). – About the polynomial solutions of homogeneous linear differential equations depending on parameters. In Dooley (Sam) (editor), *Proceedings of the 1999 International Symposium on Symbolic and Algebraic Computation (Vancouver, BC)*. pp. 261–268. – ACM, New York, 1999.
- [5] Boucher (Delphine). – Sur la non-intégrabilité du problème plan des trois corps de masses égales. *Comptes Rendus de l’Académie des Sciences. Série I. Mathématique*, vol. 331, n° 5, 2000, pp. 391–394.
- [6] Boucher (Delphine). – *Sur les équations différentielles linéaires paramétrées, une application aux systèmes hamiltoniens*. – PhD Thesis, Université de Limoges, October 2000.
- [7] Bronshteĭn (M.) and Petkovshek (M.). – Ore rings, linear operators and factorization. *Rossiĭskaya Akademiya Nauk. Programirovanie*, vol. 1994, n° 1, 1994, pp. 27–44.
- [8] Hestenes (David). – *New foundations for classical mechanics*. – D. Reidel Publishing Co., Dordrecht, 1986, xii+644p.
- [9] Ince (E. L.). – *Ordinary differential equation*. – Dover Publications, New York, 1956, viii+558p. Reprint of the 1926 edition.
- [10] Morales-Ruiz (Juan J.) and Ramis (Jean Pierre). – Galoisian obstructions to integrability of Hamiltonian systems: statements and examples. In *Hamiltonian systems with three or more degrees of freedom (S’Agaró, 1995)*, pp. 509–513. – Kluwer Acad. Publ., Dordrecht, 1999.
- [11] Morales-Ruiz (Juan J.) and Ramis (Jean-Pierre). – Galoisian obstructions to integrability of Hamiltonian systems. *Methods and Applications of Analysis*, vol. 8, n° 1, 2001.
- [12] Morales-Ruiz (Juan J.) and Ramis (Jean-Pierre). – Galoisian obstructions to integrability of Hamiltonian systems. II. *Methods and Applications of Analysis*, vol. 8, n° 1, 2001.
- [13] Morales-Ruiz (Juan J.) and Ramis (Jean-Pierre). – A note on the non-integrability of some Hamiltonian systems with a homogeneous potential. *Methods and Applications of Analysis*, vol. 8, n° 1, 2001.

- [14] Singer (Michael F.). – An outline of differential Galois theory. In *Computer algebra and differential equations*, pp. 3–57. – Academic Press, London, 1990.
- [15] Tournier (Évelyne). – *Solutions formelles d'équations différentielles*. – Doctorat d'État, Université scientifique, technologique et médicale de Grenoble, 1987.
- [16] Tsarëv (S. P.). – Some problems that arise in the factorization of linear ordinary differential operators. *Rossiiskaya Akademiya Nauk. Programirovanie*, n° 1, 1994, pp. 45–48. – (Russian). Translation in *Programming and Computer Software*, vol. 20, n° 1, 1994, pp. 27–29.
- [17] Tsygvintsev (Alexei). – La non-intégrabilité méromorphe du problème plan des trois corps. *Comptes Rendus de l'Académie des Sciences. Série I. Mathématique*, vol. 331, n° 3, 2000, pp. 241–244.
- [18] Tsygvintsev (Alexei). – Sur l'absence d'une intégrale première méromorphe supplémentaire dans le problème plan des trois corps. *Comptes Rendus de l'Académie des Sciences. Série I. Mathématique*, vol. 333, n° 2, 2001, pp. 125–128.
- [19] van Hoeij (Mark). – Factorization of differential operators with rational functions coefficients. *Journal of Symbolic Computation*, vol. 24, n° 5, 1997, pp. 537–561.
- [20] van Hoeij (Mark). – Formal solutions and factorization of differential operators with power series coefficients. *Journal of Symbolic Computation*, vol. 24, n° 1, 1997, pp. 1–30.



## Effective Algebraic Analysis in Linear Control Theory

Alban Quadrat

University of Leeds (United Kingdom) & Projet CAFÉ, INRIA Sophia-Antipolis (France)

December 4, 2000

Summary by Frédéric Chyzak

### Abstract

In the 1960's, Malgrange made use of D-module theory for studying linear systems of PDEs [2]. Several aspects of this approach, now called *algebraic analysis*, have then been made effective in the 1990's, owing to the extension of the theory of Gröbner bases to rings of differential operators. Correspondingly, algorithms have also been implemented in several systems. Recently, the introduction of algebraic analysis to control theory has allowed to classify linear multidimensional control systems according to algebraic properties of associated D-modules, to redefine their structural properties in a more intrinsic fashion, and to develop effective tests for deciding these structural properties [3, 6, 7, 8, 9, 10, 12, 14].

### 1. From Linear Multidimensional Control Systems to Algebraic Analysis

A control system relates the state  $x$  of a physical process with an external command  $u$  and some output  $y$ . Each of  $u$ ,  $x$ , and  $y$  is a vector of functions of the time  $t$ , and the system describes their evolution with  $t$ . Several classes of such systems can be represented by matrices with coefficients in a ring of operators. Sample classes are the following:

1. *Kalman systems* are first-order linear (ordinary) differential systems

$$\dot{x} = Ax + Bu, \quad y = Cx + Du,$$

where  $A$ ,  $B$ ,  $C$ , and  $D$  are matrices with real entries [5]. For example, RLC circuits can be described by Kalman systems.

2. *Polynomial systems* are higher-order differential systems expressed without the help of any state variable, in the form

$$(1) \quad P(d/dt)y(t) + Q(d/dt)u(t) = 0.$$

Here  $P$  and  $Q$  are matrices with coefficients that are scalar linear differential operators with real coefficients [5]. For example, a harmonic oscillator commanded by a force is described by a second-order polynomial system. By Laplace transform, an equivalent formulation of (1) is

$$P(s)\hat{y}(s) + Q(s)\hat{u}(s) = 0;$$

the matrices  $P$  and  $Q$  are now matrices of polynomials in  $s$  with real coefficients [5].

3. *Differential-delay systems with constant delays* are a generalization common to Kalman systems and polynomial systems by introducing the constant-delay operators  $\delta_i$  defined by

$(\delta_i f)(t) = f(t - t_i)$  for some real  $t_i$ . The generalized forms are

$$\dot{x}(t) = \sum_{i=0}^r A_i x(t - t_i) + B_i u(t - t_i), \quad y(t) = \sum_{i=0}^r C_i x(t - t_i) + D_i u(t - t_i),$$

and

$$P(d/dt, \delta_1, \dots, \delta_r)y + Q(d/dt, \delta_1, \dots, \delta_r)u = 0,$$

respectively. A typical occurrence of delay is when transmitting a signal  $u$  through a channel.

4. *Multivariate linear differential systems with real coefficients* appear frequently to describe physical phenomena, like electromagnetism, (linear) elasticity, hydrodynamism, and so on [7, 8, 12].

In each case, the column vector  $\xi = (y, x, u)^T$  satisfies  $R\xi = 0$  for a (rectangular) matrix  $R$  with coefficients in some ring  $\mathbb{A}$ . Thus, we henceforth consider a linear control system as defined by a matrix  $R$  with coefficients in an entire ring  $\mathbb{A}$ . To give simple examples, the matrix forms corresponding to Kalman and polynomial systems respectively are

$$R = \begin{pmatrix} 0 & A - d/dt \text{ Id} & B \\ \text{Id} & C & D \end{pmatrix} \quad \text{and} \quad R = \begin{pmatrix} P & Q \end{pmatrix}.$$

In these differential cases, the ring  $\mathbb{A}$  is  $\mathbb{R}[d/dt]$  or a multivariate generalization, but more general rings of coefficients are also considered in place of  $\mathbb{R}$  in applications, like the ring  $\mathbb{R}(t)$  of rational function, or the ring  $C^\infty(I)$  of infinitely differentiable functions over some real interval  $I$ . In the equivalent formulation by Laplace transform or in the mixed differential-delay situation with *constant* coefficients, the ring is isomorphic to the polynomial ring  $\mathbb{R}[s]$  or a multivariate analogue. Here again, more general rings of functions often appear in applications, like:  $\mathbb{R}[s, \exp(-s)]$ , for situations related to the wave equation; or the ring  $H_\infty(\mathbb{C}_+)$  of complex-analytic functions bounded in the right half complex plane  $\mathbb{C}_+$  (Hardy space) and its subring  $\text{RH}_\infty(\mathbb{C}_+)$  of real rational functions with no pole on the right half complex plane, for the study of the stability of some distributed systems [11].

Several structural properties of systems are all-important in control theory. An *observable* of a control system is any scalar function of its command  $u$ , state  $x$ , and output  $y$  and of their derivatives up to a certain order. An observable is called *autonomous* if it satisfies a non-trivial PDE. A control system is called *controllable* if no observable is autonomous. The study of structural properties of a system turns out to lead to linear algebra: controllability and observability are related to various notions of primeness of the linear maps

$$z \mapsto Rz \quad \text{and} \quad z \mapsto zR;$$

in the polynomial systems case, stability is related to poles and zeroes of the system, that are invariant factors of the matrix  $R$ ; similarly with the existence of generalized Bézout identities and flatness of a control system; etc.

By associating an  $A$ -module  $M$  to the matrix  $R$ , another interpretation of the structural properties is in terms of module-theoretic and homological properties of  $M$  (torsion, torsion-free, reflexive, and projective modules; extension and torsion functors). In fact, a full classification of modules by homological algebra methods translates into a classification of linear control systems.

## 2. Duality Between Differential Operators and D-Modules

Let us turn to the formal theory of PDEs [13]. Starting with a naive viewpoint on differential operators (so as to avoid the formalism of jet bundles), we introduce *formally exact sequences of differential operators*. For each  $k$ , let  $F^k$  denote the algebra of functions in  $k$  variables, and consider

a differential operator  $\mathcal{D}$  from  $F^m$  to  $F^l$  (of finite order). Given  $\eta \in F^n$ , the necessary conditions for the existence of  $\xi \in F^m$  such that  $\mathcal{D}\xi = \eta$  are called *compatibility conditions* of  $\mathcal{D}$ ; they take the form  $\mathcal{D}_1\eta = 0$  for some differential operator  $\mathcal{D}_1$ . Writing  $\mathcal{D}_0 = \mathcal{D}$ , we have  $\mathcal{D}_1 \circ \mathcal{D}_0 = 0$ . When  $\mathcal{D}_1$  encapsulates all compatibility conditions, the sequence

$$F^m \xrightarrow{\mathcal{D}_0} F^{l_0} \xrightarrow{\mathcal{D}_1} F^{l_1}$$

of differential operators is called *formally exact* (at  $F_{l_0}$ ). Formally exact sequences can always be extended (to the right) into longer sequences, so that denoting the solution set of  $\mathcal{D} = \mathcal{D}_0$  in  $F^m$  by  $\Theta$ , we obtain a formally exact sequence

$$0 \rightarrow \Theta \rightarrow F^m \xrightarrow{\mathcal{D}_0} F^{l_0} \xrightarrow{\mathcal{D}_1} F^{l_1} \xrightarrow{\mathcal{D}_2} F^{l_2} \rightarrow \dots$$

(at  $\Theta$  and each  $F^{l_k}$ ) where the first two maps denote inclusions. Under technical conditions (*regularity* and *involutivity*), the formal theory of PDEs proves the existence of a *finite* formally exact sequence for  $\mathcal{D}$ , in the sense that  $F^{l_n} = 0$  from some  $n$  on, by exhibiting a canonical, formally exact sequence

$$(2) \quad 0 \rightarrow \Theta = \ker \mathcal{D}_0 \rightarrow F^m \xrightarrow{\mathcal{D}_0} F^{l_0} \xrightarrow{\mathcal{D}_1} F^{l_1} \xrightarrow{\mathcal{D}_2} F^{l_2} \rightarrow \dots \xrightarrow{\mathcal{D}_r} F^{l_r} \rightarrow 0$$

called the *Janet sequence* of  $\mathcal{D}$ , in which each (non-zero)  $\mathcal{D}_i$  is of order 1 (and involutive) for  $i \geq 1$ , and  $r$  is the number of derivatives.

A dual, more algebraic counterpart to this differential viewpoint is in terms of *exact sequences of  $D$ -modules*. To this end, we now view each  $\mathcal{D}_i$  as defined by an  $l_i \times l_{i-1}$  matrix  $R_i$  of multivariate linear differential operators in

$$\mathbb{A} = \mathbb{R}(x_1, \dots, x_r)[\partial_1, \dots, \partial_r].$$

(We set  $l_{-1} = m$ .) In terms of matrices,

$$\mathcal{D}_i = R_i \cdot = (\xi \mapsto R_i \xi),$$

so that  $R_{i+1}R_i \cdot = 0$ . We then consider the maps  $\cdot R_i$  from  $\mathbb{A}^{l_i}$  to  $\mathbb{A}^{l_{i-1}}$ , whose elements are viewed as row vectors. To start with, the map  $\cdot R_0$  defines an algebraic representation of a generic solution  $\xi$  the PDE system  $\mathcal{D}_0\xi = 0$  in the following way. Let  $(e_1, \dots, e_m)$  be the canonical basis of  $\mathbb{A}^m$  and consider the maps

$$(3) \quad 0 \leftarrow M = \mathbb{A}^m / \mathbb{A}^{l_0} R_0 \xleftarrow{\pi} \mathbb{A}^m \xleftarrow{\cdot R_0} \mathbb{A}^{l_0},$$

where  $\pi$  denotes the canonical projection  $\pi(v) = v + \mathbb{A}^{l_0} R_0$ . The *cokernel*

$$M = \text{coker}(\cdot R_0) = \mathbb{A}^m / \mathbb{A}^{l_0} R_0$$

of  $\cdot R_0$  contains the announced generic solution: setting

$$\xi_i = \pi(e_i) = e_i + \mathbb{A}^{l_0} R_0,$$

we get  $\mathcal{D}_0\xi = \xi R_0 = 0$ . Other members of  $M$  correspond to linear combinations of the  $\xi_i$  and their derivatives, i.e., to the observables defined above. We now proceed to follow up with the next  $\mathcal{D}_i$ 's. A sequence

$$L \xrightarrow{u} L' \xrightarrow{v} L''$$

of linear maps (between modules) is said to be *exact* (at  $L'$ ) if  $\text{im } u = \ker v$ . (Thus (3) is exact at  $M$  and  $\mathbb{A}_{l_0}$ .) It can be shown that any Janet sequence (2) gives rise to the exact sequence

$$(4) \quad 0 \leftarrow M \xleftarrow{\pi} \mathbb{A}^m \xleftarrow{\cdot R_0} \mathbb{A}^{l_0} \xleftarrow{\cdot R_1} \mathbb{A}^{l_1} \xleftarrow{\cdot R_2} \mathbb{A}^{l_2} \leftarrow \dots \xleftarrow{\cdot R_r} \mathbb{A}^{l_r} \leftarrow 0$$

(at  $M$  and each  $\mathbb{A}^{l_k}$ ). Here,  $\cdot R_{i+1}R_i = 0$  by exactness. Since  $\mathbb{A}$  has no zero divisor, this means that  $R_{i+1}R_i = 0$ . The sequence (4) of (left)  $\mathcal{D}$ -modules is called a *free resolution* of  $M$ : it encapsulates the obstruction of  $M$  to be free (as the module  $\ker \pi = \text{im}(\cdot R_0)$ ), then the obstruction of  $\ker \pi$  to be free (as the module  $\ker(\cdot R_0) = \text{im}(\cdot R_1)$ ), etc. (A module is called *free* when it is isomorphic to some  $\mathbb{A}^r$ , whence the name “free resolution.”)

### 3. Parametrization and Controllability

A problem dual to the search of compatibility conditions is, for a given differential equation  $\mathcal{D}\xi = 0$ , to determine whether the solutions can be parametrized by certain arbitrary functions which, in physical systems, play the role of *potentials*. In other words, the problem is to determine whether there exists another operator

$$\mathcal{D}_{-1} : F^{l-1} \rightarrow F^{l_0}$$

whose compatibility conditions are described by  $\mathcal{D} = \mathcal{D}_0$ , i.e., to look for a formally exact sequence

$$F^{l-1} \xrightarrow{\mathcal{D}_{-1}} F^{l_0} \xrightarrow{\mathcal{D}_0} F^{l_1}.$$

In this situation, for any  $\xi \in F^{l_0}$  the existence of  $\pi \in F^{l-1}$  satisfying  $\mathcal{D}_{-1}\pi = \xi$  is equivalent to the fact that  $\xi$  solves the differential equation  $\mathcal{D}_0\xi = 0$ , and so  $\mathcal{D}_{-1}$  “parametrizes”—in the usual sense—all its solutions.

The existence of a parametrization has a nice application to *optimal command*: assume one needs to minimize a cost function provided by the integral  $\int_0^T F(t) dt$  of an observable  $F$  of some system  $\mathcal{D}_0$ . The optimization problem is then to minimize over all tuples  $\xi = (y, x, u)^T$  of functions constrained by  $\mathcal{D}_0\xi = 0$ . On the other hand, once the solutions  $\xi$  are given by a parametrization  $\xi = \mathcal{D}_{-1}\pi$ , the optimization problem reduces to the non-constrained problem of minimizing the integral  $\int_0^T G(t) dt$  of a new observable  $G$  of  $\mathcal{D}_{-1}$  over unconstrained  $\pi$  [12].

To study the control-theoretic properties of the differential operator  $\mathcal{D}$ , starting with the existence of a parametrization, we in fact study the module-theoretic properties of  $M$ , which in turn are derived from the study of the right  $\mathcal{D}$ -module defined by

$$(5) \quad \mathbb{A}^{l-1} \xrightarrow{R_0} \mathbb{A}^{l_0} \rightarrow N = \text{coker}(R_0 \cdot) = \mathbb{A}^{l_0} / R_0 \mathbb{A}^{l-1} \rightarrow 0$$

(recall that  $l_{-1} = m$  and compare with (3)). The key ingredient to be used comes from linear algebra: dualization, which maps a left  $\mathbb{A}$ -module  $L$  to the right module  $\text{hom}_{\mathbb{A}}(L, \mathbb{A})$  of  $\mathbb{A}$ -linear applications from  $L$  to  $\mathbb{A}$ . Correspondingly, any linear map  $L \xrightarrow{u} L'$  induces a map from the dual of  $L'$  to the dual of  $L$ : to  $\lambda \in \text{hom}_{\mathbb{A}}(L', \mathbb{A})$ , one associates  $\lambda \circ u \in \text{hom}_{\mathbb{A}}(L, \mathbb{A})$ . This takes a simple form when the modules are free and of finite rank (i.e.,  $L = \mathbb{A}^m$  and  $L' = \mathbb{A}^l$ , viewed as left modules of row vectors). Indeed, the linear map  $u$  is just the application of an  $m \times l$  matrix  $U$ :  $u = \cdot U$ . Elements  $\mu \in \text{hom}_{\mathbb{A}}(\mathbb{A}^k, \mathbb{A})$  are defined by their values on the canonical basis  $(e_i)$  of  $\mathbb{A}^k$  by

$$\mu = \cdot (\mu(e_1), \dots, \mu(e_k))^T,$$

so that the dual of  $\mathbb{A}^k$  is isomorphic to  $\mathbb{A}^k$  (now viewed as a right module of column vectors). In this setting, the dual of a map  $\mathbb{A}^m \xrightarrow{\cdot U} \mathbb{A}^l$  is  $\mathbb{A}^m \xleftarrow{U \cdot} \mathbb{A}^l$ . The same ideas apply mutatis mutandis for the dual of right modules.

To search for a parametrization, one thus extends the exact sequence (5) into an exact sequence

$$\mathbb{A}^{l-2} \xrightarrow{R_{-1} \cdot} \mathbb{A}^{l-1} \xrightarrow{R_0 \cdot} \mathbb{A}^{l_0} \rightarrow N \rightarrow 0.$$

An algorithm for this purpose will be given in Section 5. By dualization (i.e., application of the  $\text{hom}_{\mathbb{A}}(\cdot, \mathbb{A})$  functor), it becomes a sequence

$$\mathbb{A}^{l-2} \xleftarrow{\cdot R_{-1}} \mathbb{A}^{l-1} \xleftarrow{\cdot R_0} \mathbb{A}^{l_0} \leftarrow \text{hom}_{\mathbb{A}}(N, \mathbb{A}) \leftarrow 0$$

of left D-modules that is *usually no longer exact*. In particular, we may well have  $\ker(\cdot R_{-1})$  strictly larger than  $\text{im}(\cdot R_0)$ . Upon forgetting the map  $\cdot R_0$  and prolonging  $\cdot R_{-1}$  into

$$\mathbb{A}^{l-2} \xleftarrow{\cdot R_{-1}} \mathbb{A}^{l-1} \xleftarrow{\cdot R'_0} \mathbb{A}'_0,$$

we obtain an “exact” representation of  $\ker(\cdot R_{-1})$  as  $\text{im}(\cdot R'_0)$ . It can be proved that the quotient

$$\text{im}(\cdot R'_0)/\text{im}(\cdot R_0) \subseteq M$$

is the *torsion module*  $t(M)$  of  $M$ , i.e., the set of all its members  $m$  for which there exists a non-zero scalar  $a \in \mathbb{A}$  such that  $am = 0$ . Thus we have obtained that a (linear) control system is controllable if and only if its associated module  $M$  of observables is torsion-free, which can be tested algorithmically. Moreover, a basis for the module  $t(M)$  of autonomous elements is obtained from the rows of  $R'_0$  (that are elements of  $\text{im}(\cdot R'_0)$ ), viewed modulo  $\text{im}(\cdot R_0)$ .

#### 4. More Structural Properties of Control Systems as Extension Modules

Other structural properties of  $\mathcal{D}$  will be described in terms of the *extension modules* of  $N$ , a central tool in homological algebra. Consider a free resolution

$$(6) \quad \dots \xrightarrow{R_{-n}} \mathbb{A}^{l-n} \xrightarrow{R_{-n+1}} \dots \xrightarrow{R_{-2}} \mathbb{A}^{l-2} \xrightarrow{R_{-1}} \mathbb{A}^{l-1} \xrightarrow{R_0} \mathbb{A}^{l_0} \rightarrow N \rightarrow 0$$

(as obtained, for example, with the algorithms of Section 5). This is an exact sequence of right D-modules. By dualization it becomes a sequence

$$(7) \quad \dots \xleftarrow{\cdot R_{-n}} \mathbb{A}^{l-n} \xleftarrow{\cdot R_{-n+1}} \dots \xleftarrow{\cdot R_{-2}} \mathbb{A}^{l-2} \xleftarrow{\cdot R_{-1}} \mathbb{A}^{l-1} \xleftarrow{\cdot R_0} \mathbb{A}^{l_0} \leftarrow \text{hom}_{\mathbb{A}}(N, \mathbb{A}) \leftarrow 0$$

of left D-modules that, again, is usually no longer exact. By dropping  $\text{hom}_{\mathbb{A}}(N, \mathbb{A})$  from (7), we obtain another non-exact sequence, but of *free* modules only,

$$\dots \xleftarrow{\cdot R_{-n}} \mathbb{A}^{l-n} \xleftarrow{\cdot R_{-n+1}} \dots \xleftarrow{\cdot R_{-2}} \mathbb{A}^{l-2} \xleftarrow{\cdot R_{-1}} \mathbb{A}^{l-1} \xleftarrow{\cdot R_0} \mathbb{A}^{l_0} \leftarrow 0.$$

Its defects of exactness are encapsulated by its *cohomology* sequence, that is to say, by the quotients

$$\ker(\cdot R_{-i})/\text{im}(\cdot R_{-i+1}).$$

An all-important fact is that this family depends on  $N$  only, and not of the choice of a free resolution (6). This motivates the notation

$$\text{ext}_{\mathbb{A}}^i(N, \mathbb{A}) = \ker(\cdot R_{-i})/\text{im}(\cdot R_{-i+1})$$

for extension modules (with in particular  $\text{ext}_{\mathbb{A}}^0(N, \mathbb{A}) = \ker(\cdot R_0) = \text{hom}_{\mathbb{A}}(N, \mathbb{A})$ ).

The nullity or non-nullity of the  $\text{ext}^i$ 's provides with the classification of modules in Theorem 1 below; in turn this classification provides with the classification of control systems in Theorem 3 below. Here are two more module-theoretic notions missing to state Theorem 1. A module  $L$  is *projective* whenever there exists a module  $L'$  such that  $L \oplus L'$  is free; it is *reflexive* whenever it is isomorphic to the dual of its dual through the linear map

$$\epsilon : M \rightarrow \text{hom}_{\mathbb{A}}(\text{hom}_{\mathbb{A}}(M, \mathbb{A}), \mathbb{A})$$

defined by

$$\epsilon(m)(f) = f(m).$$

Then, a free module is always projective, a projective module always reflexive, and a reflexive module always torsion-free. (For modules over a principal ideal, these notions coincide; for modules over a multivariate polynomial ring with coefficients over a field, free and projective are equivalent, a theorem by Quillen and Suslin.)

The following theorems [1, 4] make the link between properties of a module and the nullity of the extension modules of its transposed module.

**Theorem 1** (Palamodov, Kashiwara). *For the modules  $M$  and  $N$  defined by (3) and (5), we have:*

1.  $M$  is torsion-free if and only if  $\text{ext}_{\mathbb{A}}^1(N, \mathbb{A}) = 0$ ;
2.  $M$  is reflexive if and only if  $\text{ext}_{\mathbb{A}}^1(N, \mathbb{A}) = \text{ext}_{\mathbb{A}}^2(N, \mathbb{A}) = 0$ ;
3.  $M$  is projective if and only if  $\text{ext}_{\mathbb{A}}^1(N, \mathbb{A}) = \cdots = \text{ext}_{\mathbb{A}}^r(N, \mathbb{A}) = 0$ .

**Theorem 2** (Palamodov, Kashiwara). *Let  $M$  and  $N$  be the two modules defined by (3) and (5). Then there exists an exact sequence*

$$0 \rightarrow M \rightarrow \mathbb{A}^{p_1} \rightarrow \mathbb{A}^{p_2} \rightarrow \cdots \rightarrow \mathbb{A}^{p_r}$$

*if and only if  $\text{ext}_{\mathbb{A}}^i(N, \mathbb{A}) = 0$  for  $i = 1, \dots, r$ .*

We finally obtain the following classification of linear control systems, which admits some refinements in the case of differential operators with constant coefficients, i.e., matrices with entries in  $\mathbb{R}[\partial_1, \dots, \partial_r] \subset \mathbb{A}$  [7, 8, 12].

**Theorem 3.** *For a control system defined by the differential operator  $\mathcal{D} = R \cdot$  where  $R$  is an  $l \times m$  matrix with  $l \leq m$  and entries in*

$$\mathbb{A} = \mathbb{R}(x_1, \dots, x_r)[\partial_1, \dots, \partial_r],$$

*introduce the two left  $D$ -modules  $M = \text{coker}(\cdot R)$  and  $N = \text{coker}(R \cdot)$  of the maps between the free modules  $\mathbb{A}^m$  and  $\mathbb{A}^l$ . Then:*

1. *if  $M$  has torsion, the control system has autonomous elements, and in the event  $R$  has constant coefficients and full row module, it has no primality property;*
2.  *$M$  is torsion-free if and only if  $\text{ext}_{\mathbb{A}}^1(N, \mathbb{A}) = 0$ . In this case, the control system is controllable, and in the event  $R$  has constant coefficients and full row module, it is prime in the sense of minors, i.e., there is no common factor between the minors of  $R$  of order  $l$ ;*
3.  *$M$  is reflexive if and only if  $\text{ext}_{\mathbb{A}}^1(N, \mathbb{A}) = \text{ext}_{\mathbb{A}}^2(N, \mathbb{A}) = 0$ ;*
4. *in the event  $R$  has constant coefficients and full row module, and if*

$$\text{ext}_{\mathbb{A}}^1(N, \mathbb{A}) = \cdots = \text{ext}_{\mathbb{A}}^{r-1}(N, \mathbb{A}) = 0 \quad \text{while} \quad \text{ext}_{\mathbb{A}}^r(N, \mathbb{A}) \neq 0,$$

*the control system is weakly prime in the sense of zeroes, i.e., all minors of order  $l$  simultaneously vanish at finitely many points only;*

5.  *$M$  is projective if and only if*

$$\text{ext}_{\mathbb{A}}^1(N, \mathbb{A}) = \cdots = \text{ext}_{\mathbb{A}}^r(N, \mathbb{A}) = 0.$$

*In this case the control system has an inverse generalized Bézout identity, and in the event  $R$  has constant coefficients and full row module, it is prime in the sense of zeroes, i.e., all minors of order  $l$  simultaneously vanish at no point;*

6. *if  $M$  is free, the control system is flat and has direct and inverse generalized Bézout identities.*

Further intermediate situations,  $\text{ext}_{\mathbb{A}}^1(N, \mathbb{A}) = \cdots = \text{ext}_{\mathbb{A}}^{k-1}(N, \mathbb{A}) = 0$  and  $\text{ext}_{\mathbb{A}}^k(N, \mathbb{A}) \neq 0$ , correspond to further intermediate primeness conditions (described in terms of the dimension of the algebraic variety defined by the  $l \times l$  minors of  $R$ ).

## 5. Gröbner Basis Calculations for Compatibility Conditions and Parametrizations

The whole machinery of the previous sections crucially bases on prolongations of exact sequences. A point that is important in view of computations is that these can be obtained by Gröbner basis calculations for free modules over  $\mathbb{A}$ .

The prolongation of a map  $\mathbb{A}^m \xleftarrow{R} \mathbb{A}^l$  into an exact sequence  $\mathbb{A}^m \xleftarrow{R} \mathbb{A}^l \xleftarrow{S} \mathbb{A}^k$  is done in the following fashion. Let  $(e_1, \dots, e_m)$  and  $(f_1, \dots, f_l)$  be the canonical bases of  $\mathbb{A}^m$  and  $\mathbb{A}^l$ , respectively, and denote the  $i$ th row of  $R = (r_{i,j})$  by  $\eta_i$ . Thus  $\eta_i = \sum_{j=1}^m r_{i,j} e_j$ . Prolonging the map amounts to finding non-trivial relations  $\sum_{i=1}^l s_i \eta_i = 0$ . Now introduce the submodule  $Z$  of  $\mathbb{A}^{m+l}$  generated by the formal linear combinations  $f_i - \eta_i$ . We contend that computing a Gröbner basis for this module and for a term order that *eliminates* the  $e_i$  results in linear combinations  $\sum_{i=1}^l s_i f_i \in Z$ , each of which corresponds to a relation between the  $\eta_i$ . Additionally, any relation can be obtained as a linear combination of the relations thus obtained.

In effect, consider an element  $z = \sum_{i=1}^l s_i f_i \in Z$ ; thus  $\sum_{i=1}^l s_i \eta_i$  is in  $Z$  and is a combination  $\sum_{i=1}^l \lambda_i (f_i - \eta_i)$ , which is only possible, in view of the coefficients of the  $f_i$ , if the  $\lambda_i$  are zero, thus if  $\sum_{i=1}^l s_i \eta_i = 0$ ; the converse property is also true. Since the Gröbner basis calculation precisely computes a finite generating set, say of  $k$  elements, for all the  $z$ 's free of the  $e_i$ , it suffices to consider each of those  $k$  elements as a row, and to glue them in column to obtain a new matrix  $S = (S_{i,j})$  such that the sequence  $\mathbb{A}^m \xleftarrow{R} \mathbb{A}^l \xleftarrow{S} \mathbb{A}^k$  is exact.

Now, existing packages often contain facilities to compute Gröbner bases for left modules only; some of our computations require to deal with right modules. A last ingredient, *adjunction*, enables one to turn any left module into a right module, and vice versa, in a way *that preserves the exactness of sequences*. Indeed, the adjoint map  $P \mapsto \tilde{P}$  defined by associativity from the rules  $\tilde{x}_i = x_i$ ,  $\tilde{\partial}_i = -\partial_i$ , and  $(PQ)^\sim = \tilde{Q}\tilde{P}$ , is an (anti)automorphism of the algebra  $\mathbb{A}$  which extends to matrices by mapping itself to the entries of the transpose matrix. Thus, for example, the exact sequence (5) of right D-modules of columns in Section 3 is replaced with the exact sequence

$$\mathbb{A}^{l-1} \xrightarrow{\tilde{R}_0} \mathbb{A}^{l_0} \rightarrow \tilde{N} = \text{coker}(\cdot \tilde{R}_0) \rightarrow 0$$

of left D-modules of lines, for the purpose of explicit calculations.

## Bibliography

- [1] Kashiwara (Masaki). – Algebraic study of systems of partial differential equations. *Mémoires de la Société Mathématique de France*, n° 63, 1995. – Japanese translation of the author's thesis (1970). xiv+72p.
- [2] Malgrange (Bernard). – Systèmes différentiels à coefficients constants. In *Séminaire Bourbaki*, Vol. 8, pp. 79–89. – Société Mathématique de France, Paris, 1995. Exposé n° 246.
- [3] Oberst (Ulrich). – Multidimensional constant linear systems. *Acta Applicandae Mathematicae*, vol. 20, n° 1-2, 1990, pp. 1–175.
- [4] Palamodov (V. P.). – *Linear differential operators with constant coefficients*. – Springer-Verlag, New York, 1970. viii+444 pages.
- [5] Polderman (Jan Willem) and Willems (Jan C.). – *Introduction to mathematical systems theory*. – Springer-Verlag, New York, 1998. A behavioral approach. xxx+424 pages.
- [6] Pommaret (J. F.) and Quadrat (A.). – Generalized Bézout identity. *Applicable Algebra in Engineering, Communication and Computing*, vol. 9, n° 2, 1998, pp. 91–116.
- [7] Pommaret (J. F.) and Quadrat (A.). – Algebraic analysis of linear multidimensional control systems. *IMA Journal of Mathematical Control and Information*, vol. 16, n° 3, 1999, pp. 275–297.
- [8] Pommaret (J. F.) and Quadrat (A.). – Localization and parametrization of linear multidimensional control systems. *Systems & Control Letters*, vol. 37, n° 4, 1999, pp. 247–260.
- [9] Pommaret (J. F.) and Quadrat (A.). – Formal elimination for multidimensional systems and applications to control theory. *Mathematics of Control, Signals, and Systems*, vol. 13, n° 3, 2000, pp. 193–215.

- [10] Pommaret (J. F.) and Quadrat (A.). – A functorial approach to the behaviour of multidimensional control systems. In *The Second International Workshop on Multidimensional (nD) Systems (Czocha Castle, 2000)*, pp. 91–96. – Tech. Univ. Press, Zielona Góra, 2000.
- [11] Quadrat (Alban). – A fractional representation approach of synthesis: an algebraic analysis point of view. *SIAM Journal on Optimization and Control*. – To appear.
- [12] Quadrat (Alban). – *Analyse algébrique des systèmes de contrôle linéaires multidimensionnels*. – Thèse de doctorat, École nationale de Ponts et Chaussées, September 1999.
- [13] Spencer (D. C.). – Overdetermined systems of linear partial differential equations. *Bulletin of the American Mathematical Society*, vol. 75, 1969, pp. 179–239.
- [14] Wood (Jeffrey). – Modules and behaviours in  $nd$  systems theory. *Multidimensional Systems and Signal Processing*, vol. 11, n° 1-2, 2000, pp. 11–48. – Recent progress in multidimensional control theory and applications.



## Effective Test of Local Algebraic Observability — Applications to Systems and Control Theory

*Alexandre Sedoglavic*

GAGE, École polytechnique (France)

December 4, 2000

### Abstract

In control and systems theory, the problem of structural algebraic observability consists in deciding whether the state variables involved in a model can be determined in terms of commands and measures supposed perfectly known. Structural identifiability is a variant where one tries to know whether the parameters of a model are observable.

We propose a probabilistic algorithm with polynomial complexity to answer the question in the ordinary differential framework. This algorithm relies on seminumerical techniques (modular computations, series expansions, and Newton operator) that allow the computation of the generic rank of the Jacobian matrix of measures and their derivatives with respect to time.

To conclude, we present experimental results that illustrate the notion of algebraic observability and show the efficiency of our approach.



## Part IV

# Probabilistic Methods



# Reflected Brownian Bridge Area Conditioned on its Local Time at the Origin

Guy Louchard

Université libre de Bruxelles, Bruxelles (Belgique)

June 25, 2001

Summary by Michel Nguyen-Thé

## Abstract

Using properties of the Airy functions, we analyze the reflected Brownian bridge area  $W_b$  conditioned on its local time  $b$  at the origin. We give a closed form expression of the Laplace transform of  $W_b$ , a recurrence equation for the moments, leading to an efficient computation algorithm and an asymptotic form for the density  $f(x, b)$  of  $W_b$  for  $x \rightarrow 0$ .

## 1. Introduction

Let us first introduce the standard Brownian motion denoted by  $x(t)$  and a few classical variants: the reflected Brownian motion  $x^+(t) = |x(t)|$ ; the Brownian bridge  $B(t)$ ; the reflected Brownian bridge  $B^+(t)$  on  $[0, 1]$ ; the Brownian excursion  $e(t)$ .

The object of interest in this talk is  $W_b := \int_0^1 B^+(t) dt$ , the area of the reflected Brownian bridge conditioned on having a local time at the origin equal to  $b$ . This random variable appeared in [4] as the limit law for  $m^{-3/2} D_{m, m-b\sqrt{m}}$ , where  $D_{m, m-b\sqrt{m}}$  denotes the total displacement for a hash table with  $m$  locations and  $b\sqrt{m}$  empty locations, using linear probing. It also represents the limit law for the total height of random forests with  $b\sqrt{m}$  trees and  $m$  nodes or leaves. The only description of it was given by its moments, related to the classical Airy function  $\text{Ai}(z) := \frac{1}{\pi} \int_0^{+\infty} \cos(\frac{1}{3}t^3 + zt) dt$  (recall  $\text{Ai}'' = z\text{Ai}$ ) in the following way:

$$\mathbf{E}[W_b^k] = k! \sum_{j=1}^k \left( \sum_{k_1, \dots, k_j \geq 1, \Sigma k_i = k} \prod_{i=1}^j \omega_{k_i} \right) \frac{b^{j-1}}{j!} q_{3k-j-2}(b),$$

where the  $\omega_k$  are defined by the asymptotic expansion  $\frac{\text{Ai}'(z)}{\text{Ai}(z)} \underset{z \rightarrow +\infty}{\sim} \sum_{k=0}^{+\infty} \omega_k \frac{(-1)^k z^{-3(k-1)/2}}{2^k}$ , and  $q_r(b) := \int_0^{+\infty} \frac{x^r}{r!} e^{-bx-x^2/2} dt$ .

We will provide a closed form expression for the Laplace transform of  $W_b$ , a better way to compute its moments, and an asymptotic form for the density  $f(x, b)$  of  $W_b$  when  $x \rightarrow 0$ .

## 2. Laplace Transform of $W_b$

Computing the Laplace transform of  $W_b$  essentially requires using Kac's formula [3] and a few technicalities. Eq. (30) in [5, p. 491] states that, if we denote by  $t^+(t, a)$  the local time of  $x(t)$  at  $a$ ,

$$(1) \quad \int_0^\infty e^{-\alpha t} \mathbf{E}_0 \left[ \exp\left(-\int_0^t x^+(u) du - \delta t^+(t, 0)\right) \middle| x(t) = 0 \right] \frac{dt}{\sqrt{2\pi t}} = \left( \delta - \frac{2^* \text{Ai}'(2^* \alpha)}{\text{Ai}(2^* \alpha)} \right)^{-1},$$

where  $2^* := 2^{1/3}$ . From it we can derive the following theorem:

**Theorem 1.** *The Laplace transform  $\Theta(z, b)$  of  $W_b$  has the closed form expression*

$$\Theta(z, b) = \mathbf{E}[e^{-zW_b}] = \frac{-z^{1/3}e^{b^2/2}}{i2^{1/6}\sqrt{\pi}} \int_{-i\infty}^{i\infty} e^{bz^{1/3}2^{1/3}Ai'(u)/Ai(u)} (Ai'(u)/Ai(u))' e^{uz^{2/3}/2^{1/3}} du.$$

*Proof.* Given  $[\int_0^t x^+(u) du | x(t) = 0] \stackrel{\mathcal{D}}{=} t^{3/2}Y$  and  $t^+(t, 0) \stackrel{\mathcal{D}}{=} \sqrt{t}t^+(1, 0)$  (scaling property), Eq. (1) leads to

$$\mathbf{E}_0 \int_0^\infty e^{-\alpha t} \int_0^\infty e^{-t^{3/2}W_b} b e^{-b^2/2} e^{-\delta\sqrt{t}b} \frac{db dt}{\sqrt{2\pi t}} = [\delta - 2^*\Lambda(\alpha)]^{-1},$$

where  $\Lambda(\alpha) := \frac{Ai'(2^*\alpha)}{Ai(2^*\alpha)}$ . The change of variable  $v = \sqrt{t}b$  and an inversion on  $\delta$  delivers

$$(2) \quad \int_0^\infty e^{-b^2/2} e^{-\alpha v^2/b^2} \mathbf{E} \left[ e^{-v^3/b^3 W_b} \right] \frac{2 db}{\sqrt{2\pi}} = e^{v2^*\Lambda(\alpha)}.$$

After setting  $b = \frac{v}{\sqrt{2^*\alpha}}$ ,  $u = 2^*\alpha$ , differentiating with respect to  $u$  and using  $(\frac{Ai'}{Ai})' = u - (\frac{Ai'}{Ai})^2$ :

$$\frac{1}{\sqrt{2\pi}} \int_0^\infty e^{-u\sigma} \mathbf{E} \left[ e^{-\sqrt{2}\sigma^{3/2}W_{v/\sqrt{2^*\sigma}}} \right] e^{-v^2/(2^{4/3}\sigma)} \frac{d\sigma}{\sqrt{2\sigma}} = -e^{v2^*Ai'(u)/Ai(u)} (Ai'(u)/Ai(u))'.$$

The inversion formula for Laplace transforms then writes:

$$(3) \quad \mathbf{E} \left[ e^{-\sqrt{2}\sigma^{3/2}W_{v/\sqrt{2^*\sigma}}} \right] e^{-v^2/(2^{4/3}\sigma)} / \sqrt{4\pi\sigma} = \frac{-1}{2\pi i} \int_{-i\infty}^{i\infty} e^{v2^*Ai'(u)/Ai(u)} (Ai'(u)/Ai(u))' e^{u\sigma} du.$$

Now set  $v = b\sqrt{2^*\sigma}$ ,  $z = \sqrt{2}\sigma^{3/2}$ ,  $\Theta(z, b) = \mathbf{E}[e^{-zW_b}]$ . Eq. (3) becomes

$$\frac{2^{1/6}\Theta(z, b)e^{-b^2/2}}{2\sqrt{\pi}} = \frac{-z^{1/3}}{2\pi i} \int_{-i\infty}^{i\infty} e^{bz^{1/3}2^*Ai'(u)/Ai(u)} (Ai'(u)/Ai(u))' e^{uz^{2/3}/2^*} du$$

which proves the theorem.  $\square$

### 3. Recurrence Formulae

Using Laplace transforms and inversions of Laplace transforms, we show here how to find an algorithm to compute the moments  $\psi_k(b) := \mathbf{E}[W_b^k]$  by recurrence. We first need:

**Lemma 1.** *Define  $G(\eta) := 2^*\Lambda(\alpha)/\sqrt{\alpha}$  and  $s = 1/b^2$ ; we have*

$$(4) \quad \int_0^\infty e^{-1/(2s)} e^{-ws} (-1)^k s^{3/2k} \psi_k(b) \frac{ds}{s^{3/2}\sqrt{2\pi k!}} = [\eta^k] \frac{e^{\sqrt{w}G_0}}{w^{3/2k}} \sum_{i=1}^\infty \frac{(\sqrt{w}(G(\eta) - G_0))^i}{i!}.$$

*Proof.* Set  $s := 1/b^2$ ,  $w = \alpha v^2$ , and  $\eta = \alpha^{-3/2}$ . Eq. (2) becomes

$$\int_0^\infty e^{-1/(2s)} e^{-ws} \mathbf{E} \left[ e^{-\eta w^{3/2} s^{3/2} W_b} \right] \frac{ds}{s^{3/2}\sqrt{2\pi}} = e^{\sqrt{w}G(\eta)},$$

Set  $G_0 := G(0)$ . Eq. (3) leads to

$$(5) \quad \int_0^\infty e^{-1/(2s)} e^{-ws} \mathbf{E} \left[ e^{-\eta w^{3/2} s^{3/2} W_b} - 1 \right] \frac{ds}{s^{3/2}\sqrt{2\pi}} = e^{\sqrt{w}G(\eta)} - e^{\sqrt{w}G_0} \\ = e^{\sqrt{w}G_0} \sum_{i=1}^\infty \frac{(\sqrt{w}(G(\eta) - G_0))^i}{i!}.$$

Upon expanding both sides of (5) with respect to  $\eta$ , this gives the desired formula.  $\square$

To invert the Laplace transforms of the form  $e^{-\sqrt{2w}}/w^{(j+1)/2}$ , we will use the following lemmas:

**Lemma 2.** Set  $\phi^{(1)}(x) := \phi(x) := \frac{1}{\sqrt{2\pi}} \int_{-\infty}^x e^{-t^2/2} dt$  (classical Gaussian distribution function) and  $\phi^{(j+1)}(x) := \int_{-\infty}^x \phi^{(j)}(u) du$ . Then

$$\int_0^\infty \phi^{(j)}(-b)e^{-ws} \frac{(2s)^{(j+1)/2}}{s} ds = \frac{e^{-\sqrt{2w}}}{w^{(j+1)/2}}, \quad j \geq 1, \quad \text{where } b = 1/\sqrt{s}.$$

*Sketch of proof.* One proves the lemma by induction and uses an integration by part and an integration with respect to  $w$  to prove it at rank  $k + 1$  from rank  $k$ .  $\square$

**Lemma 3.** The  $\phi^{(j)}(x)$  can be expressed in the form:

$$\phi^{(k)}(z) = p_1(k, z)\phi(z) + p_2(k, z)e^{-z^2/2}/\sqrt{2\pi},$$

where  $p_1(k, z)$  is of degree  $k - 1$ ,  $p_2(k, z)$  is of degree  $k - 2$ .

Using integration by parts on  $\int_{-\infty}^z x^j \phi(x) dx$  and identification of coefficients, it is possible to prove the following proposition, enabling us to compute nice expressions of the  $\phi^{(j)}(x)$ :

**Proposition 1.** Define, for  $k \geq 1$ ,  $j \geq 0$ ,  $P_1[k, j] := [z^j]p_1(k, z)$ , and  $P_2[k, j] := [z^j]p_2(k, z)$ . Then the sequences  $(P_1[k, j])_{k \geq 1, j \geq 0}$  and  $(P_2[k, j])_{k, j \geq 0}$  are defined by the initial values  $P_1[1, 0] = 1$ ,  $P_2[1, 0] = 0$ ,  $P_1[1, j] = P_2[1, j] = 0$  for  $j \geq 1$ , and the recurrence relations, for  $k \geq 1$ :

$$\begin{aligned} P_1[k+1, j] &:= P_1[k, j-1]/j, \quad j = 1, \dots, k, \\ P_2[k+1, j] &:= \sum_{l=0}^{\lfloor (k-1-j)/2 \rfloor} P_1[k, j+2l]/(j+2l+1)(j+2l+1)_l \\ &\quad - \sum_{l=0}^{\lfloor (k-3-j)/2 \rfloor} P_2[k, j+2l+1](j+2l+1)_l, \quad j = 0, \dots, k-1, \\ P_1[k+1, 0] &:= - \sum_{l=1,3,\dots,k-1} P_1[k, l]/(l+1)(l+1)_{(l+1)/2} + \sum_{l=0,2,\dots,k-2} P_2[k, l](l)_l. \end{aligned}$$

Determining a recurrence relation for the moments  $\psi_k(b)$  hence amounts to determining a recurrence relation for the  $Z_j$  defined by (see (4)):

$$(-1)^j b^{-3j} \frac{Z_j}{j!} = [\eta^j] \frac{1}{w^{3/2j}} \sum_{i=1}^{\infty} \frac{(\sqrt{w}(G(\eta) - G_0))^i}{i!}.$$

Indeed, along the mechanical transfer rule  $\frac{1}{w^{(l+1)/2}} \rightarrow \frac{\phi^{(l)}(-b)}{b^{l+1}} b^2 2^{(l+1)/2}$ ,  $\psi_j(b)$  is equivalent to  $Z_j \sqrt{2\pi} e^{b^2/2} / b^3$ . To get a recurrence formula giving  $Z_k$  in function of the  $Z_1, \dots, Z_j$ , we introduce

$$S_k(\eta) := \sum_{j=1}^k (-1)^j b^{-3j} \frac{Z_j}{j!} w^{3j/2} \eta^j = \sum_{j=1}^k \eta^j [\eta^j] \left( \frac{\sum_{l=1}^k (-1)^l (d_l - c_l) \left(\frac{3\eta}{2^{3/2}}\right)^l}{\sum_{l=0}^k (-1)^l c_l \left(\frac{3\eta}{2^{3/2}}\right)^l} \right)^j \frac{(-\sqrt{2w})^j}{j!},$$

where the coefficients  $c_l$  and  $d_l$  are defined in [1, Eq. (10.4.59) and (10.4.61)] by asymptotic expansions of Ai and Ai' for  $|z|$  large,  $|\arg(z)| < \pi$ :

$$\text{Ai}(z) \sim \frac{1}{2\sqrt{\pi}} z^{-1/4} e^{-\zeta} \sum_{k=0}^{\infty} (-1)^k c_k \zeta^{-k}, \quad \text{Ai}'(z) \sim -\frac{1}{2\sqrt{\pi}} z^{1/4} e^{-\zeta} \sum_{k=0}^{\infty} (-1)^k d_k \zeta^{-k},$$

with  $\zeta := \frac{2}{3}z^{3/2}$ . More explicitly:  $c_0 = 1$ ,  $c_k = \Gamma(3k + 1/2)/(\Gamma(k + 1/2) \cdot 54^k k!)$ ,  $d_0 = 1$ ,  $d_k = -\frac{6k+1}{6k-1}c_k$ . The relation

$$[\eta^k] \sum_{j=1}^k (-1)^j b^{-3j} \frac{Z_j}{j!} w^{3j/2} \eta^j \left( \sum_{l=0}^k (-1)^l c_l \left( \frac{3\eta}{2^{3/2}} \right)^l \right)^k \\ [\eta^k] \sum_{j=1}^k \left( \frac{-\sqrt{2}}{z} \right)^j \frac{1}{j!} \left( \sum_{l=1}^k (-1)^l (d_l - c_l) \left( \frac{3\eta}{2^{3/2}} \right)^l \right)^j \left( \sum_{l=0}^k (-1)^l c_l \left( \frac{3\eta}{2^{3/2}} \right)^l \right)^{k-j}$$

provides an algorithm that can easily be implemented in Maple and proves more tractable than the general expressions of the moments given by Janson.

#### 4. Asymptotic Form of Density

**4.1. Asymptotics of  $f(x, b)$  as  $b \rightarrow \infty$ .** Using  $\mathbf{E}[W_b] \sim \frac{1}{2b}$  and  $\mathbf{Var}[W_b] \sim \frac{1}{4b^4}$  as  $b \rightarrow \infty$ , already mentioned in [4], asymptotics of  $(\log \text{Ai})'$  and  $(\log \text{Ai})''$ , and a saddle point method, we recover the fact that we obtain a density of a Gaussian distribution when  $b \rightarrow \infty$ .

**4.2. Asymptotics of  $\Theta(z, b)$  as  $|z| \rightarrow \infty$ .** Using a saddle point again, setting  $z = \kappa^6$ , we obtain

$$\Theta \sim e^{\kappa^3 \mu_1} e^{-\alpha_1 \kappa^4 / 2^*} \left( \frac{2^{1/2} \kappa^{3/2}}{2b^{3/4}} + \frac{b^{1/4} 2^{1/6} \alpha_1}{4\kappa^{1/2}} + \mathcal{O}\left(\frac{1}{\kappa^{3/2}}\right) \right).$$

**4.3. Asymptotics of  $f(x, b)$  as  $x \rightarrow 0$ .** The formula  $f(x, b) = \frac{1}{2\pi i} \Re \int_{c-i\infty}^{c+i\infty} e^{xz} \Theta(z, b) dz$ ,  $c > 0$ , the former asymptotics and a saddle point method lead to:

$$f(x, b) \sim e^{\mu_2/x^2} \frac{\sqrt{2}}{\sqrt{\pi}} \left( \frac{3^{1/4} \alpha_1^{9/4}}{9x^{11/4} b^{3/4}} - \frac{3^{3/4} \alpha_1^{3/4}}{3x^{9/4} b^{1/4}} + \frac{b^{1/4} 3^{1/4} (27 + 16\alpha_1^3)}{x^{7/4} \alpha_1^{3/4}} + \mathcal{O}\left(\frac{1}{x^{5/4}}\right) \right).$$

#### 5. Open Questions

It remains to find an asymptotic form for the density  $f(x, b)$  as  $x \rightarrow \infty$ —this not even known for the classical Airy density—and an explicit form for the density  $f(x, b)$ . Are also missing an analysis of the local time  $t^+(t, a)$  of  $B^+(t)$  at  $a$ , conditioned on its local time  $b$  at the origin, and some analytic variations on  $W_b$  (see [2] for the classical Airy distribution).

#### Bibliography

- [1] Abramowitz (Milton) and Stegun (Irene A.) (editors). – *Handbook of mathematical functions, with formulas, graphs, and mathematical tables*. – Dover Publications Inc., New York, 1966, xiv+1046p.
- [2] Flajolet (P.) and Louchard (G.). – Analytic variations on the Airy distribution. *Algorithmica*, vol. 31, n° 3, 2001, pp. 361–377. – Mathematical analysis of algorithms.
- [3] Itô (Kiyosi) and McKean, Jr. (Henry P.). – *Diffusion processes and their sample paths*. – Springer-Verlag, Berlin, 1974, xv+321p. Second printing, corrected, Die Grundlehren der mathematischen Wissenschaften, Band 125.
- [4] Janson (Svante). – Asymptotic distribution for the cost of linear probing hashing. *Random Structures & Algorithms*, vol. 19, n° 3-4, 2001, pp. 438–471.
- [5] Louchard (G.). – Kac's formula, Levy's local time and Brownian excursion. *Journal of Applied Probability*, vol. 21, n° 3, 1984, pp. 479–499.



## Cover Time and Favourite Points for Planar Random Walks

Amir Dembo

Mathematics and Statistics Department, Stanford University (USA)

June 18, 2001

Summary by Christine Fricker and Pierre Nicodème

### Abstract

In this talk, Amir Dembo considers random walks on  $\mathbb{Z}^2$  and presents a proof of the Erdős–Taylor conjecture related to frequently covered points. The Kesten–Révész conjecture on the covering time of the two-dimensional torus  $\mathbb{Z}_n^2 = \mathbb{Z}^2/n\mathbb{Z}^2$  is also solved. These results are a common work of Amir Dembo, Yuval Peres, Jay Rosen, and Ofer Zeitouni.

### 1. Introduction

Let  $(X_n)$  be a simple random walk on  $\mathbb{Z}^2$  and  $T_n(x) = \sum_{j=1}^n 1_{\{X_j=x\}}$  be the number of visits to  $x$  before time  $n$ . Let  $T_n^* = \max_{x \in \mathbb{Z}^2} T_n(x)$  be the number of visits to the most visited point. The *Erdős–Taylor conjecture* asserts that

$$(1) \quad \lim_{n \rightarrow \infty} \frac{T_n^*}{(\log n)^2} = \frac{1}{\pi}, \quad \text{almost surely.}$$

Erdős and Taylor [7] proved the upper bound  $1/\pi$  and a lower bound  $1/(4\pi)$ . The main result of the talk is that the Erdős–Taylor conjecture is true.

Let  $(\tilde{X}_j)$  be a simple random walk on the two-dimensional torus  $\mathbb{Z}_n^2 = \mathbb{Z}^2/n\mathbb{Z}^2$ . Consider  $\mathcal{T}(x) = \min\{j \geq 0 \mid \tilde{X}_j = x\}$ , the time to attain the point  $x$  for the first time and

$$\mathcal{T}_n = \max_{x \in \mathbb{Z}_n^2} \mathcal{T}(x),$$

the covering time of the torus. The *Aldous–Lawler conjecture* asserts that

$$(2) \quad \lim_{n \rightarrow \infty} \frac{\mathcal{T}_n}{(n \log n)^2} = \frac{4}{\pi}, \quad \text{in probability.}$$

Kesten, Révész, Lawler, and Aldous proved an upper bound  $4/\pi$  (see [1, Corollary 25, Chapter 7]) and a lower bound  $2/\pi$ . A related question is the Kesten–Révész conjecture for the simple random walk on  $\mathbb{Z}^2$  (see [4]).

The proofs for the upper bounds rely on the second moment method, the approximation of random walks by Brownian motions, and an underlying tree structure for the occupation of small disks by a Brownian motion. We give here a sketch of the proofs; see [4, 5] for complete proofs.

### 2. The Second Moment Method

Janson gives a short account of the second moment method in [2]. Basically, we consider a sequence of non-negative random variables  $X_n$ , and we want to estimate  $\mathbf{P}(X_n > 0)$ . The second

moment method asserts that if

$$(3) \quad \frac{\mathbf{Var}(X_n)}{(\mathbf{E}X_n)^2} \rightarrow 0, \quad \text{or equivalently,} \quad \frac{\mathbf{E}X_n^2}{(\mathbf{E}X_n)^2} \rightarrow 1 \quad (\text{as } n \rightarrow \infty),$$

then

$$(4) \quad \mathbf{P}(X_n > 0) \rightarrow 1.$$

The method is frequently used in the context of random graphs; for example, this method proves the existence of a Hamilton cycle in random graphs satisfying suitable conditions.

The second moment method is a consequence of the Chebyshev inequality,

$$\mathbf{P}(|X| > t) \leq \frac{1}{t^2} \mathbf{E}(X^2).$$

As a consequence of the latter,

$$\mathbf{P}(X = 0) \leq \mathbf{P}(|X - \mu| \geq \mu) \leq \frac{\mathbf{Var}(X)}{\mu^2}, \quad \text{for } \mu = \mathbf{E}X.$$

### 3. Proof of the Erdős–Taylor Conjecture

**3.1. Upper bound.** By definition, the truncated Green function  $G_n(x, y)$  is the expectation of the number of passages at  $y$  in  $n$  steps, when starting from  $x$ .

We have

$$G_n(0, 0) = \sum_{j=0}^n \mathbf{E} \left( 1_{\{X_j=0\}} \right) = \sum_{j=0}^n \mathbf{P}(X_j = 0) \sim \frac{\log n}{\pi}.$$

(See Feller [8, p. 361].) Applying [3, Theorem 8.7.3] for the renewal sequence  $u_n = \mathbf{P}(X_n = 0)$ , we deduce that for large  $n$ , and fixed small  $\delta > 0$ ,

$$\mathbf{P}(X_j \neq 0 \text{ for } j = 1, \dots, n-1) \leq \frac{(1-\delta)\pi}{\log n}.$$

This implies by the strong Markov property that

$$(5) \quad \mathbf{P}(T_n(0) \geq \alpha\pi(\log n)^2) \leq \left( 1 - \frac{(1-\delta)\pi}{\log n} \right)^{\alpha(\log n)^2} \leq e^{-\alpha\pi(\log n)(1-\delta)} = n^{-(1-\delta)\alpha\pi}.$$

We now consider the disk of center zero and radius  $n^{(1+\delta)/2}$ . The probability that the random walk exits this disk before time  $n$  tends to zero as  $n$  tends to infinity, and the number of points of  $\mathbb{Z}^2$  inside this disk is close to  $\pi n^{(1+\delta)}$ . From Equation (5), we then get

$$(6) \quad \mathbf{P}_n^\alpha \leq \mathbf{P} \left( \max_{0 \leq i \leq n} |X_i| > n^{(1+\delta)/2} \right) + \pi n^{(1+\delta)} n^{-(1-\delta)\alpha\pi},$$

where  $\mathbf{P}_n^\alpha = \mathbf{P}(T_n^* \geq \alpha(\log n)^2)$ . The first term of the right member of Equation (6) vanishes as  $n$  tends to infinity. Therefore, applying the Borel–Cantelli lemma to the subsequence  $\mathbf{P}_{2^m}^\alpha$ , for  $\alpha > 1/\pi$ , and using interpolation for all  $n$ , we have  $\mathbf{P}(\overline{\lim} T_n^* \geq \alpha\pi(\log n)^2) \rightarrow 0$ . This gives an upper bound  $1/\pi$ .

**3.2. Lower bound.** We can try to adapt the proof from the upper bound and use the second moment method. Let  $D(x, r)$  be the disk of center  $x$  and radius  $r$  and

$$Z_n = \sum_{x \in D(0, \sqrt{n})} 1_{\{T_n(x) \geq \beta(\log n)^2\}}.$$

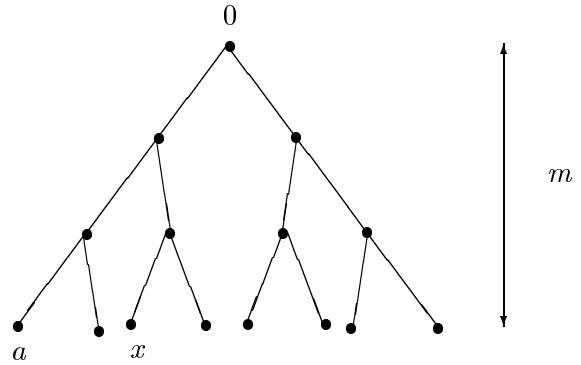
Adapting the proof of the upper bound (Equation (6)) gives  $\mathbf{E}Z_n \approx n^{(1-\beta\pi)}$ . Therefore,

$$\frac{\mathbf{E}Z_n^2}{(\mathbf{E}Z_n)^2} = \frac{1}{\mathbf{E}(Z_n)} + \frac{\Sigma_{x,y}}{\Sigma_{x,y} + \Sigma_x}, \quad \text{where} \quad \Sigma_x = \sum_{x \in D(0, \sqrt{n})} \left( \mathbf{P}(T_n(x) \geq \beta(\log n)^2) \right)^2$$

and  $\Sigma_{x,y} = \sum_{x \neq y \in D(0, \sqrt{n})} \mathbf{P}(T_n(x) \geq \beta(\log n)^2) \mathbf{P}(T_n(y) \geq \beta(\log n)^2).$

A naive approach would say the following: the number of summand in  $\Sigma_{x,y}$  is  $O(n^{2(1-\beta\pi)})$  while it is only  $O(n^{(1-\beta\pi)})$  in  $\Sigma_x$ . Therefore, for  $\beta < 1/\pi$ ,  $\mathbf{E}Z_n^2/(\mathbf{E}Z_n)^2 \rightarrow 1$  and  $\mathbf{P}(T_n^* \geq \frac{1}{\pi}(\log n)^2) = 1$  almost surely. However, Erdős and Taylor [7] show that the correlation structure between points  $x$  such that  $\mathbf{P}(T_n(x) \geq \beta(\log n)^2)$  is too strong to get this result. They obtain an upper limit  $1/(4\pi)$ . We move in the following section to a tree model to overcome this difficulty.

*Modelling by a (toy) tree problem.* We<sup>1</sup> consider a complete binary tree  $B_m$  of height  $m$  and a (nearest neighbor) random walk  $X$  starting from the left-most leaf  $a$ , with probability  $1/3$  of choosing any direction when being at an internal node. In this model, the starting point  $a$  and the root  $0$  respectively represent the origin  $(0,0)$  and the boundary of a “disk” of radius  $m$  on  $\mathbb{Z}^2$ . Let  $L_m$  be the set of leaves of  $B_m$ . We consider  $T_m(x)$ , the time spent at leaf  $x$  before hitting the root  $0$ , and



$$T_m^* = \max_{x \in L_m} T_m(x),$$

its maximum over all leaves.

Let us denote by  $0, 1, 2, \dots, a = m$  the nodes of the ray going from the root  $0$  to  $a$  and let  $\mathbf{P}^y$  denote probability for walks starting from node  $y$ . We consider

$$H_y = H_y(u) = \sum_{u \geq 0} \mathbf{P}^y (X \text{ spends time } k \text{ at } a \text{ before hitting } 0) u^k.$$

For any node  $i$  of the ray  $(0, a)$ , and for any node  $y$  of the subtree rooted at the right child of  $i$ , the probability of  $k$  visits to  $a$  before hitting  $0$  of the walk starting from  $y$  is the same as if the walk starts from  $i$ ; this implies  $H_y = H_i$ . This last result is true for all  $i$  from  $1$  to  $m - 1$ .

We can therefore consider only the nodes of the ray  $(0, a)$ , which provide the set of equations

$$H_1 = \frac{H_2}{3} + \frac{H_1}{3} + \frac{1}{3}, \quad H_k = \frac{H_{k-1}}{3} + \frac{H_k}{3} + \frac{H_{k+1}}{3} \quad (2 \leq k \leq m-2), \quad H_{m-1} = \frac{H_{m-2}}{3} + \frac{(1+u)H_{m-1}}{3}.$$

<sup>1</sup>The elementary proof leading to Equation (7) was not presented by the speaker and is due to the authors of the summary.

Solving yields

$$(7) \quad H_a(u) = H_m = \frac{1}{m} \times \frac{1}{1 - (1 - \frac{1}{m})u}, \quad \text{and} \quad H_1(u) = \frac{m-1 - (m-2)u}{m - (m-1)u}.$$

The random variable  $T_m(a)$  therefore has a geometric distribution with mean  $m-1$ , which induces (for large  $m$ )

$$\mathbf{P}(T_m(a) > \alpha m^2) = \left( \left(1 - \frac{1}{m}\right)^m \right)^{\alpha m} \simeq e^{-\alpha m} \quad \text{and} \quad \mathbf{P}(T_m^* > \alpha m^2) \leq e^{-\alpha m} 2^m = e^{-(\alpha - \log 2)m}.$$

This implies the same upper bound as precedently (up to the change of model).

We now consider a variation of the second moment method. We fix some  $K$  large. We denote by  $x$ -ray the ray from the root  $0$  to a leaf  $x$  and  $N_i(x)$  counts the number of excursions from level  $i$  to level  $i+1$  on the ray  $x$ . We define the  $x$ -ray as  $\alpha$ -successful if

$$N_i(x) \simeq \alpha i^2, \quad \text{for } i = 0, K, 2K, \dots, K \left\lfloor \frac{m}{K} \right\rfloor.$$

We have

$$\mathbf{P}(N_{i+K}(x) \simeq \alpha(i+K)^2 \mid N_i(x) \simeq \alpha i^2) \simeq e^{-\alpha K} \quad \Rightarrow \quad \mathbf{P}(x\text{-ray is } \alpha\text{-successful}) \simeq e^{-\alpha m}.$$

We now have

$$\mathbf{P}(x\text{-ray and } y\text{-ray are } \alpha\text{-successful}) \simeq e^{-2\alpha m} e^{\alpha r(x,y)},$$

where  $r(x, y)$  is the depth of the first common ancestor of  $x$  and  $y$ . This induces a reduction of variance. Considering now the random variable  $Z_m$  defined by

$$Z_m = \sum_{x \in L_m} 1_{\{x\text{-ray } \alpha\text{-successful}\}},$$

we have

$$\frac{\mathbf{E}Z_m^2}{(\mathbf{E}Z_m)^2} \simeq \sum_{s=1}^{m/K} e^{(\alpha - \log 2)Ks} \rightarrow 1 \quad \text{for } \alpha < \log 2,$$

when first  $m$  and then  $K$  tend to infinity. There is no obvious way to adapt this result to the standard random walk, but it is possible to adapt it to the planar Brownian motion that we denote  $w = (w_t)$ .

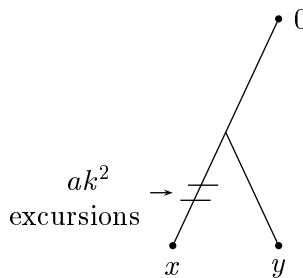
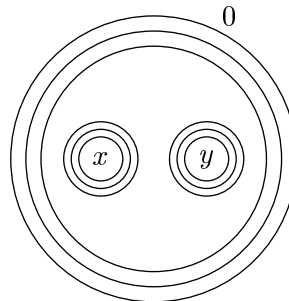
Define  $\theta$  as the first time where the Brownian motion  $w$  hits the circle of radius 1 and  $\mu_\theta^w(A)$  as the occupation time of a subset  $A$  of the disc  $D(0, 1)$  until this time. We have

$$\theta = \min\{t \mid |w_t| = 1\}$$

$$\text{and } \mu_\theta^w(A) = \int_0^\theta 1_A(w_t) dt.$$

The Perkins–Taylor conjecture states for the Brownian motion that

$$(8) \quad \limsup_{\epsilon \rightarrow 0} \sup_{|x| < 1} \frac{\mu_\theta^w(D(x, \epsilon))}{\epsilon^2 (\log \epsilon)^2} = 2.$$



We shall in a first time sketch a proof of this conjecture and apply then the KMT approximation theorem of the Brownian motion by the standard random walk.

*Sketch of proof for the Perkins–Taylor conjecture.* In the following, let  $\partial D(x, r)$  be the boundary of the disk  $D(x, r)$ .

The proof of the upper bound of the conjecture follows the same line as for the standard random walk. When considering the lower bound, the difficulty relies again in the correlation structure.

Let  $\epsilon_k = e^{-k}$  and define a point  $x$  of  $D(0, 1)$  as  $k$ -successful if the number of excursions of the Brownian motion between  $\partial D(x, \epsilon_k)$  and  $\partial D(x, \epsilon_{k+1})$  is  $ak^2$  for fixed  $a$ . We remark that if  $x$  is successful, the time spent at the ball  $D(x, \epsilon_{k+1})$  is  $ak^2\epsilon^2 \simeq a\epsilon^2(\log \epsilon)^2$ , where  $\epsilon = \epsilon_{k+1}$ , with probability close to 1.

*KMT approximation theorem.* The Komlós–Major–Tusnády (KMT) approximation theorem [9] states that for each  $n$  it is possible to construct a random walk  $\{X_k\}_{k=1}^n$  and the Brownian motion  $\{w_t\}_{0 \leq t \leq 1}$  on the same probability space so that for any  $\delta > 0$  and any  $\eta > 0$

$$(9) \quad \lim_{n \rightarrow \infty} \mathbf{P} \left( \max_{k=1, \dots, n} \left| w_{k/n} - \frac{\sqrt{2}}{\sqrt{n}} S_k \right| > \delta n^{\eta-1/2} \right) = 0.$$

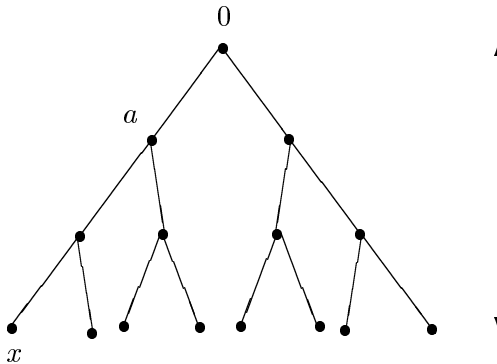
(The original one-dimension KMT approximation has been extended to the multivariate case by Einmahl [6]).

Note that the Brownian motion between two successful points  $x$  and  $y$  before reaching the boundary may again be modeled by a tree structure, and that the same technique as for trees works once more (with many technical issues).

*Application of the KMT approximation theorem.* The proof follows by considering the lattice points inside the circle  $\{z : |\sqrt{2}z - y| < \sqrt{n}(1 + 2\delta)\epsilon_n\}$  whose number is less than

$$\frac{\pi}{2} n(1 + 2\delta)^3 \epsilon_n^2.$$

#### 4. Covering Time of the Torus



First, we once again consider the “toy” problem of the covering time of the binary tree  $B_m$ . Let  $X = (X_n)$  be the first neighbor random walk starting from the left son  $a$  of the root, and consider hits to  $x$ , the leftmost leaf.  $\mathbf{P}^x$  again refers to walks starting at point  $x$ .

**4.1. Upper bound.** From Section 3.2 we get

$$\mathbf{P}^a(X \text{ hits } x \text{ before } 0) = 1 - H_1(0) = \frac{1}{m}.$$

This implies that

$$\mathbf{P}^0(X \text{ does not cover } x \text{ during first } N \text{ visits to } 0) \simeq \left(1 - \frac{1}{2m}\right)^N.$$

Let  $\Pi^0$  be the probability that the random walk starting at zero does not cover the binary tree  $B_m$  during  $N$  visits to 0. We have

$$\Pi^0 \leq 2^m \left(1 - \frac{1}{2m}\right)^N \quad \text{so that} \quad \Pi^0 \rightarrow 0 \quad \text{for} \quad N = 2(1 + \delta)m^2 \log 2.$$

The time needed for  $N$  visits to the root is  $2^{m+1}N$ ; this implies that

$$\mathbf{P}^0(X \text{ does not cover } B_m \text{ before time } 2(1 + \delta) \log 2 \times m^2 2^{m+1}) \rightarrow 0.$$

**4.2. Lower bound.** A ray  $x$  is called *successful* if the number of excursions from level  $i$  to level  $i+1$  in the ray is  $a(m-i)^2$ . Dembo *et al.* apply a second moment analysis to the successful rays to show that, with probability one, before  $2(1-\delta)m^2 \log 2$  visits to the root, there are points which are not covered. Then, the time needed to visit the root that many times is about  $2(1-\delta)m^2(\log 2)2^{m+1}$ .

To solve the standard random walk problem on  $\mathbb{Z}^2$ , Dembo *et al.* first solve the equivalent problem for the Brownian motion on the torus  $\mathbb{T}^2$ , where  $\mathbb{T}^2$  is identified with the set  $(-1/2, 1/2]^2$ .

Let  $\mathcal{T}(x, \epsilon)$  denote the time needed by the Brownian motion to enter the ball  $D(x, \epsilon)$ ,

$$\mathcal{T}(x, \epsilon) = \inf\{t > 0 \mid w_t \in D(x, \epsilon)\}, \quad \text{and} \quad C_\epsilon = \sup_{x \in \mathbb{T}^2} \mathcal{T}(x, \epsilon).$$

Therefore,  $C_\epsilon$  is the minimum time needed for the Brownian motion  $W_t$  to come within  $\epsilon$  of each point of  $\mathbb{T}^2$ . Equivalently,  $C_\epsilon$  is the amount of time needed for the Wiener sausage of radius  $\epsilon$  to completely cover  $\mathbb{T}^2$ . Dembo *et al.* [4] prove that

$$\lim_{\epsilon \rightarrow 0} \frac{C_\epsilon}{(\log \epsilon)^2} = \frac{2}{\pi}, \quad \text{almost surely.}$$

Using the KMP strong approximation theorem again provides the result for the standard random walk on  $\mathbb{T}^2$ .

### Bibliography

- [1] Aldous (D.) and Fill (J.). – Reversible Markov chains and random walks on graphs. – Monograph in preparation.
- [2] Aldous (David) and Pemantle (Robin) (editors). – *Random discrete structures*. – Springer-Verlag, New York, 1996, xviii+225p. Papers from the workshop held in Minneapolis, Minnesota, November 15–19, 1993.
- [3] Bingham (N. H.), Goldie (C. M.), and Teugels (J. L.). – *Regular variation*. – Cambridge University Press, Cambridge, 1987, xx+491p.
- [4] Dembo (Amir), Peres (Yuval), Rosen (Jay), and Zeitouni (Ofer). – Cover times for Brownian motion and random walks in two dimensions. – To appear.
- [5] Dembo (Amir), Peres (Yuval), Rosen (Jay), and Zeitouni (Ofer). – Thick points for planar Brownian motion and the Erdős-Taylor conjecture on random walk. *Acta Mathematica*, vol. 186, n° 2, 2001, pp. 239–270.
- [6] Einmahl (Uwe). – Extensions of results of Komlós, Major, and Tusnády to the multivariate case. *Journal of Multivariate Analysis*, vol. 28, n° 1, 1989, pp. 20–68.
- [7] Erdős (P.) and Taylor (S. J.). – Some problems concerning the structure of random walk paths. *Acta Mathematica Academiae Scientiarum Hungaricae*, vol. 11, 1960, pp. 137–162. (Unbound insert).
- [8] Feller (W.). – *An introduction to probability theory and its applications*. – Wiley international edition, 1970. Volume 1.
- [9] Komlós (J.), Major (P.), and Tusnády (G.). – An approximation of partial sums of independent rv's and the sample df. I. *Zeitschrift für Wahrscheinlichkeitstheorie und Verwandte Gebiete*, vol. 32, 1975, pp. 111–131.

## Introduction to Random Walks on Groups

Yves Guivarc’h

IRMAR, Université Rennes 1 (France)

March 5, 2001

Summary by Philippe Robert

### Abstract

In this talk simple examples are presented to illustrate some aspects of random walks on groups from the point of view of probability theory, statistical physics, ergodic theory, harmonic analysis, and group theory.

### 1. Shuffling Cards

A deck of cards is described by  $J = (a_1, \dots, a_r)$ , where  $a_i$  indicates the position of the  $i$ th card in the deck. The cards are shuffled so that the state of the deck of cards is  $(\sigma(a_1), \dots, \sigma(a_r))$ , where  $\sigma \in \Sigma$  is some permutation on  $J$ . Another shuffle would give the deck  $(\tau(\sigma(a_1)), \dots, \tau(\sigma(a_r)))$ , and so on. Of course, the permutation is likely to be different from one shuffle to another, but the habits of a given player will be such that he will choose at random among a given set  $A$  of permutations. For  $\alpha \in A$ , the permutation  $\alpha$  is chosen with probability  $p(\alpha) > 0$ . After a shuffle, the next permutation is chosen independently of the past. The position of the  $i$ th card is  $j$  after the first shuffle with probability

$$\sum_{\alpha \in A: \alpha(a_i)=j} p(\alpha),$$

after two shuffles the probability will be

$$\sum_{(\alpha, \beta) \in A: \beta(\alpha(a_i))=j} p(\alpha)p(\beta).$$

If  $p^n$  denotes the  $n$ th convolution of  $p$ ,

$$p^n(\sigma) = \sum_{\alpha_i \in A: \alpha_n \circ \alpha_{n-1} \circ \dots \circ \alpha_1 = \sigma} p(\alpha_n)p(\alpha_{n-1}) \cdots p(\alpha_1),$$

the distribution of the position of the  $i$ th card after the  $n$ th shuffle is given by

$$\mu_n^i = \sum_{\sigma \in \Sigma} p^n(\sigma) \delta_{\sigma(a_i)},$$

where  $\delta_x$  is the Kronecker symbol at  $x$ :  $\delta_x(x) = 1$  and  $\delta_x(y) = 0$  when  $y \neq x$ . A natural question in this setting is: provided that the set  $A$  is rich enough, is the position of the card  $a_j$  uniformly distributed on  $\{1, \dots, r\}$  when  $n$  gets large?

The distribution  $\mu_n$  on  $\Sigma$  of the configuration of the deck of cards after  $n$  shuffles is given by

$$\mu_n = \sum_{\sigma \in \Sigma} p^n(\sigma) \delta_\sigma,$$

with this notation  $\mu_n^i(j) = \mu_n(\sigma : \sigma(i) = j)$ . Does the distribution  $\mu_n$  on  $\Sigma$  converges to the uniform distribution on the group of permutations as  $n$  gets large? The answer to both questions is positive if the probability  $p$  satisfies some assumptions. It can then be shown that the convergence to the uniform distribution is exponentially fast with  $n$  (see Diaconis [2]).

This simple problem gives an illustration of the ergodic principle introduced in statistical physics after the work of Boltzmann and Gibbs:

- the limit is independent of the initial state;
- the limit is independent of the particular choice of the probability  $p$ ;
- the limit is the most disordered distribution  $m$  on  $\Sigma$ , i.e., the distribution with the maximal entropy  $H(m)$ , with

$$H(m) = \sum_{\sigma \in \Sigma} -m(\sigma) \log(m(\sigma)).$$

## 2. Random Walks in $\mathbb{Z}^d$

This random walk is defined as follows: starting from  $x \in \mathbb{Z}^d$ , it jumps to  $x \pm e_i$  with probability  $1/2d$ , where  $e_i$  is the  $i$ th unit vector. If  $S_n$  denotes the position after  $n$  steps it is well known that when  $d \leq 2$ , the sequence  $(S_n)$  almost surely visits the origin infinitely often; the random walk is then said to be recurrent. When  $d \geq 3$  the random walks visits 0 only a finite number of times; the random walk is transient. These results can be expressed in terms of electrical networks: each edge of  $\mathbb{Z}^d$  is assumed to have resistance 1,  $R_d$  is the effective resistance of  $\mathbb{Z}^d$  when the potential at 0 is 1 and 0 at infinity. It turns out that for  $d \leq 2$ ,  $R_d$  is infinite and  $R_d$  is finite when  $d \geq 3$ .

The Laplacian  $\Delta$  of the random walk is given by

$$\Delta(f)(x) = \frac{1}{2d} \left( \sum_{i=1}^d (f(x + e_i) + f(x - e_i)) \right) - f(x),$$

where  $f$  is some function on  $\mathbb{Z}^d$ . The potential function  $v(x)$  for the electrical network should satisfy  $\Delta(v) = 0$  with  $v(0) = 1$  and  $\lim_{x \rightarrow +\infty} v(x) = 0$ .

## 3. Polymer Dynamics in the Plane

A simplified model of a polymer in the plane is given by a broken line  $A_0 A_1 \dots A_n$  where each segment  $A_i A_{i+1}$  has length 1 and the angle between  $A_{i-1} A_i$  and  $A_i A_{i+1}$  is  $\pm \alpha \in [0, 2\pi)$  with probability  $1/2$ . If  $A_0 = (0, 0)$  and  $A_1 = (1, 0)$ , the vector  $Z_n = A_0 A_n$  can be represented in the complex plane as

$$Z_n = 1 + \sum_{k=1}^n e^{i\alpha S_k},$$

where  $S_k = \epsilon_1 + \dots + \epsilon_k$  and the  $\epsilon_i$  are independent Bernoulli random variables with  $\mathbf{P}(\epsilon_i = 1) = \mathbf{P}(\epsilon_i = -1) = 1/2$ ;  $(S_n)$  is the simple random walk on  $\mathbb{Z}$ . The average quadratic length of the polymer with  $N$  segments is given by

$$l_n = \sqrt{\mathbf{E}(Z_n^2)}.$$

It has been shown by Eyring that  $l_n/\sqrt{n}$  converges to a constant as  $n$  tends to infinity. The average length is conjectured to grow like  $n^\xi$  with  $\xi > 1/2$ .



#### 4. Random Rotations on the Sphere

This problem has been considered by Arnold and Krylov [1]. The action of two rotations  $a$  and  $b$  of  $\mathbb{R}^3$  on the unit sphere  $S^2$  centered at 0 is analyzed. If  $\lambda(a, b)$  is a product of  $n$  such rotations, one writes  $|\lambda| = n$ . For  $p \in S^2$ , the distribution  $\mu_n$  of  $\lambda(a, b)(p)$  is given by

$$\mu_n = \frac{1}{2^n} \sum_{|\lambda|=n} \delta_{\lambda(a,b)(p)}.$$

The problem is to determine when  $\mu_n$  converges to the uniform distribution on  $S^2$  and, if it occurs, the rate of convergence. The answer to the first point is positive under mild assumptions. The question concerning the speed is, for the moment, unsolved. This example is in some sense, a continuous analogue of the example of card shuffling.

#### 5. Random Walks on the Free Group

The free group with two generators  $a$  and  $b$  is denoted by  $\Gamma$ . An element  $\gamma$  is a string of letters  $a, a^{-1}, b$  and  $b^{-1}$  where a letter cannot be the inverse of the previous letter or the next letter in the string (otherwise the two letters cancel). The distance  $d(\gamma, \gamma')$  is given by the length of the string  $\gamma^{-1}\gamma'$ . The group  $\Gamma$  can be compactified by adding the set  $\partial\Gamma$  of infinite strings. If  $\xi$  is such a string and  $\gamma \in \Gamma$ , it is easily seen that, if  $(x_n)$  is a sequence of  $\Gamma$  and  $e$  is the empty string (the neutral element of the group), the quantity

$$\beta(\gamma, \xi) = \lim_{x_n \rightarrow \xi} (d(\gamma, x_n) - d(e, x_n))$$

is well defined.

The random walk considered here just adds  $a, a^{-1}, b$  or  $b^{-1}$  at the end of the string, with the convention that the inverse of the last letter suppresses this letter. This random walk is equivalent to a random walk on a homogeneous tree with degree 4. In particular it is transient and the length of the string almost surely converges to infinity. The Laplacian  $\Delta$  of this random walk is given by

$$\Delta(f)(\gamma) = \frac{1}{4}(f(\gamma a) + f(\gamma a^{-1}) + f(\gamma b) + f(\gamma b^{-1})) - f(\gamma),$$

for  $\gamma \in \Gamma$  and  $f$  a function on  $\Gamma$ . For  $\xi \in \partial\Gamma$ ,  $h_\xi(\gamma) = (1/3)^{\beta(\gamma, \xi)}$  is harmonic with respect to this Laplacian, i.e.,  $\Delta(h_\xi) = 0$ . Dynkin and Maljutov [4] have shown that every positive harmonic function  $f$  can be expressed as an integral of the elementary functions  $h_\xi$ ,  $\xi \in \partial\Gamma$ , i.e.,

$$f(\gamma) = \int_{\partial\Gamma} h_\xi(\gamma) \nu(d\xi),$$

where  $\nu$  is a positive measure on  $\partial\Gamma$ .

This situation has to be compared with the case of the random walks on  $\mathbb{Z}^d$  with  $d \geq 3$  which are also transient but without non-constant positive harmonic functions. Similarly, in a continuous setting, there does not exist any non-constant positive harmonic function  $f$  on  $\mathbb{R}^d$ , i.e., such that

$$\sum_{i=1}^d \frac{\partial^2 f}{\partial x_i^2} = 0.$$

But restricted to the unit disc of  $\mathbb{R}^2$ , such functions exist and can be represented as

$$\frac{1}{2\pi} \int_0^{2\pi} P(z, \theta) \nu(d\theta),$$

where  $\nu$  is some finite measure on  $[0, 2\pi)$  and  $P$  is the Poisson kernel

$$P(z, \theta) = \frac{1 - |z|^2}{|e^{i\theta} - z|^2}.$$

One can check that  $z \mapsto P(z, \theta)$  is harmonic: it is the equivalent of the function  $h_\xi$  for the unit disc.

### Bibliography

- [1] Arnol'd (V. I.), Kozlov (V. V.), and Neishtadt (A. I.). – *Dynamical systems. III.* – Springer-Verlag, Berlin, 1988, xiv+291p. Mathematical Aspects of Classical and Celestial Mechanics. Translated from the Russian.
- [2] Diaconis (Persi). – *Group representations in probability and statistics.* – Institute of Mathematical Statistics, Hayward, CA, 1988, vi+198p.
- [3] Doyle (Peter G.) and Snell (J. Laurie). – *Random walks and electric networks.* – Mathematical Association of America, Washington, DC, 1984, xiv+159p.
- [4] Dynkin (E. B.) and Maljutov (M. B.). – Random walk on groups with a finite number of generators. *Doklady Akademii Nauk SSSR*, vol. 137, 1961, pp. 1042–1045.
- [5] Guivarc'h (Yves). – Marches aléatoires sur les groupes. In *Development of mathematics 1950–2000*, pp. 577–608. – Birkhäuser, Basel, 2000.
- [6] Woess (Wolfgang). – *Random walks on infinite graphs and groups.* – Cambridge University Press, Cambridge, 2000, xii+334p.

## Random Matrices and Queues in Series

Yuliy Baryshnikov

LAMA, Université de Versailles – Saint-Quentin-en-Yvelines (France)

December 11, 2000

*Summary by Marianne Durand*

### Abstract

Queues in series are defined as an infinite sequence of clients queuing in front of an infinite sequence of servers where each time a client is served by a server, it immediately enters the next queue. Simple questions about this model are very hard to solve directly. This talk describes the centralized and normalized law of the departure of the  $k$ th client from the  $n$ th server, as  $n$  tends to infinity while  $k$  remains bounded; this law is related to a sequence of largest eigenvalues of random matrices. This relation allows us to use the numerous asymptotic results known regarding the spectra of random matrices and gain useful informations about the queuing processes.

### 1. Queues and Brownian Motions

Consider an infinite series of queues corresponding to servers, and an infinity of jobs. At first all the jobs are in the first queue; then when a job leaves the server  $Q_i$ , it immediately enters the queue corresponding to the server  $Q_{i+1}$ . The question asked is: When does the  $i$ th job leave the  $j$ th server? This can be modeled by pathweights in an infinite matrix. Let  $w_{k,l}$  denote the time needed to process the  $k$ th job on the  $l$ th server. The cost of the maximal weight of a path from  $(0,0)$  to  $(i,j)$  in the matrix  $(w_{k,l})$  is noted  $c(i,j)$ . The path is made of steps of size one where only one component increases. Then one observes that  $c(i,j)$  is equal to the time when the  $i$ th job leaves the  $j$ th server. This equality illustrates the fact that server  $j$  can process job  $i$  if it has already processed job  $i-1$  and if job  $i$  has left queue  $j-1$ .

The problem of queues in series can thus be modeled by an infinite matrix, where we assume from now on that the entries are independent identically distributed random variables, with finite variance. For the main theorem and for Section 3 the distribution is assumed to be geometric with parameter  $q$ . The aim of the talk [1] is to link the queue problem to the distribution of the largest eigenvalues of random Hermitian matrix with appropriate distribution. An infinite matrix of weights is also a model for a physical problem, the interacting particle process, see [7]; there we assume that all the integers corresponds to sites that are capable of containing one particle, and that at first all the sites with negative positions are full. In this model the weight  $w_{i,j}$  is the time taken by a particle to move from  $i$  to  $i+j$ .

A preliminary remark links this queue problem to Brownian motion [3]. Given  $(B_k)_{k=1,2,\dots}$  independent standard Brownian motions, and  $D_k^{(n)} = \frac{c(k,n)-en}{\sqrt{vn}}$ , where  $e$  is the expectation of  $w_{1,1}$  and  $v$  its variance, the following theorem holds:

**Theorem 1.** *The processes  $D^{(n)} = \left(D_k^{(n)}\right)_{k=1,2,\dots}$  converge in law as  $n \rightarrow \infty$  to the stochastic process  $D = (D_k)_{k=1,2,\dots}$  where  $D_k = \sup_{0=t_0 < t_1 < \dots < t_k=1} \sum_{i=0}^{k-1} (B_i(t_{i+1}) - B_i(t_i))$ .*

This can easily be seen by modeling a path from  $(0, 0)$  to  $(k, N)$  for large  $N$  by  $k$  long vertical lines, where the path uses the  $i$ th vertical line from the time  $t_i$  to the time  $t_{i+1}$ .

We now introduce the Gaussian Unitary Ensemble (GUE) [8] as the probability distribution on the Hermitian matrices with the density  $r_{\text{GUE}}(H) = Z^{-1} e^{-\text{tr} H^2/2}$  where  $Z$  is a normalizing constant equal to  $\int e^{-\text{tr} H^2/2} dH$ . A useful property is that a Hermitian matrix  $H$  is drawn from GUE if  $\Re(h_{ij})$  and  $\Im(h_{ij})$  are i.i.d. Gaussian random variables with mean 0 and variance 1. Given a matrix  $H$ , let  $H_k = (h_{(i,j)})_{1 \leq i, j \leq k}$  be the main minor of size  $k$  of  $H$  and  $\sigma_k$  the largest eigenvalue of  $H_k$ .

**Theorem 2.** *The laws of both processes  $\sigma = \{\sigma_k\}$  for  $H$  drawn from GUE, and  $D = \{D_k\}$  coincide.*

This theorem is proven in the next sections. We first exhibit a bijection between a finite restriction of size  $M$  of the queue problem and a subspace of  $\mathbb{N}^{M(M+1)/2}$  via Young tableaux. The second part of the proof is to relate this subspace of  $\mathbb{N}^{M(M+1)/2}$  to the dominant eigenvalues of minors of the matrix  $H$ .

## 2. Combinatorics

The bijection between the matrix  $H$  and Young tableaux is a generalization of the Robinson–Schensted–Knuth correspondence (see [5]) between Young tableaux and permutations. The matrix  $W$  of size  $N \times M$  with coefficients the weights  $w_{i,j}$  can be represented as a generalized permutation  $\alpha$ ,

$$\alpha = \begin{pmatrix} i_1 & i_2 & \dots & i_k \\ j_1 & j_2 & \dots & j_k \end{pmatrix},$$

where  $i_l \in \mathbb{N}_N$ , the integers between 1 and  $N$ ,  $j_l \in \mathbb{N}_M$  and  $j_l$  represents  $\alpha(i_l)$ . The integers  $i_l$  are not necessarily distinct, this is why the permutation is said to be generalized. The number of columns of type  $\binom{i}{j}$  is equal to  $w_{i,j}$ . As the generalized permutation  $\alpha$  is written in a lexicographically sorted fashion, the bijection is quite obvious. Indeed, given a matrix, the set of columns is well defined, and sorting gives the uniqueness of the image; conversely given a generalized permutation, one simply has to count the numbers of columns of type  $\binom{i}{j}$  to reconstruct the matrix. Recall that a *Young diagram*  $\lambda$  is a decreasing sequence  $(\lambda_1, \lambda_2, \dots, \lambda_r)$  that can be represented as  $r$  rows of boxes of heights  $\lambda_i$ . A *semi-standard Young tableau* is a filling of the boxes by positive integers such that the filling is increasing rightwards in rows and strictly increasing in columns. The Young diagram underlying a Young tableau  $P$  is called the shape and is denoted by  $\text{sh}(P)$ . The Robinson–Schensted–Knuth (RSK) correspondence is a bijection between the set of generalized permutations with  $k$  columns and the set of pairs of semi-standard Young tableaux  $(P, Q)$  having the same shape  $\lambda$  such that  $\sum_i \lambda_i = k$ . With the previous bijection between matrices and generalized permutations, we thus have a bijection between matrices of size  $N \times M$  and pairs of semi-standard Young tableaux  $(P, Q)$  with shape  $\lambda$ , such that  $\sum_i \lambda_i = k$ . We denote  $(P(w), Q(w))$  Young tableaux obtained by the RSK correspondence from a matrix  $w$ .

Some properties make this correspondence really interesting. We denote by  $W_{N,M,k}$  the set of matrices of size  $N \times M$  whose coefficients add up to  $k$ , to state a result on the way the distribution of the Young tableaux behaves through this correspondence.

**Lemma 1.** *If the set  $W_{N,M,k}$  is given the uniform distribution, then the distribution of  $P(w)$  for  $w \in W_{N,M,k}$  given  $\text{sh}(P(w)) = \lambda$  is uniform on the semi-standard Young tableaux of shape  $\lambda$ .*

Moreover the shape  $\text{sh}(P)$  encodes a few characteristics of  $w$ : the length  $\lambda_1$  of the first row is the maximal weight of the monotonous paths from  $(0, 0)$  to  $(M, N)$  in the table  $w$ , because it is the length of the longest increasing subsequence of the second line of  $\alpha$ . This is a direct consequence of the RSK algorithm. It is then possible to consider the Young tableau  $P$  as an embedding of several Young tableaux  $P_1, \dots, P_M$  where  $P_M$  equals  $P$  and  $P_i$  is obtained from  $P_{i+1}$  by removing all boxes filled with  $i + 1$ . A nice consequence of the way the tableau  $P$  is built is that the sequence of the lengths of the first rows of the embedded Young tableaux coincides with the sequence of maximal paths from  $(0, 0)$  to  $(i, N)$  as  $i$  goes from 1 to  $M$ . Thus we now focus on Young tableaux instead of the weight matrix. We introduce a representation of the Young tableaux that encodes the shape of the tableaux, and also the shape of all the embedded tableaux inside. To describe a Young tableau  $P$  filled with  $\mathbb{N}_M$ , let  $x_j^i$  be the coordinate of the rightmost box filled with a number at most  $i$  in the  $j$ th row of the tableau. Equivalently, this is just the length of the  $j$ th row in the tableau  $P_{i-1}$  defined by the embedding above. The elements  $x_j^i$ ,  $1 \leq i \leq M$ ,  $1 \leq j \leq i$  can be seen as a triangular array of size  $\frac{M(M+1)}{2}$ . The image  $x$  of a tableau  $P$  by this transformation has the property that its last line is equal to  $\lambda = (\lambda_1, \dots, \lambda_M)$ , and that its first column is equal to  $c(k, N)$ ,  $k \in \mathbb{N}_M$ , corresponding to the length of maximal paths from  $(0, 0)$  to  $(k, N)$ . This correspondence is formalized in the following lemma.

**Lemma 2.** *Let the Gelfand–Cetlin cone  $C_{GC}$  be the set of triangular arrays  $(x_j^i)$  of size  $\frac{M(M+1)}{2}$  such that  $x_{j-1}^i \geq x_{j-1}^{i-1} \geq x_j^i$  for  $1 \leq i \leq M$ ,  $1 \leq j \leq i$ . Then the Young tableaux filled with  $\mathbb{N}_M$  are in one-to-one correspondence with the integer points in the Gelfand–Cetlin cone.*

What we have now is a mapping from  $W_{N,M,k}$ , the set of matrices of size  $N \times M$  whose coefficients add up to  $k$ , to the set of integers point in the Gelfand–Cetlin cone. This mapping has the property that if  $W_{N,M,k}$  is given the uniform distribution, then the distribution of  $x$ , given that the last line is equal to  $\lambda = (\lambda_1, \dots, \lambda_M)$ , is uniform on the integers points of the Gelfand–Cetlin cone.

### 3. Gaussian Unitary Ensemble

From now on we add the restriction that the distributions of the coefficients  $w(i, j)$  are i.i.d. geometric with parameter  $q$ , that is  $\mathbf{P}(w_{i,j} = k) = (1 - q)q^k$ . All the results of the previous section still hold in this context, as they were obtained in full generality. The aim of this section is to finish the proof of Theorem 2. For this we describe the distribution of  $x(w)$  by its distribution conditioned upon values of its last line, and the distribution of its last line. This is then linked to the distribution of the limit processes  $D_k$ .

The probability that the RSK correspondence applied to a random matrix  $w$  with i.i.d. geometric entries with parameter  $q$  yields a pair of Young tableaux of shape  $\lambda = (\lambda_1, \dots, \lambda_M)$  is

$$(1) \quad \frac{(1 - q)^{MN}}{M!} \prod_{j=0}^{M-1} \frac{1}{j!(N - M + j)!} \prod_{1 \leq i < j \leq M} (\lambda_i - \lambda_j - i + j)^2 \prod_{i=1}^M \frac{(\lambda_i + N - i)!}{(\lambda_i + M - i)!} q^k$$

where  $k = \sum \lambda_i$ . The proof of this formula can be found in [4] and is based on the fact that there are  $\prod_{1 \leq i < j \leq M} \frac{\lambda_i - \lambda_j - i + j}{j - i}$  semi-standard tableaux of shape  $\lambda$  filled with  $\mathbb{N}_N$ . The vector of centered and normalized variables  $\xi_i = \frac{\lambda_i - eN}{\sqrt{vN}}$  (with  $e$  the average and  $v$  the variance of  $w_{1,1}$ ), is noted  $\xi$ . Plugging the Stirling approximation in Equation (1) yields that for fixed  $q$  such that  $0 < q < 1$ , fixed  $M$  and  $N \rightarrow \infty$ , the distribution of  $\xi$  converges weakly to  $Z^{-1} \prod_{i < j} (\xi_i - \xi_j)^2 \prod_i e^{-\xi_i^2/2}$  ( $Z$  is a normalizing constant). This is the distribution of the vector of ordered eigenvalues of a random matrix drawn from GUE. This property leads to the following theorem:

**Theorem 3.** *The distribution of the sequence  $(D_1, \dots, D_M)$  defined in Theorem 1 is the distribution of the first column of the random triangular array  $x$  of size  $M(M+1)/2$  distributed uniformly for a fixed last line, and the distribution of its last line is the distribution of the eigenvalues of matrices drawn from GUE [6].*

The distribution of  $(D_1, \dots, D_M)$  is the same as the distribution of the first column of  $x(w)$ , up to a proper normalization. And if the distribution on  $w$  is uniform, then the distribution on  $x$ , knowing its last line, is uniform. As the probability of getting a Young tableau of shape  $\lambda$ , and thus an array  $x$  of last line  $\lambda$ , is the same as getting the vector of ordered eigenvalues of a random matrix drawn from GUE equal to  $\lambda$ , Theorem 3 is proved.

The Gelfand–Cetlin polyhedron  $GC(\lambda)$  is defined as a subset of  $C_{GC}$  such that the last line of the array is equal to  $\lambda = (\lambda_1, \dots, \lambda_M)$ . Theorem 3 means that the distribution of  $(D_1, \dots, D_M)$  is uniform on  $GC(\lambda)$ .

This allows us to state the theorem below, which is a major step in the proof of Theorem 2.

**Theorem 4.** *Let  $H = (h_{ij})$ ,  $i, j \leq M$  be a random matrix drawn from GUE with eigenvalues  $(\lambda_1, \dots, \lambda_M)$ , and*

$$x(H) = \begin{pmatrix} x_1^1 & & & & \\ \cdots & \cdots & & & \\ x_1^{M-1} & \cdots & x_{M-1}^{M-1} & & \\ x_1^M & x_2^M & \cdots & x_M^M & \end{pmatrix}$$

where  $x_j^i$  is the  $j$ th eigenvalue of the main minor of size  $i$  of  $H$ . Then the triangular array  $x(H)$  is uniformly distributed in the polyhedron  $GC(\lambda)$ .

The proof of this theorem is based on the fact that the last line of the array  $x$  corresponds to the eigenvalues of the matrix  $H$ , which is drawn from GUE, and that given this last line, the distribution is uniform. The last line of  $x$  is equal to  $\lambda$  and its first column to the vector  $(\sigma_k)$  by definition. This together with Theorems 3 and 4 proves Theorem 2.

Theorem 2 can be used to prove the conjecture of Peter Glynn and Ward Witt [3] stating that  $D_k/\sqrt{k} \rightarrow 2$ . The proof uses the already known fact that the largest eigenvalue of random Hermitian matrix drawn from GUE rescaled by  $\sqrt{k}$  converges in distribution to 2, see [2].

### Bibliography

- [1] Baryshnikov (Yu.). – GUEs and queues. *Probability Theory and Related Fields*, vol. 119, n° 2, 2001, pp. 256–274.
- [2] Geman (Stuart). – A limit theorem for the norm of random matrices. *Annals of Probability*, vol. 8, n° 2, 1980, pp. 252–261.
- [3] Glynn (Peter W.) and Whitt (Ward). – Departures from many queues in series. *Annals of Applied Probability*, vol. 1, n° 4, 1991, pp. 546–572.
- [4] Johansson (Kurt). – Shape fluctuations and random matrices. *Communications in Mathematical Physics*, vol. 209, n° 2, 2000, pp. 437–476.
- [5] Knuth (Donald E.). – Permutations, matrices, and generalized Young tableaux. *Pacific Journal of Mathematics*, vol. 34, 1970, pp. 709–727.
- [6] Kuperberg (Greg). – Random words, quantum statistics, central limits, random matrices, September 2000. To appear in *Methods Appl. Anal.*. Available from <http://front.math.ucdavis.edu/math.PR/9909104>.
- [7] Liggett (Thomas M.). – *Interacting particle systems*. – Springer-Verlag, New York, 1985, xv+488p.
- [8] Mehta (M. L.). – *Random matrices and the statistical theory of energy levels*. – Academic Press, New York, 1967, x+259p.

# Information Theory by Analytic Methods: The Precise Minimax Redundancy

Wojciech Szpankowski

Department of Computer Science, Purdue University (USA)

March 5, 2001

Summary by Thomas Klausner

## 1. Introduction

The redundancy-rate problem of universal coding is concerned with determining by how much the actual code length (representation of a word in a code) exceeds the optimal code length. Revisiting the theme of his last year's seminar talk [1], Szpankowski went into more detail explaining different models for redundancy, and introduced the *generalized Shannon code* in order to solve the minimax redundancy problem for a single memoryless source.

A code is defined as follows:

**Definition 1.** A code  $C_n$  is a mapping from the set  $\mathcal{A}^n$  of all sequences of length  $n$  over the alphabet  $\mathcal{A}$  to the set  $\{0, 1\}^*$  of binary sequences.

Most of the time we use source models which specify probabilities for specific messages. For these,  $\mathcal{P}(x_1^n)$  is the probability of the message  $x_1^n$ , the code length of a message  $x_1^n = x_1 \dots x_n$ , with  $x_i \in \mathcal{A}$ , in the code  $C_n$  will be denoted by  $L(C_n, x_1^n)$ , and  $H_n(\mathcal{P}) = -\sum_{x_1^n} \mathcal{P}(x_1^n) \log \mathcal{P}(x_1^n)$  is the entropy of the probability distribution, where  $\log$  is taken to base 2.

## 2. Basic Results

A *prefix code* or *instantaneous code* is a code in which no codeword is a prefix for another codeword; in other words, if you present the codewords as a binary trie, the valid codewords are only in the leaves (not in the internal nodes).

For prefix codes the following inequality holds:

**Lemma 1** (Kraft's inequality). *For any prefix code (over a binary alphabet), the codeword lengths  $l_1, l_2, \dots, l_m$  satisfy the inequality*

$$\sum_{i=1}^m 2^{-l_i} \leq 1.$$

A related problem is to find out how many tuples  $l_1, \dots, l_m$  exist where equality holds. This has been tackled and solved by Flajolet and Prodinger [2]. Asymptotically, it grows as  $\alpha \phi^m$ , where  $\alpha \approx 0.254$  and  $\phi \approx 1.794$ .

Another important result is Shannon's classic lower bound on the average code length (see [3]):

**Lemma 2** (Shannon). *For any code, the average code length  $\mathbf{E}[L(C_n, X_1^n)]$  cannot be smaller than the entropy of the source  $H_n(\mathcal{P})$ :*

$$\mathbf{E}[L(C_n, X_1^n)] \geq H_n(\mathcal{P})$$

Trivially, one can see that there must exist at least one  $\tilde{x}_1^n$  with

$$L(\tilde{x}_1^n) \geq -\log \mathcal{P}(\tilde{x}_1^n).$$

A lemma by Barron deals with the individual lengths of the code words:

**Lemma 3** (Barron). *Let  $L(X_1^n)$  be the length of a codeword in a code satisfying Kraft's inequality, where  $X_1^n$  is generated by a stationary ergodic source. For any sequence of positive constants  $a_n$  satisfying  $\sum 2^{-a_n} < \infty$ , the following holds:*

$$\mathbf{P}\{L(X_1^n) \leq -\log \mathcal{P}(X_1^n) - a_n\} \leq 2^{-a_n}.$$

From this we immediately get

$$L(X_1^n) \geq -\log \mathcal{P}(X_1^n) - a_n \quad (\text{almost surely}).$$

### 3. Redundancy

Redundancy measures the distance to the optimal code state, reaching the lower bound given by the entropy. Since there are different ways to define the “worst case,” we define three types of redundancy: pointwise  $R_n(C_n, \mathcal{P}; x_1^n)$ , average  $\bar{R}_n(C_n, \mathcal{P})$  and maximal  $R^*(C_n, \mathcal{P})$ :

$$\begin{aligned} R_n(C_n, \mathcal{P}; x_1^n) &= L(C_n, x_1^n) + \log \mathcal{P}(x_1^n) && (\geq -a_n(a.s.)), \\ \bar{R}_n(C_n, \mathcal{P}) &= \mathbf{E}_{X_1^n} [R_n(C_n, \mathcal{P}; X_1^n)] \\ &= \mathbf{E}[L(C_n, X_1^n)] - H_n(\mathcal{P}), \\ R^*(C_n, \mathcal{P}) &= \max_{x_1^n} [R_n(C_n, \mathcal{P}; x_1^n)]. \end{aligned}$$

The redundancy-rate problem consists in finding the rate of growth of the corresponding minimax quantities

$$\begin{aligned} \bar{R}_n(\mathcal{S}) &= \min_{C_n} \sup_{\mathcal{P} \in \mathcal{S}} \mathbf{E}[R_n(C_n, \mathcal{P}; x_1^n)], \\ R_n^*(\mathcal{S}) &= \min_{C_n} \sup_{\mathcal{P} \in \mathcal{S}} \max_{x_1^n} [R_n(C_n, \mathcal{P}; x_1^n)], \end{aligned}$$

as  $n \rightarrow \infty$  for a class  $\mathcal{S}$  of source models.

There are also other measures of optimality, e.g. for coding, gambling, or predictions. For these, the following functions, called minimax regret functions, are used:

$$\begin{aligned} \bar{r}_n &= \min_{C_n} \sup_{\mathcal{P} \in \mathcal{S}} \sum_{x_1^n} \mathcal{P}(x_1^n) [L_i + \log \sup_{\mathcal{P}} \mathcal{P}(x_1^n)], \\ r_n^* &= \min_{C_n} \max_{x_1^n} [L_i + \log \sup_{\mathcal{P}} \mathcal{P}(x_1^n)]. \end{aligned}$$

Note that  $r_n^* = R_n^*$ . Sometimes, the maximin regret is of interest:

$$\tilde{r}_n = \sup_{\mathcal{P} \in \mathcal{S}} \min_{C_n} \sum_{x_1^n} \mathcal{P}(x_1^n) [L_i + \log \sup_{\mathcal{P}} \mathcal{P}(x_1^n)].$$

These functions are sometimes called the average minimax regret ( $\bar{r}_n$ ), the maximal minimax regret ( $r_n^*$ ), and the average maximin regret ( $\tilde{r}_n$ ). One can interpret these functions as target functions for the game theoretical problem of choosing  $L$  so that for all  $x_1^n$ , the value of the function gets as good as possible, that is,  $-\log \sup \mathcal{P}(x_1^n)$ .

In the following, we will only look at the redundancy functions.



#### 4. Precise Maximal Redundancy

In 1978, Shtarkov proved the following bounds for the minimax redundancy:

$$\log \left( \sum_{x_1^n} \sup_{\mathcal{P} \in \mathcal{S}} \mathcal{P}(x_1^n) \right) \leq R_n^*(\mathcal{S}) \leq \log \left( \sum_{x_1^n} \sup_{\mathcal{P} \in \mathcal{S}} \mathcal{P}(x_1^n) \right) + 1.$$

We want to find a precise result for  $R_n^*(\mathcal{S})$ . We start with the easier problem of finding the optimal code for maximal redundancy for a known source  $\mathcal{P}$

$$R_n^*(\mathcal{P}) = \min_{C_n \in \mathcal{C}} R_n^*(C_n, \mathcal{P}).$$

We already know that for the average redundancy of one known source

$$\bar{R}_n(\mathcal{P}) = \min_{C_n \in \mathcal{C}} \mathbf{E}_{x_1^n} [R_n(C_n, \mathcal{P}; x_1^n)],$$

the Huffman code is optimal—indeed, it is designed so as to solve this optimization problem. For the maximal redundancy problem we introduce a new code, the *generalized Shannon code*.

In the ordinary *Shannon code*, the length of its symbol in the code for a given  $\mathcal{P}$  is  $\lceil 1/\mathcal{P}(x_1^n) \rceil$ . In the generalized Shannon code, on the other hand, we set the length to be  $\lfloor 1/\mathcal{P}(x_1^n) \rfloor$  for some symbols  $x_1^n \in \mathcal{L}$  and  $\lceil 1/\mathcal{P}(x_1^n) \rceil$  for the others in such a way that Kraft's inequality holds. For non-dyadic codes (dyadic ones fulfill  $R_n^*(\mathcal{P}) = 0$ ), we sort the probabilities  $\mathcal{P}(x_1^n)$ :

$$0 \leq \langle -\log p_1 \rangle \leq \langle -\log p_2 \rangle \leq \dots \leq \langle -\log p_{|\mathcal{A}|^n} \rangle \leq 1 \quad (\text{where } \langle x \rangle = x - \lfloor x \rfloor)$$

and choose  $j_0$  to be the maximal  $j$  such that Kraft's inequality still holds:

$$\sum_{i=0}^{j-1} p_i 2^{\langle -\log p_i \rangle} + \sum_{i=j}^{|\mathcal{A}|^n} p_i 2^{\langle -\log p_i \rangle - 1} \leq 1.$$

Then  $R_n^*(\mathcal{P}) = 1 - \langle -\log p_{j_0} \rangle$  and the generalized Shannon code with  $\mathcal{L} = \{1, \dots, j_0\}$  is optimal.

Now we generalize to systems of probability distributions  $\mathcal{S}$ . Let

$$Q^*(x_1^n) = \frac{\sup_{\mathcal{P} \in \mathcal{S}} \mathcal{P}(x_1^n)}{\sum_{y_1^n \in \mathcal{A}^n} \sup_{\mathcal{P} \in \mathcal{S}} \mathcal{P}(y_1^n)}.$$

Then

$$R_n^*(\mathcal{S}) = R_n^*(Q^*) + \log \left( \sum_{x_1^n \in \mathcal{A}^n} \sup_{\mathcal{P} \in \mathcal{S}} \mathcal{P}(x_1^n) \right),$$

with

$$R_n^*(Q^*) = 1 - \langle -\log q_{j_0} \rangle$$

as above.

If we now take the generalized Shannon code that minimizes the maximal redundancy, we get for a sequence generated by a single memoryless source, for  $n \rightarrow \infty$ , and  $\alpha = \log \frac{1-p}{p}$  irrational:

$$R_n^*(\mathcal{P}_p) = -\frac{\log \log 2}{\log 2} + o(1) = 0.5287 + o(1).$$

## 5. Average Minimax Redundancy

In the simple case where  $\mathcal{S}$  consists of one distribution  $\mathcal{P}$ , the computation of  $\bar{R}_n^H$  is the Huffman problem:

$$\bar{R}_n^H(\mathcal{P}) = \min_{C_n \in \mathcal{C}} \sum_{x_1^n} \mathcal{P}(x_1^n) R_n(C_n, \mathcal{P}; x_1^n).$$

From known results (where we have  $\bar{R}_n^H \approx R_n^*$ ), we conjecture:

**Conjecture 1.** *Under certain additional conditions, we have, as  $n \rightarrow \infty$ ,*

$$\bar{R}_n = R_n^* + \Theta(1) = \log \left( \sum_{x_1^n \in \mathcal{A}^n} \sup_{\mathcal{P} \in \mathcal{S}} \mathcal{P}(x_1^n) \right) + \Theta(1).$$

## 6. Average Redundancy for Particular Codes

For single memoryless sources, we have explicit results for  $n \rightarrow \infty$  for some codes. In particular, we have for the Huffman code

$$\bar{R}_n = \begin{cases} \frac{3}{2} - \frac{1}{\ln 2} & \text{if } \alpha \text{ irrational,} \\ \frac{3}{2} - \frac{1}{M} (\langle Mn\beta \rangle - \frac{1}{2}) - (M(1 - 2^{-1/M}))^{-1} 2^{-\langle Mn\beta \rangle / M} & \text{if } \alpha = \frac{N}{M}, \end{cases}$$

for the Shannon code

$$\bar{R}_n = \begin{cases} \frac{1}{2} & \text{if } \alpha \text{ irrational,} \\ \frac{1}{2} - \frac{1}{M} (\langle Mn\beta \rangle - \frac{1}{2}) & \text{if } \alpha = \frac{N}{M}, \end{cases}$$

and for the generalized Shannon code

$$\bar{R}_n = \frac{3}{2} - 2 \ln 2 + o(1) \approx 0.113705639.$$

For more basics and in-depth knowledge regarding analytic information theory, the interested reader is referred to Szpankowski's book [4].

## Bibliography

- [1] Flajolet (Philippe). – Analytic information theory and the redundancy rate problem [summary of a talk by Wojciech Szpankowski]. In Chyzak (Frédéric) (editor), *Algorithms Seminar, 1999–2000*, pp. 133–136. – Institut National de Recherche en Informatique et en Automatique, November 2000. Research Report n° 4056.
- [2] Flajolet (Philippe) and Prodinger (Helmut). – Level number sequences for trees. *Discrete Mathematics*, vol. 65, n° 2, 1987, pp. 149–156.
- [3] Shannon (C. E.). – A mathematical theory of communication. *Bell System Technical Journal*, vol. 27, 1948, pp. 379–423 and 623–656.
- [4] Szpankowski (Wojciech). – *Average-case analysis of algorithms on sequences*. – John Wiley & Sons, Chichester, New York, March 2001, *Wiley-Interscience Series in Discrete Mathematics*.

## Part V

# Asymptotics and Analysis



## On Jackson's $q$ -Bessel Functions

Changgui Zhang

Université de La Rochelle (France)

October 30, 2000

Summary by Bruno Salvy

### Abstract

An analytic study of linear  $q$ -difference equations leads to a simple derivation of some connection formulae, generalizing the asymptotic expansion of the Bessel  $J_\nu$  functions.

### 1. Differential and $q$ -Difference Equations

Linear differential operators are polynomials in  $x$  and  $\partial_x = d/dx$ . These operators can be discretized using  $q$ -difference operators expressed in terms of  $q$ ,  $x$ , and  $\sigma_q$  where  $\sigma_q(f)(x) := f(qx)$ . When  $q \rightarrow 1$ ,  $(\sigma_q - 1)(f)(x)/(q - 1)$  tends to  $xf'(x)$ . This discretization is not unique. It gives rise to several generalizations of classical functions and identities relating them. C. Zhang's work is an *analytic* study of these operators, of the asymptotics of their solutions and the divergence of their series expansions.

A simple example of a  $q$ -difference equation is given by  $(x\sigma_q - 1)y(x) = 0$ . For  $|q| < 1$  and  $x \in \mathbb{C}^* := \mathbb{C} \setminus \{0\}$ , a solution of this equation is the Jacobi function

$$\theta_q(x) := \sum_{n \in \mathbb{Z}} q^{n(n-1)/2} x^n = (q; q)_\infty (-x; q)_\infty (-q/x; q)_\infty$$

where the last equality is *Jacobi's triple product identity*, using the notation

$$(a; q)_\infty = (1 - a)(1 - aq)(1 - aq^2) \cdots$$

The product form shows that  $\theta_q(x)$  is analytic in  $\mathbb{C}^*$ , and that its set of zeroes is  $-q^{\mathbb{Z}}$ .

Another important solution of the same equation is  $e_q(x) := q^{\log_q x (\log_q x - 1)/2}$ , equivalent to  $\theta_q(x)$  when  $x \rightarrow 0$ . In the asymptotic behaviour of solutions in the neighbourhood of irregular singular points, the function  $e_q$  plays the same role as the exponential in the differential case. Another simple equation is  $(\sigma_q - x)y(x) = 0$ , which has  $q^{-\log_q x (\log_q x - 1)/2}$  and  $1/\theta_q(x)$  as solutions. As opposed to the differential case, inverses of these analogues of the exponential are not obtained by changing  $x$  into  $-x$ .

A complete classification of the possible formal local behaviours of solutions of linear  $q$ -difference equations was obtained by Carmichael in 1912. For an equation of order  $m$  in  $\sigma_q$  with analytic coefficients at the origin, there exists a family of  $m$  formal solutions, each of which is of the form

$$(1) \quad y_j(x) = x^{r_j} e_q^{-k_j}(x) \sum_{\nu=0}^{m_j-1} (\log_q x)^\nu f_{j,\nu}(x), \quad j = 1, \dots, m,$$

where  $r_j \in \mathbb{C}$ ,  $k_j \in \mathbb{Q}$ ,  $m_j \in \mathbb{N}^*$ , and  $f_{j,\nu}(x) \in \mathbb{C}[[x^{1/d}]]$  for some  $d \in \mathbb{N}^*$ . Each of these can be computed from the equation.

## 2. Hypergeometric and $q$ -Hypergeometric Connection Formulae

The connection problem lies in expressing (the analytic continuation of) one of the above  $y_j$ 's that are defined at the origin as a linear combination in terms of a similar basis at another singular point. There is no general method to compute "closed forms" for these constants, except in special cases such as the hypergeometric case.

Hypergeometric series in the classical (differential) case are series  $F(x) = \sum_{n \geq 0} a(n)x^n$  such that  $a(n+1)/a(n) =: r(n) = P(n)/Q(n)$  is a fixed rational function in  $n$ . In terms of the shift operator  $S_n$  this means that the sequence  $a(n)$  cancels  $Q(n) - P(n)S_n^{-1}$  from which it follows that the generating series  $F$  cancels the linear differential operator  $Q(x\partial_x) - P(x\partial_x)x$ . Introducing the roots of  $P$  and  $Q$ , hypergeometric series are classically denoted

$${}_pF_q \left( \begin{matrix} a_1, \dots, a_p \\ b_1, \dots, b_q \end{matrix} \middle| x \right) := \sum_{n \geq 0} \frac{(a_1)_n \cdots (a_p)_n}{(b_1)_n \cdots (b_q)_n} \frac{x^n}{n!},$$

where  $(a)_n = a(a+1)\cdots(a+n-1)$ . This series is convergent for  $q > p$  and has only regular singularities if and only if  $p = q + 1$ .

The  $q$ -analogue of this function is known as the  ${}_r\phi_s$  *basic hypergeometric series*. In this case the ratio of two consecutive coefficients is a fixed rational function in  $q^n$ . The general form is

$${}_r\phi_s \left( \begin{matrix} a_1, \dots, a_r \\ b_1, \dots, b_s \end{matrix} ; q, x \right) := \sum_{n \geq 0} \frac{(a_1; q)_n \cdots (a_r; q)_n}{(b_1; q)_n \cdots (b_s; q)_n} \left( (-1)^n q^{n(n-1)/2} \right)^{s+1-r} x^n,$$

where  $(a; q)_n = (1-a)(1-aq)\cdots(1-aq^{n-1})$ .

A simple example is Heine's  ${}_2\phi_1(a, b; c; q, x)$ , which has Gauss's  ${}_2F_1(\alpha, \beta; \gamma; x)$  as a limit when  $a = q^\alpha$ ,  $b = q^\beta$ ,  $c = q^\gamma$ , and  $q \rightarrow 1$ . Heine's function satisfies a second-order  $q$ -difference equation. This equation has no irregular singularity (it is a *Fuchsian* equation). A general technique to relate solutions of such equations at 0 and infinity in the classical hypergeometric case is based on a Mellin–Barnes integral representation. This approach was extended to the  $q$ -difference case by Watson in 1910, who found that for  $ab \neq 0$ ,

$$(2) \quad {}_2\phi_1(a, b; c; q, x) = C_1(x) {}_2\phi_1(a, aq/c; aq/b; q, cq/abx) + C_2(x) {}_2\phi_1(b, bq/c; bq/a; q, cq/abx),$$

where

$$C_1(x) = \frac{(b, c/a; q)_\infty (ax, q/ax; q)_\infty}{(c, b/a; q)_\infty (x, q/x; q)_\infty}, \quad C_2(x) = \frac{(a, c/b; q)_\infty (bx, q/bx; q)_\infty}{(c, a/b; q)_\infty (x, q/x; q)_\infty}.$$

This method is presented in detail in Slater's book [4]. The connection "constants"  $C_1(x)$  and  $C_2(x)$  are annihilated by  $\sigma_q - 1$  and are uniform (they satisfy  $C_k(xe^{2i\pi}) = C_k(x)$ ). Thus they are *elliptic*, since when expressed in  $(u, \tau)$  defined by  $x = \exp(2i\pi u)$  and  $q = \exp(-2i\pi\tau)$  with  $\Im(\tau) > 0$  they are doubly periodic.

## 3. Jackson's $q$ -Bessel Functions

Bessel functions are classically defined as solutions of the Bessel equation

$$((x\partial_x - \nu)(x\partial_x + \nu) + x^2) y(x) = 0.$$

When  $\nu \notin \mathbb{Z}$ , a basis of solutions is given by the Bessel  $J_\nu(x)$  and  $J_{-\nu}(x)$  functions, which can be expressed in terms of the hypergeometric series by

$$J_{\pm\nu}(x) = \frac{(x/2)^{\pm\nu}}{\Gamma(\pm\nu + 1)} {}_2F_1(1, 1; \pm\nu + 1; -x^2/4).$$

The Bessel equation can be derived from the differential equation of the  ${}_2F_1$  by *confluence*: this is achieved by considering  ${}_2F_1(\nu + 1/2, \beta; 2\nu + 1; x/\beta)$  and letting  $\beta$  tend to infinity. In this process, the singularity at infinity becomes irregular.

Similarly, Jackson introduced in 1905 two  $q$ -analogues of the Bessel functions,

$$(3) \quad \begin{aligned} J_\nu^{(1)}(x; q) &= \frac{(q^{\nu+1}; q)_\infty}{(q; q)_\infty} \left(\frac{x}{2}\right)^\nu {}_2\phi_1\left(0, 0; q^{\nu+1}; q, -\frac{x^2}{4}\right), \\ J_\nu^{(2)}(x; q) &= \frac{(q^{\nu+1}; q)_\infty}{(q; q)_\infty} \left(\frac{x}{2}\right)^\nu {}_0\phi_1\left(; q^{\nu+1}; q, -\frac{x^2 q^{\nu+1}}{4}\right). \end{aligned}$$

The classical  $J_\nu$  function is recovered in two ways by letting  $q$  tend to 1 in  $J_\nu^{(k)}(x(1-q); q)$  for  $k \in \{1, 2\}$ . The radiuses of convergence of the basic hypergeometric series (in  $q$ ) given here are respectively finite for  $J_\nu^{(1)}$  (provided  $|x| < 2$ ) and infinite for  $J_\nu^{(2)}$ .

These functions are solutions of two  $q$ -difference equations of order 2 in  $\sigma_p$  with  $p = \sqrt{q}$  that are easily derived from (3). These equations can be seen as arising from the equation of the  ${}_2\phi_1$  by confluence, but it is not clear how to use this process in order to obtain a connection formula by a limiting process from (2). As in the classical case, both  $J_\nu^{(k)}$  and  $J_{-\nu}^{(k)}$  are independent solutions of their respective  $q$ -difference equation, for  $k = 1, 2$ . The equations have a regular singularity at the origin and an irregular singularity at infinity.

#### 4. Derivation of Connection Formulae

Connection formulae between the series expansions (3) and the (unique) basis of formal solutions at infinity of the form given by (1) generalize the classical asymptotic expansion

$$(4) \quad J_\nu(x) = \frac{e^{-i\frac{\pi}{4}(2\nu+1)}}{\sqrt{2\pi x}} e^{ix} {}_2F_0\left(-\nu + \frac{1}{2}, \nu + \frac{1}{2}; ; \frac{2i}{x}\right) + \frac{e^{i\frac{\pi}{4}(2\nu+1)}}{\sqrt{2\pi x}} e^{-ix} {}_2F_0\left(-\nu + \frac{1}{2}, \nu + \frac{1}{2}; ; -\frac{2i}{x}\right).$$

(A nice application of this formula is the derivation of an asymptotic expansion of the location of the zeroes of  $J_\nu(x)$ ; this generalizes to those of  $J_\nu^{(2)}$ .)

We start with  $J_\nu^{(1)}$  and its  $q$ -difference equation

$$(\sigma_p^2 - (p^\nu + p^{-\nu})\sigma_p + (1 + x^2/4)) y(x) = 0.$$

By changing  $x$  into  $1/t$  and  $y(x)$  into  $z(1/t)$ , the equation becomes

$$\left( \left(1 + \frac{1}{4p^4 t^2}\right) \sigma_p^2 - (p^\nu + p^{-\nu})\sigma_p + 1 \right) z(t) = 0.$$

The exponential part of the behaviour (see Eq. (1)) is sought in terms of  $\mathcal{E}_\alpha(t) = 1/\theta_p(-\alpha t)$ , which is cancelled by  $\sigma_p + \alpha t$ . The change of unknown function  $z(t) = \mathcal{E}_\alpha(t) f_\alpha(t)$  leads to

$$\left( \left(1 + \frac{1}{4p^4 t^2}\right) \alpha^2 p t^2 \sigma_p^2 - \alpha t (p^\nu + p^{-\nu})\sigma_p - 1 \right) f_\alpha(t) = 0.$$

Thus, by choosing  $\alpha$  such that  $\alpha^2 = 4p^3$ , one gets an equation for  $f_\alpha$  which has power series solutions. A further simplification is achieved by considering the “ $p$ -Borel transform” of the series  $f_\alpha$ :

$$(5) \quad g_\alpha(\tau) := \mathcal{B}_p f_\alpha(\tau) = \sum_{n \geq 0} a_n p^{-n(n-1)/2} \tau^n,$$

where  $a_n$  are the coefficients of  $f_\alpha$ . By the commutation rule  $\mathcal{B}_p(t^m \sigma_p^\ell) = p^{-m(m-1)/2} \tau^m \sigma_p^{\ell-m} \mathcal{B}_p$ ,  $g_\alpha$  is solution of a two-term  $q$ -difference equation. This is easily solved to find

$$g_\alpha(\tau) = \frac{1}{(-\alpha p^\nu \tau; q)_\infty (-\alpha p^{-\nu} \tau; q)_\infty}.$$

It follows that  $g_\alpha$  is meromorphic in  $\mathbb{C}$  with (simple) poles at  $\{-p^{\nu-2n}/\alpha, -p^{-\nu-2n}/\alpha\}$  for  $n \in \mathbb{N}$ , which implies that  $f_\alpha$  is an entire function.

In order to recover  $f_\alpha$  from  $g_\alpha$ , the  $p$ -Borel transform of (5) is reverted by means of a Hadamard product of  $g_\alpha$  with  $\theta_p$ . This leads to a Cauchy integral representation from which a residue computation yields the connection formula. The Cauchy integral is

$$f_\alpha(t) = \frac{1}{2\pi i} \int_{|\tau|=r} g_\alpha(\tau) \theta_p(t/\tau) \frac{d\tau}{\tau},$$

where  $r < \min(|p^\nu/\alpha|, |p^{-\nu}/\alpha|)$ . The only residues come from the poles of  $g_\alpha$ . The asymptotic behaviour of  $g_\alpha$  implies that this integral is equal to the sum of the residues and an actual computation of these residues leads to

$$f_\alpha(t) = \frac{\theta_p(-\alpha q^{\nu/2} t)}{(q; q)_\infty (q^{-\nu}; q)_\infty} {}_2\phi_1(0, 0; q^{\nu+1}; q, -x^2/4) + \frac{\theta_p(-\alpha q^{-\nu/2} t)}{(q; q)_\infty (q^\nu; q)_\infty} {}_2\phi_1(0, 0; q^{-\nu+1}; q, -x^2/4),$$

where  $xt = 1$  and  $|x| < 2$ . With very little rewriting, this is the desired connection formula. The limiting behaviour of this formula when  $q \rightarrow 1$  is studied in [5].

The second family of  $q$ -Bessel functions is actually related to the first one by a relation discovered by Hahn in 1949:

$$J_\nu^{(2)}(x; q) = (-x^2/4; q)_\infty J_\nu^{(1)}(x; q), \quad |x| < 2.$$

Another way of viewing the relation between these functions is through the  $p$ -Laplace transform that sends  $x^n$  to  $p^{n(n-1)/2} x^n$ . Then the transform of the  ${}_2\phi_1$  in the definition of  $J_\nu^{(1)}$  is the  ${}_0\phi_1$  in that of  $J_\nu^{(2)}$ . From there, a Cauchy integral representation follows and again a residue computation gives the connection formula, thanks to extra considerations about the asymptotic behaviour of the integrand.

## 5. Comments

It has been observed that the connection ‘‘constants’’ possess the nice property that they are elliptic in the case of Heine’s function. This is a general phenomenon [3]. The formulae in the  $q$ -world imply important identities (after all, Jacobi’s triple product can be seen as a connection formula). Recent work by Changgui Zhang shows that the limiting behaviour of these  $q$ -connection formulae when  $q \rightarrow 1$  yields the Stokes phenomenon of the differential world.

## Bibliography

- [1] Andrews (George E.), Askey (Richard), and Roy (Ranjan). – *Special functions*. – Cambridge University Press, Cambridge, 1999, xvi+664p.
- [2] Gasper (George) and Rahman (Mizan). – *Basic hypergeometric series*. – Cambridge University Press, Cambridge, 1990, xx+287p. With a foreword by Richard Askey.
- [3] Sauloy (Jacques). – Systèmes aux  $q$ -différences singuliers réguliers : classification, matrice de connexion et monodromie. *Annales de l’Institut Fourier*, vol. 50, n° 4, 2000, pp. 1021–1071.
- [4] Slater (Lucy Joan). – *Generalized hypergeometric functions*. – Cambridge University Press, Cambridge, 1966, xiii+273p.
- [5] Zhang (Changgui). – *Sur la sommabilité des séries entières solutions formelles des équations aux  $q$ -différences linéaires et à coefficients analytiques*. – Mémoire d’habilitation, Université de La Rochelle, 2000.



## On the Convergence of Borel Approximants

Donald Lutz

San Diego State University (USA) & Université d'Angers (France)

October 30, 2000

Summary by Marianne Durand

### Abstract

For some “irregular singular” problems coming from differential equations, there exist formal power series solutions that are everywhere divergent. These power series turn out to make sense as asymptotic expansions of actual solutions. The Borel summation technique is used to recover convergent representations for these actual functions solutions.

### 1. Resummation

Some “irregular singular” problems coming from differential equations have formal power series solutions that are everywhere divergent. By resummation techniques, one can obtain convergent solutions [7, 10]. We consider a power series, solution of a linear differential equation, that is everywhere divergent, noted  $\tilde{x}(z) = \sum_1^\infty x_n z^{-n}$ . We assume that it has *Gevrey order* equal to one, which means that there exist constants  $A$  and  $c$  such that  $|x_n| \leq Ac^n n!$ . For a function  $f(z)$ , holomorphic in an angular sector  $S$  extending to infinity and containing the real positive axis, we say that  $\tilde{x}(z)$  is the *Gevrey expansion of order 1* of  $f(z)$  if there exist constants  $K$  and  $C$  such that

$$\left| f(z) - \sum_1^{N-1} x_n z^{-n} \right| \leq CK^N N! |z|^{-N} \quad \text{when } z \in S \text{ and } N \geq 0.$$

This function  $f$  is a resummation of  $\tilde{x}$ , and it exists if the opening angle of  $S$  is smaller than  $\pi$ .

The formal Borel transform of  $\tilde{x}(z)$  is defined by  $y(z) = \sum_1^\infty \frac{x_n z^{n-1}}{(n-1)!}$ . It converges for  $|z| < \frac{1}{c}$ . We assume that the function  $y$  can be continued analytically along a line that does not meet a singularity. In the particular case when  $x$  is a solution of a linear differential equation with rational coefficients, so does  $y$ , as this property is stable under the Borel transform. Thus  $y$  has a finite number of singularities and verifies the above hypothesis. Up to a possible linear change of variable, we may assume that there is no singularity on the real axis, which implies that  $y$  can be continued analytically on the positive real axis. If  $y$  satisfy the expected growth conditions at infinity, we apply the Laplace transform. This transform is defined by

$$x(z) = \mathcal{L}(y) = \int_0^\infty e^{-zt} y(t) dt,$$

and is convergent for  $\Re(z^a)$  greater than a certain positive constant, the constant  $a$  being made precise later. The asymptotic expansion of  $x(z)$  when  $z \rightarrow 0^+$  is equal to  $\sum_1^\infty x_n z^{-n}$ . The function  $x$  is a solution of the initial differential equation [2, 8].

## 2. Balsler, Lutz, and Schäfke's Technique

The next step is to find a way to compute this function  $x$  quickly and in a large domain. For this, Lutz *et al.* [1] reformulate  $x$  as a convergent series of the type  $x(z) = \sum_0^\infty d_n q_n(z)$ . This series is obtained by introducing a mapping function  $\phi$  that maps  $[0, 1]$  onto  $[0, \infty]$ , so as to write the equation

$$(1) \quad x(z) = \int_0^\infty e^{-zt} y \circ \phi \circ \phi^{-1}(t) dt = \int_0^\infty e^{-zt} \sum_0^\infty d_n \phi^{-1}(t)^n dt,$$

where for the second equality we have used the re-expansion  $y \circ \phi(u) = \sum_0^\infty d_n u^n$  in terms of the sequence  $d_n$ . The sequence  $q_n$  is thus determined by  $q_n = \int_0^\infty e^{-zt} \phi^{-1}(t)^n dt$ , under the assumption that the interversion of the integral and the sum holds, permitting termwise integration. We observe that  $q_n$  does not depend on  $x$  and on the initial problem, but only on the mapping function  $\phi$ . This means that these coefficients can be precomputed. On the other hand the coefficients  $d_n$  correspond to a composition of the function  $\phi$  with the Borel transform  $y$ . This is formalized in the following theorem.

**Theorem 1** (Balsler, Lutz and Schäfke). *Let  $x(z) = \int_0^\infty e^{(-zt)}y(t) dt$  where the function  $y$  is holomorphic in the domain*

$$\mathcal{D} \supset \left\{ \left| \text{Arg}(1 + t/a) \right| < \pi/2p \right\}$$

*and satisfies  $|y(t)|e^{-b|t|} \rightarrow 0$  as  $|t| \rightarrow \infty$  in  $\mathcal{D}$ . Choose  $\phi$  holomorphic in  $\Delta = \{|\tau| < 1\}$  so that  $\phi(\Delta) \subset \mathcal{D}$ ,  $\phi([0, 1]) = [0, \infty]$ , and  $(1 - \tau)^c \phi(\tau) \rightarrow A$  as  $\tau \rightarrow 1$  in  $\Delta$ . Define  $(d_n)$  by its generating series  $y(\phi(\tau)) = \sum_0^\infty d_n \tau^n$ , and  $(q_n)$  by*

$$q_n(z) = \int_0^1 e^{-z\phi(\tau)} \tau^n \phi'(\tau) d\tau \quad \text{for } z \text{ such that } |\text{Arg}(z)| < \pi(1 + c)/2.$$

*Then for suitable positive constants (independent of  $n$ )*

$$|d_n| \leq K e^{Ln^{c/(c+1)}} \quad \text{and} \quad |q_n(z)| \leq \tilde{K} e^{-An^{c/(c+1)}\Re(z^{1/(c+1)})}.$$

*So we have  $x(z) = \sum_0^\infty d_n q_n(z)$  for  $\Re(z^{1/(c+1)})$  large.*

*Proof.* Starting from Equation (1), we obtain  $x(z) = \int_0^\infty \sum_{n=0}^\infty e^{-zt} d_n \phi^{-1}(t)^n dt$ . The saddle-point method gives upper bounds for  $d_n$  and  $q_n$  that allows us to interchange the order of integrand and summation in the equation above for  $\Re(z^{1/(c+1)})$  large enough. This interchange yields the expected result  $x(z) = \sum_{n=0}^\infty d_n q_n(z)$ . □

Some other classical conformal mappings can be found in [6]. Here is an example. The mapping

$$(2) \quad \tau = 1 - \frac{2}{(1 + t/a)^p + 1} \quad \text{with } a \in \mathbb{R} \text{ and } p \geq 1/2$$

takes the sectorial domain defined by  $|\text{Arg}(1 + t/a)| < \pi/(2p)$  onto the unit disk. The choice of the conformal mapping  $\phi$  is important because it has an effect on the speed of convergence and on the area of convergence.

In the particular case where the differential equation is linear with polynomial coefficients, some efficient computation can be done using recurrences. We also suppose now that the function  $\phi$  is algebraic. The precomputation of the coefficients  $q_n$  is based on the fact that they follow a linear recurrence. We first note that the coefficients  $q_n$  are equal to  $\int_0^1 e^{-z\phi(u)} u^n \phi'(u) du$  as shown

by a simple change of variable  $t = \phi(u)$ . The function  $e^{-z\phi(u)}\phi'(u)$  satisfies the first-order linear differential equation

$$(3) \quad G'(t) = \left( \frac{\phi''(t)}{\phi'(t)} - z\phi'(t) \right) G(t).$$

If we note  $\sum_{k=0}^K p_k(n)a(n+k) = 0$  the linear recurrence satisfied by the Taylor coefficients at the origin  $a(n)$  of a power series solution of the equation (3), then the integrals  $q_n(z)$  satisfy the recurrence  $\sum_{k=0}^K p_k(-n)q_{n-k-1}(z) = 0$ . Once we have the recurrence satisfied by the coefficients  $q_n$  and the initial conditions that are given by  $q_n = \int_0^\infty e^{-zt}\phi^{-1}(t)^n dt$ , all the  $q_n$  can be computed quickly. A problem is that we seek for numerical and not exact computations, and so we have, on each example, to seek for numerical stability. This point uses a backward scheme which is developed on an example below.

The computation of the coefficients  $d_n$  can be done efficiently by finding a recurrence for example using the gfun package [9], because it is a composition of a known algebraic function  $\phi$  and a function  $y$  known by its differential equation. The initial conditions for the  $d_n$  derive directly from the initial conditions of the differential equation satisfied by  $y$  and so from the initial conditions of the differential equation satisfied by  $\tilde{x}$ . This is illustrated by the example of the Heun equation.

### 3. Heun Equation

The Heun equation is the generic differential equation with four regular singular points located at 0, 1,  $c$ , and  $\infty$ ; see [5]. The double confluent Heun equation is obtained by letting the singularity located at  $c$  tend to the one located at  $\infty$ , and the singularity located at 1 tend to 0. The equation obtained then has two irregular singular points located at 0 and  $\infty$ . The example we study [3] is the confluent Heun equation in the form

$$(4) \quad z^2 f''(z) + (z + \alpha z^2 + \alpha) f'(z) + \frac{(2\alpha z^2 \beta_1 + \alpha z^2 + \alpha^2 z - 2\gamma z + 2\alpha\beta_{-1} - \alpha)}{2z} f(z) = 0.$$

The acceleration is realised by the function  $\phi = \frac{1}{(1-z)^2} - 1$  which maps from  $[0, 1]$  onto  $[0, \infty]$ . The recurrence satisfied by  $q_n$  is thus

$$(5) \quad q(n) = \frac{(-6 + 3n)q(n-1) + (-2z + 6 - 3n)q(n-2) + (n-2)q(n-3)}{n-2}.$$

The initial conditions, that are easily computed, using the definition of  $q_n$ , correspond to a dominated solution, so any numerical error makes the dominating solution appear. A solution to this problem is to compute the recurrence backwards, which exchanges the roles of dominating and dominated regimes. The idea is to choose arbitrary values for  $q_{N-d}, \dots, q_N$  where  $d$  is the order of the recurrence and  $N$  is a sufficiently large integer. All the values of  $q_n$  for  $n \leq N$  are then computed from these “final” values backwards. This technique is developed in [11]. The dominating solution of Recurrence 5 disappears and so the initial values found differ only by a multiplicative constant  $\lambda$  from the actual initial values. The sequence  $q_n$  thus found has to be multiplied by this constant  $\lambda$  to give the expected sequence  $q_n$ .

For the coefficients  $d_n$ , the recurrence is found easily using gfun. For parameters  $\alpha = -1$ ,  $\beta_{-1} = 1/2$ ,  $\beta_1 = 1/2$ , and  $\gamma = 1/3$ , it is

$$\begin{aligned} 0 = & (6n^2 + 3n^3)a_n + (-93n - 36 - 75n^2 - 18n^3)a_{n+1} \\ & + (568n + 404 + 267n^2 + 42n^3)a_{n+2} + (-1193n - 1176 - 411n^2 - 48n^3)a_{n+3} \\ & + (1042n + 1240 + 291n^2 + 27n^3)a_{n+4} + (-78n^2 - 336n - 480 - 6n^3)a_{n+5} \end{aligned}$$

with initial conditions  $a_0 = 0$ ,  $a_1 = 1$ ,  $a_2 = 1/3$ ,  $a_3 = -23/108$ , and  $a_4 = -2749/3888$ .

Now for each fixed  $z$ , we can compute the value of  $x(z)$  to arbitrary precision, by choosing the number of terms we take into account. The backwards computation for the  $q_n$  coefficients implies that the number of computable terms is limited by the starting point. If it is too low, we have to choose a larger starting point to get more terms. It is generally not possible to decide where a good starting point for the computation of the backward computation would be. This can be done on particular examples, but the starting point strongly depends on  $z$ .

#### 4. Applications

Many problems related to differential equations yield formal power series of Gevrey order one. Whenever the Borel–Laplace transform applies, the results of Section 2 also applies. A concrete application coming from physics is the one-dimensional complex heat equation:

$$u_\tau(\tau, z) = u_{zz}(\tau, z), \quad u(0, z) = \phi(z).$$

The Cauchy data  $\phi(z)$  is assumed to be holomorphic near the origin. A formal solution is

$$\tilde{u}(\tau, z) = \sum_0^\infty \phi^{(2n)}(z) \frac{\tau^n}{n!}.$$

Lutz *et al.* have shown that either  $\tilde{u}(\tau, z)$  is convergent, or the method of Section 2 applies. If  $v(\tau, z)$  is the Borel transform of  $\tilde{u}(\tau, z)$  with respect to  $\tau$ , then applying the Laplace transform in the variable  $\tau$  to  $v(\tau, z)$  for fixed  $z$  gives a convergent solution  $u(\tau, z)$  of the Cauchy problem. Better knowledge on the function  $\phi$  may easily lead to fast rate convergence possibly using the mapping function (2). Another application is about convergent Liouville–Green expansions for second order linear differential equations [4].

#### Bibliography

- [1] Balser (W.), Lutz (D. A.), and Schäfke (R.). – On the convergence of Borel approximants. – Preprint.
- [2] Borel (Émile). – Leçons sur les séries divergentes. In *Collection de monographies sur la théorie des fonctions, publiée sous la direction de M. Émile Borel*. – Gauthiers-Villars, Paris, 1901. Second edition, 1928. Reprinted by J. Gabay, 1988.
- [3] Chyzak (Frédéric), Durand (Marianne), and Salvy (Bruno). – Borel resummation of divergent series using GFUN. – Maple worksheet. Available from <http://algo.inria.fr/libraries/autocomb/>.
- [4] Dunster (T. M.), Lutz (D. A.), and Schäfke (R.). – Convergent Liouville–Green expansions for second-order linear differential equations, with an application to Bessel functions. *Proceedings of the Royal Society. London. Series A*, vol. 440, n° 1908, 1993, pp. 37–54.
- [5] Duval (Anne) and Loday-Richaud (Michèle). – Kovačič’s algorithm and its application to some families of special functions. *Applicable Algebra in Engineering, Communication and Computing*, vol. 3, n° 3, 1992, pp. 211–246.
- [6] Kober (H.). – *Dictionary of conformal representations*. – Dover Publications, New York, N. Y., 1952, xvi+208p.
- [7] Loday-Richaud (Michèle). – Solutions formelles des systèmes différentiels linéaires méromorphes et sommation. *Expositiones Mathematicae*, vol. 13, n° 2-3, 1995, pp. 116–162.
- [8] Ramis (Jean-Pierre). – Séries divergentes et théories asymptotiques. *Bulletin de la Société Mathématique de France (Panoramas et Synthèses)*, vol. 121, 1993, p. 74.
- [9] Salvy (Bruno) and Zimmermann (Paul). – Gfun: a Maple package for the manipulation of generating and holonomic functions in one variable. *ACM Transactions on Mathematical Software*, vol. 20, n° 2, 1994, pp. 163–177.
- [10] Thomann (Jean). – Resommation des séries formelles. Solutions d’équations différentielles linéaires ordinaires du second ordre dans le champ complexe au voisinage de singularités irrégulières. *Numerische Mathematik*, vol. 58, n° 5, 1990, pp. 503–535.
- [11] Wimp (Jet). – *Computation with recurrence relations*. – Pitman (Advanced Publishing Program), Boston, MA, 1984, xii+310p.

**Part VI**

**ALEA'2001 Lecture Notes**



# Enumerative Combinatorics: Combinatorial Decompositions and Functional Equations<sup>†</sup>

*Mireille Bousquet-Mélou*

LABRI, Université Bordeaux 1 (France)

March 26 and 27, 2001

*Summary by Mathias Vandenbogaert*

## Abstract

The primary purpose of this course is the elaboration of methods for providing answers to problems that arise in enumerative combinatorics. The main tool to be used in this respect are (ordinary) generating functions. The objects that will be dealt with are 2-dimensional walks (for which several convexity constraints will be taken into account) and trees. These objects are more generally described as “decomposable” objects. A description of the principal combinatorial decompositions by means of functional equations of generating functions will be presented as an equivalent but more synthetic approach to the use of recurrences. The modelling by generating functions of combinatorial structures like trees and walks will be discussed. The same principles hold for maps, animals, and polyominoes. The “kernel method” and “quadratic method” techniques will be presented. The course will be illustrated by numerous examples.

## 1. Enumeration Problems and the Way of Solving Them

The approach in solving an enumerative problem consists in a combinatorial step that examines the structure of the objects under consideration, and a step that resolves the recurrence relations or functional equations. By observing the structure of the objects, some (recursively definable) property can be translated into a mathematical, non-tautological information on  $a_n$ , the number of objects of size  $n$ . Instead of manipulating recurrence relations, generating functions describing the corresponding functional equations are used:

$$A(t) = \sum_{n \geq 0} a_n t^n = \sum_{A \in \mathcal{A}} t^{|A|}$$

is called the ordinary generating function of the combinatorial class  $\mathcal{A}$  endowed with the size function  $|\cdot|$ , where the number of objects  $a_n$  are to be finite. A power series with coefficients in  $A$  can be written  $\sum_{n \geq 0} a_n t^n$  with  $a_n \in \mathbb{N}$ . Using counting generating functions it can be noticed that paths of various sorts are invariably *algebraic functions*, which are defined as solutions of a polynomial equation [11].

There is a simple correspondence between operations on combinatorial classes of objects and combinations of the associated generating functions. This allows us to derive directly functional relations between generating functions starting from definitions of combinatorial objects.

---

<sup>†</sup>Lecture notes for a course given during the workshop ALÉA'01 in Luminy (France).

1. *Union:*  $A(z) + B(z)$  is the enumerative ordinary generating function of  $\mathcal{A} \sqcup \mathcal{B}$ . If  $a_n$  and  $b_n$  are the numbers of objects of size  $n$  in  $\mathcal{A}$  and  $\mathcal{B}$  respectively, then  $a_n + b_n$  is the number of objects of size  $n$  in  $\mathcal{A} \sqcup \mathcal{B}$ .
2. *Cartesian Product:*  $A(z)B(z)$  is the enumerative ordinary generating function of  $\mathcal{A} \times \mathcal{B}$ . The number of objects of size  $n$  in  $\mathcal{A} \times \mathcal{B}$  equals the simple convolution  $\sum_{0 \leq k \leq n} a_k b_{n-k}$ . Alternatively:

$$\sum_{\gamma \in \mathcal{A} \times \mathcal{B}} z^{|\gamma|} = \sum_{\alpha \in \mathcal{A}} \sum_{\beta \in \mathcal{B}} z^{|\alpha|+|\beta|} = A(z)B(z).$$

3. *Sequences:*  $(1 - A(z))^{-1}$  is the enumerative ordinary generating function of sets of objects of

$$\mathcal{A} = \{(B_1, B_2, \dots, B_k) \mid B \in \mathcal{B}, k \geq 0\}.$$

The cardinality of  $A$  is  $|A| = \sum_{i=1}^k |B_i|$ , and the generating function of  $A(z)$  is

$$A(t) = \sum_{k \geq 0} B(t)^k = (1 - B(t))^{-1}$$

according to the statements of union and cartesian product.

It can be proven that a strong algebraic *decomposability* prevails for directed lattice paths, which is obtained by a specific technique, the “*kernel method*” [2, 6]. The decomposability enables us to determine the location and nature of dominant singularities.

## 2. Enumeration Example

Fix a finite set of vectors of  $\mathbb{Z} \times \mathbb{Z}$ ,  $S = \{(a_1, b_1), \dots, (a_m, b_m)\}$ . A *lattice path* or *walk* relative to  $S$  is a sequence  $v = (v_1, \dots, v_n)$  such that each  $v_j$  is in  $S$ . The geometric realization of a lattice path  $v = (v_1, \dots, v_n)$  is the sequence of points  $(P_0, P_1, \dots, P_n)$  such that  $P_0 = (0, 0)$  and  $\overrightarrow{P_{j-1}P_j} = v_j$ . The quantity  $n$  is referred to as the size of the path. The elements of  $S$  are called *steps* or *jumps*. For these paths, the solution  $F(t, u)$  (which is always an algebraic function of  $t$  and  $u$ ), and combinatorial explanations for the simple formulae obtained from the recurrence relations can be found in [9].

**2.1. Dyck paths.** A classical example can be given with Dyck paths. A Dyck path of length  $2n$  is a path in the plane from  $(0, 0)$  to  $(2n, 0)$  which uses only steps  $(1, 1)$  (North-East), called *rises*, and  $(1, -1)$  (South-East), called *falls*. A Dyck path ends on the  $x$ -axis and does not go below the  $x$ -axis. A Dyck path therefore has even length, with the number of North-East steps equal to the number of South-East steps. A lattice point on the path is called a *peak* if it is immediately preceded by a North-East step and immediately followed by a South-East step [10]. A peak is at height  $k$  if its  $y$ -coordinate is  $k$ . By  $D_n$  we denote the set of all Dyck paths of half-length  $n$ . Obviously,  $D_0 = \{\epsilon\}$ . Every nonempty Dyck path  $\alpha$  can be decomposed uniquely in the following manner [7]:

$$\alpha = u\beta_1 d\gamma_1,$$

when writing  $u$  for a North-East step, and  $d$  for a South-East step, and where  $\beta_1$  and  $\gamma_1$  are possibly empty Dyck paths. This relation implies that

$$D_n = uD_0 dD_{n-1} \cup uD_1 dD_{n-2} \cup \dots \cup uD_{n-2} dD_1 \cup uD_{n-1} dD_0, \quad n \geq 1.$$

Alternatively, we can write  $\alpha = \beta_2 u\gamma_2 d$  in a unique manner, where  $\beta_2$  and  $\gamma_2$  are possibly empty Dyck paths. This relation implies that

$$D_n = D_0 uD_{n-1} d \cup D_1 uD_{n-2} d \cup \dots \cup D_{n-2} uD_1 d \cup D_{n-1} uD_0 d, \quad n \geq 1.$$



Both equations have disjoint unions. Thus we obtain

$$|D_n| = |D_0||D_{n-1}| + |D_1||D_{n-2}| + \cdots + |D_{n-2}||D_1| + |D_{n-1}||D_0|, \quad n \geq 1.$$

As  $|D_0| = 1$ , this sequence with  $n \geq 0$  satisfies the same recurrence relation as the sequence  $(c_n)_{n \geq 0}$  of Catalan numbers.

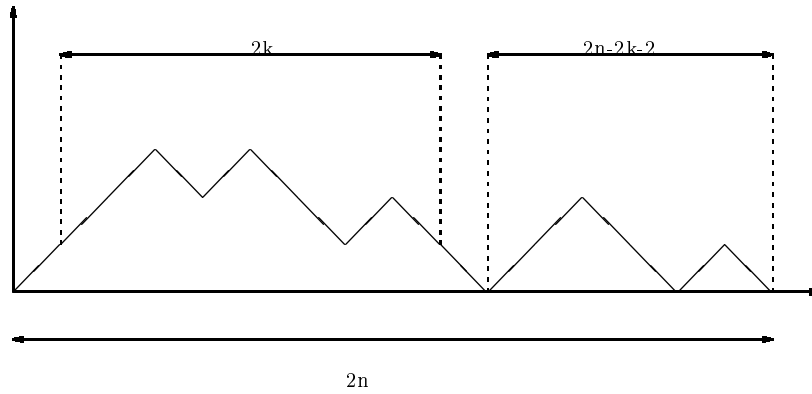
**2.2. Enumeration of Dyck paths.** Let  $p$  be a fixed nonnegative integer-valued parameter of a Dyck path, i.e., a mapping from  $\bigcup_{n \geq A} D_n$  into  $\{0, 1, 2, \dots\}$ . If  $D$  is a finite set of Dyck paths, then by  $D(t)$  we denote the *enumerating polynomial* of  $D$  relative to the parameter  $p$  given by

$$D(t) = \sum_{n \geq 0} d_n t^n \quad \text{with} \quad d_n = \sum_{\delta \in D} t^{p(\delta)}.$$

$D(t)$  is the generating function for the enumeration of Dyck paths according to semi-length (coded by  $t$ ). Thus,  $d_n$  is the enumerating polynomial of the set of all Dyck paths of length  $n$ .

The recurrence relation for Dyck paths satisfies

$$\begin{cases} d_{2n} = \sum_{k=0}^{n-1} d_{2k} d_{2n-2k-2}, & n \geq 1, \\ d_0 = 1. \end{cases}$$



This gives on summation:

$$D(t) = \sum_{n \geq 0} d_{2n} t^{2n} = 1 + \left( \sum_{n \geq 1} t^{2n} \right) \left( \sum_{k=0}^{n-1} d_{2k} d_{2n-2k-2} \right) = 1 + t^2 D(t)^2.$$

This quadratic equation is easily solved for  $D(t)$ :

$$D(t) = \frac{1 \pm \sqrt{1 - 4t^2}}{2t^2}.$$

The solution  $\frac{1 - \sqrt{1 - 4t^2}}{2t^2}$  is chosen in order to ascertain the existence of a Taylor series expansion at  $t = 0$ . It is known [2, 7, 8, 10, 11] that the number of Dyck paths of length  $2n$  is  $c_n$ , the  $n$ th Catalan number, given by  $c_n = \frac{1}{n+1} \binom{2n}{n}$ .

**2.3. Enumeration of Dyck prefixes.** Let  $b_{n,k}$  be the number of prefixes of length  $n$ , with final height  $k$ . Then

$$b_{n,0} = d_n \quad (\text{Dyck paths})$$

$$F(t, u) = \sum_{n, t \geq 0} b_{n,k} t^n u^k \in \mathbb{Q}[[t, u]] = \left( \sum_{n \geq 0} t^n \right) \left( \sum_{k=0}^n b_{n,k} u^k \right) \in \mathbb{Q}[u][[t]]$$

which is a series in  $t$  whose coefficients are polynomials in  $u$ . The last equation is equivalent to:

$$F(t, u) = 1 + t(u + u^{-1})F(t, u) - tu^{-1}F(t, 0),$$

which defines the generating function  $F(t, u)$  for these paths, counted by their length (variable  $t$ ) and their height (variable  $u$ ). This equation uniquely defines  $F(t, u)$  as a power series in  $t$  with polynomial coefficients in  $u$ .

More constraints can be imposed on such Dyck prefixes.

2.3.1. *Dyck Paths with no peaks at height  $m$ .* Let  $G_m(x) = \sum_{n \geq 0} g(m, n)x^n$  be the generating function for Dyck paths of length  $2n$  with no peaks at height  $m$  for some fixed  $m \geq 1$ . We proceed to show that

$$G_m(x) = \frac{1}{1 - xG_{m-1}(x)} \quad \text{for } m \geq 2.$$

This can be illustrated by a path starting with a North-East step followed by a segment which represents any Dyck path of length  $2k$ ,  $0 \leq k \leq n-1$ , with no peaks at height  $m-1$ . This segment is followed, after a South-East step, by a second segment which represents any Dyck path of length  $2n-2-2k$  with no peaks at height  $m$ . Therefore

$$g(m, 0) = 1$$

and

$$g(m, n) = \sum_{k=0}^{n-1} g(m-1, k)g(m, n-1-k) = [x^{n-1}](G_{m-1}(x)G_m(x)).$$

Thus,

$$G_m(x) = 1 + xG_{m-1}(x)G_m(x),$$

or equivalently,

$$G_m(x) = \frac{1}{1 - xG_{m-1}(x)}.$$

This way, the number of Dyck paths of length  $2n$  with no peaks at height 1 is the Fine number  $f_n$  for  $n \geq 0$ . Obviously,  $g(1, 0) = 1$  and  $g(1, 1) = 0$ . For  $n \geq 2$ , a Dyck path of length  $2n$  with no peaks at height 1 has a segment representing any Dyck path of length  $2k$ ,  $1 \leq k \leq n-1$ , and a second segment representing a Dyck path of length  $2n-2k-2$  with no peaks at height 1. Therefore, for  $n \geq 2$ , we have

$$\begin{aligned} g(w, n) &= \sum_{k=1}^{n-1} c_k g(w, n-k-1) \\ &= [x^{n-1}](C(x)G_1(x)) - g(1, n-1) \\ &= [x^n](xC(x)G_1(x)) - g(1, n-1). \end{aligned}$$

Therefore,

$$\begin{aligned} G_1(x) &= 1 + \sum_{n \geq 2} g(1, n)x^n = 1 + xC(x)G_1(x) - x - xG_1(x) + x \\ &= 1 + xG_1(x)(C(x) - 1) = 1 + xG_1(x)xC^2(x). \end{aligned}$$

That is,

$$G_1(x) = \frac{1}{1 - x^2C^2(x)}.$$

2.3.2. *No peaks at height 2.* Another extension establishes that the number of Dyck paths of length  $2n$  with no peaks at height 2 is the Catalan number  $c_{n-1}$ , for  $n \geq 1$ . This can be shown [10] using the first extension, so that

$$G_2(x) = \frac{1}{1 - xG_1(x)} = \frac{1}{1 - x\frac{C(x)}{1+xC(x)}} = 1 + xC(x).$$

2.4. **Bilateral paths or bridges.** A bridge is a path whose end point  $P_n$  lies on the  $x$ -axis. Given a class  $\mathcal{C}$  of paths, we let  $\mathcal{C}_n$  denote the subclass of paths that have size  $n$ , and  $\mathcal{C}_{n,k} \subset \mathcal{C}$  those that have final altitude equal to  $k$ . We introduce the corresponding ordinary generating functions:

$$C(z) = \sum_n C_n z^n, \quad uC(z, u) = \sum_{n,k} C_{n,k} u^k z^n.$$

By characterising these generating functions, that are algebraic in the case of bridges, a strong algebraic decomposition prevails, which renders the calculation of the generating function's effective. The decomposability of generating function's makes it possible to extract their singular structure, and to solve the corresponding asymptotic enumeration problems.

The equation corresponding to such a lattice path is:

$$B(t) = 1 + t^2 D(t)B(t) + t^2 B(t)D(t) = \frac{1}{1 - 2t^2 D(t)}.$$

For  $D(t) = \frac{1 \pm \sqrt{1-4t^2}}{2t^2}$ ,

$$B(t) = \frac{1}{1 - 1 \pm \sqrt{1 - 4t^2}} = \frac{1}{\sqrt{1 - 4t^2}} = \sum_{n \geq 0} t^{2n} \binom{2n}{n}.$$

Alternatively, since  $\sqrt{1 - 4t^2} = 1 - 2t^2 - 2t^4 + O(t^6)$ , we can find for Dyck paths:

$$D(t) = \frac{1 + 1 - 2t^2 - 2t^4 + O(t^6)}{2t^2} = \frac{1}{t^2} + 1 - t^2 + O(t^4)$$

or

$$D(t) = \frac{1 - \sqrt{1 - 4t^2}}{2t^2},$$

which is the result we found before.

### 3. Lagrange Inversion Formula

Inherently to the symbolic method, the extraction of coefficients of generating functions defined by functional equations is a frequently occurring problem. For this purpose, the Lagrange Inversion Theorem provides a tool that is commonly used and especially dedicated to the enumeration of trees. This theorem states that given the generating function  $A(z) = \sum_{n \geq 0} a_n z^n$  for which  $z = f(A(z))$ , if  $f(z)$  verifies the condition  $f(0) = 0$  and  $f'(0) \neq 0$ , then

$$a_n \equiv [z^n]A(z) = \frac{1}{n} [u^{n-1}] \left( \frac{u}{f(u)} \right)^n.$$

Additionally,

$$[z^n](A(z))^m = \frac{m}{n} [u^{n-m}] \left( \frac{u}{f(u)} \right)^n$$

and

$$[z^n]g(A(z)) = \frac{1}{n} [u^{n-1}]g'(u) \left( \frac{u}{f(u)} \right)^n.$$

By application of the reciprocal function to both sides of the equation  $z = f(A(z))$ , it can be noticed that the function  $A(z)$  is the reciprocal of  $f(z)$ . The surprising effect of the inversion theorem resides in the relation it establishes between the powers of a function and the coefficients of the reciprocal function.

**3.1. Example: Catalan numbers.** The language of Dyck words,

$$\mathcal{D} = \{\epsilon, x\bar{x}, xx\bar{x}\bar{x}, x\bar{x}x\bar{x}, \dots\},$$

satisfies the defining recurrence  $\mathcal{D} = \epsilon + x\mathcal{D}\bar{x}\mathcal{D}$ . This translates to the algebraic (non-commutative) equation

$$D(x, \bar{x}) = 1 + xD(x, \bar{x})\bar{x}D(x, \bar{x}).$$

Since we have an algebraic and non-ambiguous grammar, we can rewrite the system with commutative variables:

$$D(x, \bar{x}) = 1 + x\bar{x}D(x, \bar{x})^2.$$

As we know that the length of the words is always even, we will have  $n$  for a total length of  $2n$ , when we only count  $x$  (or  $\bar{x}$ ). Thus, we can substitute  $\bar{x}$  for  $\epsilon$ , and  $x$  for  $t$ .

$$D(t) = 1 + t(D(t))^2 \iff tD(t)^2 - D(t) + 1 = 0$$

By simply solving this second-order equation, we get  $D(t) = \frac{1 - \sqrt{1 - 4t}}{2t}$  (the other root is negative, hence not applicable). This solution is converted into the form  $D(t) = \sum_{n \geq 0} a_n t^n$ , for which  $a_n$  gives us the number of Dyck words having  $n$  letters  $t$  ( $x$ ), hence the number of Dyck words of length  $2n$ . Using Taylor series expansion and applying the Lagrange Inversion Formula, we get  $C_n$

$$\frac{1 - \sqrt{1 - 4t}}{2t} = \sum_{n \geq 0} \frac{1}{n+1} \binom{2n}{n} t^n$$

$$[z^n]C(t) = [z^n] \frac{1}{n} z^{n-1} (1+z)^{2n} = \frac{1}{n} \binom{2n}{n-1}.$$

## 4. Algebraic Structures and the Kernel Method

**4.1. Algebraic equations.** The equation describing sub-diagonal North-East paths,

$$F(t, u) = 1 + t(u + 1/u)F(t, u) - t/uF(t, 0),$$

belongs to a class of equations that share two properties [3]:

1. The equation uniquely defines  $F(t, u)$  as a power series in  $t$  with polynomial coefficients in  $u$ . There exist other, non-power-series solutions, for instance the rational function

$$F(t, u) = \frac{2tu - 1}{2t(u - t(u^2 + 1))}.$$

Hence, any method for solving the recurrence relation above must use the fact that  $F(t, u)$  is a power series.

2. When trying to derive an equation in  $F(t, 0)$  only from the recurrence relation, we end up with a tautologic expression. In other words, if we first multiply  $F(t, u)$  by  $u$  and directly set  $u = 0$ , this would give us  $0 = tF(t, 0) - tF(t, 0)$ .

It can be noticed that the recurrence relation is linear in  $F(t, u)$  and  $F(t, 0)$ , and we can strongly expect its solution to be algebraic and to satisfy

$$F(t, 0) = \frac{1 - \sqrt{1 - 4t^2}}{2t^2} = \sum_{n \geq 0} \frac{1}{n + 1} \binom{2n}{n},$$

since sub-diagonal walks ending on the main diagonal are well-known to be counted by Catalan numbers.

The generic form of equations that share the above properties, is

$$P(F(t, u), F_1(t), F_2(t), \dots, F_k(t), t, u) = 0,$$

where  $P$  is a polynomial in  $k+3$  variables with real coefficients. We assume that this equation defines uniquely all its unknowns as power series in  $t$ : the series  $F_i(t)$  have real coefficients, while  $F(t, u)$  has its coefficients in  $\mathbb{R}[u]$ . Rewriting our equation according to this generic form of equations yields:

$$F(t, u)(u - t(u^2 + 1)) - u + tF_1(t) = 0,$$

with  $F_1(t) = F(t, 0)$ , by setting  $u = 0$ .

In solving this instance, we propose to determine  $f_n$ , the number of excursions of length  $n$  and type  $\Omega$ , the set of jumps which is a finite subset of  $\mathbb{Z}$ , via the corresponding bivariate generating function

$$F(z, u) = \sum_{n,k} f_{n,k} u^k z^n,$$

where  $f_{n,k}$  is the number of walks of length  $n$  and final altitude  $k$ . In particular,  $F(z) = F(z, 0)$ . Let  $-c$  denote the smallest (negative) value of a jump, and  $d$  denote the largest (positive) jump. A functional role is played by the “characteristic polynomial” of the walk [1, 2, 11],

$$S(y) = \sum_{\omega \in \Omega} y^\omega = \sum_{j=-c}^d S_j y^j,$$

which is a Laurent polynomial. The bivariate generating function of generalised walks where intermediate values are allowed to be negative is rational:

$$G(z, u) = \frac{1}{1 - zS(u)}.$$

The main result to be proven is the following: for each finite set  $\Omega \subset \mathbb{Z}$ , the generating function of excursions is an algebraic function that is explicitly computable from  $\Omega$ . This problem is solved by an application of the *kernel method* [2].

**4.2. Kernel method.** [2]. Let  $f_n(u) = [z^n]F(z, u)$  be the generating function of walks of length  $n$  with  $u$  recording the final altitude. There is a simple recurrence relating  $f_{n+1}(u)$  to  $f_n(u)$ , namely,

$$f_{n+1}(u) = S(u)f_n(u) - r_n(u)$$

where  $r_n(u)$  is a Laurent polynomial consisting of the sum of all the monomials of  $S(u)f_n(u)$  that involve negative powers of  $u$ :

$$r_n(u) = \sum_{j=-c}^{-1} ([u^j]S(u)f_n(u))u^j = \{u^{<0}\}S(u)f_n(u),$$

where  $\{u^{<0}\}$  denotes the singular part of a Laurent expansion:

$$\{u^{<0}\}f(z) := \sum_{j < 0} ([u^j]f(u))u^j.$$

The idea behind the formula is to subtract the effect of those steps that would take the walk below the horizontal axis. Thus the generating function  $F(z, u)$  satisfies the fundamental functional equation

$$F(z, u) = 1 + zS(u)F(z, u) - z\{u^{<0}\}(S(u)F(z, u)).$$

Explicitly, we have

$$F(z, u) = 1 + zS(u)F(z, u) - z \sum_{j=0}^{c-1} \lambda_j(u) \left( \frac{\delta^j}{\delta u^j} F(z, u) \right)_{u=0},$$

for Laurent polynomials  $\lambda_j(u)$  that depend on  $S(u)$  in an effective way by  $\lambda_j(u) = \frac{1}{j!} \{u^{<0}\} u^j S(u)$  [2].

Both equations involve an unknown bivariate generating function  $F(z, u)$  and  $c$  univariate generating functions, the partial derivatives of  $F$  specialized at  $u = 0$ . In particular, the latter functional equation determines fully the  $c + 1$  unknowns. The basic technique is known as “cancelling the kernel” and relies on strong analyticity properties.

The equation to be used by the basic kernel technique starts by grouping on one side the terms involving  $F(z, u)$ . The main principle of the kernel method consists in coupling the values of  $z$  and  $u$  in such a way that  $1 - zS(u) = 0$ , so that  $F(z, u)$  disappears. Consequently, the “kernel equation”  $1 - zS(u) = 0$ , is rewritten as  $u^c = z(u^c S(u))$ . This kernel equation defines  $c + d$  branches of an algebraic function. Coupling  $z$  and  $u$  by  $u = u_l(z)$  gives that  $(z, u)$  is close to  $(0, 0)$  where  $F$  is bivariate analytic, so that substitution gives

$$1 - z \sum_{j=0}^{c-1} \lambda_j(u_l(z)) \left( \frac{\delta^j}{\delta u^j} F(z, u) \right)_{u=0}, \quad l = 0, \dots, c - 1,$$

which is a linear system of  $c$  equations in  $c$  unknowns with algebraic coefficients that determines  $F(z, 0)$ . Therefore, the generating function of excursions is expressible as

$$F(z) = \frac{(-1)^{c-1}}{zS_{-c}} \prod_{l=0}^{c-1} u_l(z), \text{ where } S_{-c} = [u^{-c}]S(u)$$

is the multiplicity of the smallest element  $-c \in \Omega$ .

More generally the bivariate generating function of nonnegative walks is bivariate algebraic and given by

$$F(z, u) = \frac{1}{u^c - z(u^c S(u))} \prod_{l=0}^{c-1} (u - u_l(z)).$$

In other words, to make explicit the solution  $F_s$  of the recurrence of the sub-diagonal North-East paths, written as  $F(t, u)(u - t(u^2 + 1)) - u + tF_1(t) = 0$ , we rewrite it as  $Q(x)F(x) = K(x) - U(x)$  [4], where  $K$  stands for the *unknown* initial conditions, and  $Q$  is the kernel:

$$F(t, u)(u - t(u^2 + 1)) = u - U(t),$$

$$F_s(t)Q(t) = K(t) - U(t).$$

Again, the *kernel method* consists in cancelling the kernel  $Q(x)$ , by handling a choice of algebraic values  $a$  of  $t$ , which yields a system of equations  $K(a) - U(a) = 0$ . Solving this system generally allows to make  $U$  explicit. This provides  $F_s$  for generic  $t$ :

$$F_s(t) = \frac{K(t) - U(t)}{Q(t)}.$$

The function  $U(t)$  is a sum of  $m$  unknown multivariate functions  $F_i(t_1, \dots, t_{d-1})$ . Cancelling the kernel with  $m$  different values for  $t_d$  (which then become functions of  $(t_1, \dots, t_{d-1})$ ) yields a system which allows to make explicit the  $F_i$ 's.

Regrouping the terms in  $F(t, u)$  by the kernel method yields:

$$F_s(t, u) = \frac{u - \frac{1 - \sqrt{1 - 4t^2}}{2t}}{u - t - tu^2}.$$

**4.3. The Quadratic Method.** An analogous approach is referred to as the “quadratic method,” used to solve equations of the form

$$z(x, y)^2 + P_1(x, y, z(x, 0))z(x, y) + P_2(x, y, z(x, 0)) = 0$$

with  $P_i \in \mathcal{F}[[x]][y, u]$ , where  $\mathcal{F}$  is an algebraically closed field of characteristic zero.

Rewrite the equation as

$$\left(z + \frac{1}{2}P_1\right)^2 = \frac{1}{4}P_1^2 - P_2 =: \Delta \in \mathcal{F}[[x]][y, u].$$

If some  $y = y_0 \in \mathcal{F}[[x]]$  is known to kill  $z + \frac{1}{2}P_1$ , then this  $y_0$  is a double root of  $\Delta(x, y, u)$ , viewed as a polynomial in  $(\mathcal{F}[[x]][u])[y]$ . The resultant  $R(x, u)$  of  $\Delta$  and  $\frac{\partial \Delta}{\partial y}$  with respect to  $y$  has to be zero. When we know by an external argument that the quadratic equation admits a series solutions  $z(x, y) \in \mathcal{F}[[x, y]]$ , for example when it has a combinatorial interpretation, and therefore that  $z(x, 0)$  is a series in  $\mathcal{F}[[x]]$ , the polynomial equation  $R(x, z(x, 0)) = 0$  delivers this value in  $\mathcal{F}[[x]]$  for  $z(x, 0)$ .

After substitution, there only remains to solve an equation of the form  $z^2 + P_1z + P_2 = 0$  with  $P_i \in \mathcal{F}[[x]][y]$ . In [5], necessary and sufficient conditions are derived in order that such an equation has a solution  $z$  in either of the rings  $\mathcal{F}[[x]][y]$  or  $\mathcal{F}[[x, y]]$ . In view of obtaining them, resume from the relation  $(z + \frac{1}{2}P_1)^2 = \Delta$ . Since  $P_1 \in \mathcal{F}[[x]][y] \subset \mathcal{F}[[x, y]]$ , we get that there is a solution in  $\mathcal{F}[[x]][y]$  or  $\mathcal{F}[[x, y]]$ , respectively, if and only if  $\Delta$  has a square root in the same ring. Again in [5], it is proved that  $U \in \mathcal{F}[[x]][y]$  has a square root if and only if it factors under the form  $U = P^2R$  for a polynomial  $P \in \mathcal{F}[x, y]$  and a series  $R \in 1 + y\mathcal{F}[[x]][y]$ . Therefore, the equation has a solution in  $\mathcal{F}[[x]][y]$  or  $\mathcal{F}[[x, y]]$ , respectively, if and only if  $\Delta$  rewrites under the form  $P^2R$  for some polynomial  $P$  and some series  $R$  of the form

$$R = 1 + y\mathcal{F}[[x]][y] \quad \text{or} \quad R = 1 + xy\mathcal{F}[[x]][y], \quad \text{respectively.}$$

### Bibliography

- [1] Banderier (C.), Bousquet-Mélou (M.), Denise (A.), Flajolet (P.), Gardy (D.), and Gouyou-Beauchamps (D.). – Generating functions for generating trees. *Discrete Mathematics*. – 25 pages. To appear.
- [2] Banderier (Cyril) and Flajolet (Philippe). – Basic analytic combinatorics of directed lattice paths. – Available from <http://algo.inria.fr/flajolet/Publications/BaFl01.ps.gz>, August 2001. 39 pages. Accepted for publication in *Theoretical Computer Science*.
- [3] Bousquet-Mélou (Mireille). – On (some) functional equations arising in enumerative combinatorics. – In preparation. Extended abstract for FPSAC'2001 available from <http://dept-info.labri.u-bordeaux.fr/~bousquet/>.

- [4] Bousquet-Mélou (Mireille) and Petkovšek (Marko). – Linear recurrences with constant coefficients: the multivariate case. *Discrete Mathematics*, vol. 225, n° 1-3, 2000, pp. 51–75. – Formal power series and algebraic combinatorics (Toronto, ON, 1998).
- [5] Brown (William G.). – On the existence of square roots in certain rings of power series. *Mathematische Annalen*, vol. 158, 1965, pp. 82–89.
- [6] Cori (Robert) and Richard (Jean). – Énumération des graphes planaires à l'aide des séries formelles en variables non commutatives. *Discrete Mathematics*, vol. 2, 1972, pp. 115–162.
- [7] Deutsch (Emeric). – Dyck path enumeration. *Discrete Mathematics*, vol. 204, n° 1-3, 1999, pp. 167–202.
- [8] Duchon (Philippe). – On the enumeration and generation of generalized Dyck words. *Discrete Mathematics*, vol. 225, n° 1-3, 2000, pp. 121–135. – Formal power series and algebraic combinatorics (Toronto, ON, 1998).
- [9] Dvoretzky (A.) and Motzkin (Th.). – A problem of arrangements. *Duke Mathematical Journal*, vol. 14, 1947, pp. 305–313.
- [10] Peart (Paul) and Woan (Wen-Jin). – Dyck paths with no peaks at height  $k$ . *Journal of Integer Sequences*, vol. 4, n° 1, 2001. – Article 01.1.3. 6 pages.
- [11] Sedgewick (Robert) and Flajolet (Philippe). – *An introduction to the analysis of algorithms*. – Addison-Wesley Publishing Co., Reading, MA, 1996.



## Symbolic Enumerative Combinatorics and Complex Asymptotic Analysis<sup>†</sup>

*Philippe Flajolet*

Projet ALGO, INRIA Rocquencourt (France)

March 26 and 27, 2001

*Summary by Yvan Le Borgne*

### Abstract

Complex analysis is a fruitful source of asymptotic estimates in enumerative combinatorics. This lecture starts with a symbolic method to encode counting sequences of combinatorial structures by complex functions. The residue theorem is then applied to extract from these functions the asymptotic behavior of the corresponding sequences.<sup>1</sup>

A *class of combinatorial structures* (often simply called a *class*) is a finite or countable set on which a *size* function is defined, the size of an element being a nonnegative integer. If  $\mathcal{A}$  is a class, the size of an element  $\alpha \in \mathcal{A}$  is denoted by  $|\alpha|$ , or  $|\alpha|_{\mathcal{A}}$  in the few cases where the underlying class needs to be made explicit. Given a class  $\mathcal{A}$  we consistently let  $\mathcal{A}_n$  be the set of elements in  $\mathcal{A}$  that have size  $n$  and use the same group of letters for the counts  $A_n = \text{Card } \mathcal{A}_n$ . We further assume that the  $\mathcal{A}_n$  are all finite. The *counting sequence* of  $\mathcal{A}$  is the sequence of integers  $\{A_n\}_{n \geq 0}$ . For instance, binary sequences are combinatorial structures that form a class  $\mathcal{S}$  when the size of a word is defined to be its length. The corresponding counting sequence is then given by  $S_n = 2^n$ .

Average-case analysis of algorithms typically reduces to counting problems for combinatorial structures. Statistical physics is another field of application of counting sequences where the free energy of a system may be expressed as the logarithm of the number of accessible states which can be described by a combinatorial structure.

There are two main approaches to estimate the asymptotic behavior of the counting sequence of a class. The first one is to embed the combinatorial structure in a stochastic model where the randomly chosen element is representative of the elements of the class. This allows to eliminate rare pathological elements. Then the asymptotic behavior of the counting sequence is deduced from the behavior of the stochastic model. The second approach, which will be described here, is based on the decomposition of elements of the class into combination of elements of simpler classes and lower size. Counting sequences are encoded by formal generating functions that can have tractable compact representations as complex functions. A restriction to certain combinations, called admissible constructions, preserves these tractable representations since they directly translate into simple operators on the complex functions of the subclasses. The extraction of the counting sequence encoded by a complex function is sometimes difficult, but complex analysis can often be used to obtain the asymptotic behavior.

This summary presents in the first section a symbolic method to compute a function encoding the counting sequence of a class. The second section is dedicated to complex analysis. The aim is

---

<sup>†</sup>Lecture notes for a course given during the workshop ALÉA'01 in Luminy (France).

<sup>1</sup>This summary is inspired by the book in preparation of Flajolet and Sedgewick [2, 3].

to give a method to extract the asymptotic behavior of a counting sequence encoded by a complex function. The final section illustrates these methods throughout two examples: clouds and  $\Omega$ -trees.

## 1. A Symbolic Method for Enumerative Combinatorics

A counting sequence  $\{A_n\}_{n \geq 0}$  can be encoded by different types of formal power series: an ordinary generating function  $\sum_{n \geq 0} A_n z^n$ , an exponential generating function  $\sum_{n \geq 0} \frac{A_n}{n!} z^n$ , a Dirichlet series  $\sum_{n \geq 0} \frac{A_n}{n^z}$ , ... The aim of these representations is to lead in some cases to a description of a counting sequence shorter than the sequence itself. For instance the class  $\mathcal{N}$  of natural integers, where the size of  $n$  is  $n$ , is such that  $N_n = 1$ . Its ordinary generating function is  $\sum_{n \geq 0} z^n = \frac{1}{1-z}$ , its exponential generating function is  $\sum_{n \geq 0} \frac{1}{n!} z^n = e^z$ , its Dirichlet series  $\sum_{n > 0} \frac{1}{n^z} = \zeta(z)$ .

Assume that  $\Phi$  is a binary construction that associates to two classes  $\mathcal{B}$  and  $\mathcal{C}$  a new class

$$\mathcal{A} = \Phi\{\mathcal{B}, \mathcal{C}\},$$

in a finite way (each  $\mathcal{A}_n$  depends on finitely many of the  $\mathcal{B}_n$  and  $\mathcal{C}_n$ ). Then  $\Phi$  is an *admissible construction* if and only if the counting sequence  $\{A_n\}$  of  $\mathcal{A}$  is a function of the counting sequences  $\{B_n\}$  and  $\{C_n\}$  of  $\mathcal{B}$  and  $\mathcal{C}$  only. In that case, this function may be translated into a simple operator relating formal power series representing  $\{A_n\}_{n \geq 0}$ ,  $\{B_n\}_{n \geq 0}$ , and  $\{C_n\}_{n \geq 0}$ . This section is devoted to some particular admissible constructions in the case of unlabeled and labeled combinatorial structures. The goal is to define a language of elementary combinatorial constructions such that any expression of a class in this language can be translated straightforwardly into a function encoding the counting sequence of the class.

**1.1. Unlabeled structures.** The principle of this representation is that an element of size  $n$  is encoded by the monomial  $z^n$ . Thus the class  $\mathcal{A}$  is mapped to the ordinary generating function

$$A(z) = \text{ogf}(\mathcal{A})(z) = \sum_{\alpha \in \mathcal{A}} z^{|\alpha|} = \sum_{n \geq 0} A_n z^n.$$

An additional assumption on the sizes is made: if an element  $\alpha$  can be decomposed into a combination of elements  $\beta_1, \beta_2, \dots, \beta_k$ , then the size of  $\alpha$  is the sum of the sizes of the  $\beta_i$ . Its translation as regards monomials is the usual product law:

$$z^{|\alpha|_{\mathcal{A}}} = z^{|\beta_1|_{\mathcal{B}_1}} z^{|\beta_2|_{\mathcal{B}_2}} \dots z^{|\beta_k|_{\mathcal{B}_k}}.$$

Let us consider the class  $\mathcal{A}$  defined as the *Cartesian product* of two given classes  $\mathcal{B}$  and  $\mathcal{C}$ . Following the additional assumption, the size of the element  $\alpha = (\beta, \gamma)$  is  $|\beta|_{\mathcal{B}} + |\gamma|_{\mathcal{C}}$ . Thus we have

$$A(z) = \sum_{(\beta, \gamma) \in \mathcal{B} \times \mathcal{C}} z^{|\beta, \gamma|_{\mathcal{A}}} = \sum_{\beta \in \mathcal{B}, \gamma \in \mathcal{C}} z^{|\beta|_{\mathcal{B}} + |\gamma|_{\mathcal{C}}} = \sum_{\beta \in \mathcal{B}} z^{|\beta|_{\mathcal{B}}} \cdot \sum_{\gamma \in \mathcal{C}} z^{|\gamma|_{\mathcal{C}}} = B(z)C(z).$$

Here is the first example of an admissible construction which has a simple translation in terms of ordinary generating functions:

$$\text{ogf}(\mathcal{B} \times \mathcal{C})(z) = \text{ogf}(\mathcal{B})(z) \cdot \text{ogf}(\mathcal{C})(z).$$

The *union* of two classes  $\mathcal{B}$  and  $\mathcal{C}$  is translated into the sum of the two ordinary generating functions in the case of a disjoint union. More generally,

$$\text{ogf}(\mathcal{B} \cup \mathcal{C})(z) = \sum_{\alpha \in \mathcal{B} \cup \mathcal{C}} z^{|\alpha|_{\mathcal{A}}} = \sum_{\beta \in \mathcal{B}} z^{|\beta|_{\mathcal{B}}} + \sum_{\gamma \in \mathcal{C}} z^{|\gamma|_{\mathcal{C}}} - \sum_{\alpha \in \mathcal{B} \cap \mathcal{C}} z^{|\alpha|_{\mathcal{B} \cup \mathcal{C}}} = \text{ogf}(\mathcal{B}) + \text{ogf}(\mathcal{C}) - \text{ogf}(\mathcal{B} \cap \mathcal{C}).$$

The additional assumption on the sizes implies that the size of an element  $\alpha$  of  $\mathcal{B} \cap \mathcal{C}$  is well defined since  $|\alpha|_{\mathcal{B}} = |\alpha|_{\mathcal{B} \cup \mathcal{C}} = |\alpha|_{\mathcal{C}}$ .

The class  $\mathcal{A}$  of finite *sequences* of elements of the class  $\mathcal{B}$  is denoted  $\text{Seq}(\mathcal{B})$ . It is well defined if and only if the class  $\mathcal{B}$  has no element of size zero, a restriction which prevents from getting an infinite number of sequences of size zero. Grouping sequences of the same length yields the relation

$$\text{Seq}(\mathcal{B}) = \{\epsilon\} \cup \mathcal{B} \cup (\mathcal{B} \times \mathcal{B}) \cup (\mathcal{B} \times \mathcal{B} \times \mathcal{B}) \cup \dots,$$

where  $\epsilon$  is an element of size zero which has essentially the same meaning as the empty word in the context of languages. Thus, using both previous constructions,

$$\text{ogf}(\text{Seq}(\mathcal{B})) = 1 + \text{ogf}(\mathcal{B}) + \text{ogf}(\mathcal{B})^2 + \text{ogf}(\mathcal{B})^3 + \dots = \sum_{k \geq 0} \text{ogf}(\mathcal{B})^k = \frac{1}{1 - \text{ogf}(\mathcal{B})}.$$

The class  $\mathcal{A}$  of *subsets* of the class  $\mathcal{B}$  is denoted  $\text{Set}(\mathcal{B})$ . The class of directed *cycles* of the class  $\mathcal{B}$  is denoted  $\text{Cycle}(\mathcal{B})$ . Directed cycles are sequences defined up to cyclic permutations: two sequences  $(\alpha_1, \dots, \alpha_k)$  and  $(\beta_1, \dots, \beta_k)$  represent the same directed cycle if and only if there exists an integer  $l$  such that for all  $i$ ,  $\alpha_i = \beta_{i+l \bmod k}$ . These two constructions admit almost reasonable translations mentioned at the end of this section.

**1.2. Labeled structures.** Many objects of classical combinatorics present themselves naturally as labeled structures whose “atom” (typically nodes in a graph or a tree) bear distinctive integer labels. For instance the cycle decomposition of a permutation represents the permutation as an unordered collection of cyclic graphs whose nodes are labeled by integers. More precisely, an element of size  $n$  of a labeled structure can be decomposed in  $n$  “atomic” elements of size 1 and these atoms are labeled by distinct elements of  $\{1, \dots, n\}$ .

Operation on labeled structures are based on a special product, the *labeled* (or *partionnal*) *product* that distributes labels between components. This operation is a natural analogue of the Cartesian product for plain unlabeled structures. The labeled product in turn leads to labeled analogues of the sequence, set, and cycle constructions.

Let us define the labeled product  $\mathcal{A} = \mathcal{B} \bowtie \mathcal{C}$  of two classes  $\mathcal{B}$  and  $\mathcal{C}$ . The ordered pair  $(\beta, \gamma)$ , for  $\beta \in \mathcal{B}$  and  $\gamma \in \mathcal{C}$ , is not a labeled structure since atoms of  $\beta$ , respectively  $\gamma$ , have labels in  $\{1, \dots, |\beta|\}$ , respectively  $\{1, \dots, |\gamma|\}$ , leading to atoms with common labels. A natural lift of these two labelings, is a labeling with labels in  $\{1, \dots, |\beta| + |\gamma|\}$  such that the order relation between labels of  $\beta$ , respectively  $\gamma$ , are preserved. These labeled structures are the elements of the labeled product. For instance, consider the class of chains which are total orderings of the elements of  $\{1, \dots, k\}$  for all integers  $k$ . The pair consisting of the two chains  $(2, 1)$  and  $(1)$  is not a labeled structure:  $((2, 1), (1))$  has two atoms labeled 1. On the other hand, three natural expansions lead to labeled structures:  $((2, 1), (3))$ ,  $((3, 1), (2))$ , and  $((3, 2), (1))$ .

Any element of  $\mathcal{A}$  has a unique decomposition into elements of  $\mathcal{B} \times \mathcal{C}$ . But conversely, the pair of an element of  $\mathcal{B}$  of size  $k$  and an element of  $\mathcal{C}$  of size  $l$ , is the decomposition of as many elements as there are possibilities to label  $(\beta, \gamma)$  by  $\{1, \dots, l + k\}$  in a way that preserves the labeling induced on  $\beta$  and  $\gamma$ . So there are  $\frac{(k+l)!}{k!l!}$  such decompositions. As regards the counting sequence, an element of size  $n$  of  $\mathcal{A}$  decomposes into a pair of elements of size  $k$  and  $l$  such that  $k + l = n$ , so that

$$A_n = \sum_{k+l=n} \frac{n!}{k!l!} B_k C_l.$$

This equation can be rewritten as

$$\frac{A_n}{n!} = \sum_{k+l=n} \frac{B_k}{k!} \frac{C_l}{l!}.$$

Construction	Unlabeled structures	Labeled structures
Product	$\text{ogf}(\mathcal{B}) \cdot \text{ogf}(\mathcal{C})$	$\text{egf}(\mathcal{B}) \cdot \text{egf}(\mathcal{C})$
Union	$\text{ogf}(\mathcal{B}) + \text{ogf}(\mathcal{C})$	$\text{egf}(\mathcal{B}) + \text{egf}(\mathcal{C})$
Sequence	$\frac{1}{1 - \text{ogf}(\mathcal{B})(z)}$	$\frac{1}{1 - \text{egf}(\mathcal{B})(z)}$
Set	$\exp\left(\sum_{k=1}^{\infty} \frac{(-1)^{k+1}}{k} \text{ogf}(\mathcal{B})(z^k)\right)$	$\exp(\text{egf}(\mathcal{B})(z))$
Cycle	$\sum_{k=1}^{\infty} \frac{\phi(k)}{k} \log\left(\frac{1}{1 - \text{ogf}(\mathcal{B})(z^k)}\right)$	$\log \frac{1}{1 - \text{egf}(\mathcal{B})(z)}$

TABLE 1. Admissible constructions and generating functions interpretations.

The use of exponential generating functions to encode the counting sequences is then natural because the previous equation characterizes the product of two such functions. So the counting sequence  $\{A_n\}_{n \geq 0}$  is represented by  $A(z) = \text{egf}(\mathcal{A})(z) = \sum_{n \geq 0} \frac{A_n}{n!}$ , which was chosen such that

$$\text{egf}(\mathcal{B} \bowtie \mathcal{C}) = \text{egf}(\mathcal{B}) \cdot \text{egf}(\mathcal{C}).$$

The same work as for unlabeled structures leads to the results summarized in Table 1.

## 2. Complex Asymptotic Analysis

Once a function encoding the counting sequence has been determined, it remains to extract the sequence from the function. The explicit expansion of the function is often too difficult. To avoid it, the crucial observation is that most of the generating functions that occur in combinatorial enumerations are also *analytic functions*: their expansions converge in a neighborhood of the origin and Cauchy's integral formula expresses Taylor coefficients of such analytic functions as contour integrals.

This section is dedicated to a short presentation of analytic functions, then to the determination of the exponential growth of the counting sequence, and finally to the subexponential factors.

**2.1. Residue theorem.** A function  $f(z)$  of the complex variable  $z$  is *analytic* at a point  $z = a$  if it is defined in a neighborhood of  $z = a$  and is given there by a convergent power series expansion

$$f(z) = \sum_{n \geq 0} f_n (z - a)^n.$$

The quotient of two analytic functions  $f(z)/g(z)$  gives the intuition of what is a meromorphic function. More precisely,  $h(z)$  is *meromorphic* at  $z = a$  if and only if in a neighborhood of  $z = a$  it is given by an expansion of the form

$$h(z) = \sum_{n \geq -M} h_n (z - a)^n \quad \text{for } z \neq a.$$

If  $M \geq 1$  and  $h_{-M} \neq 0$  then  $h(z)$  is said to have a *pole* of order  $M$  at  $z = a$ . When  $h(z)$  has a pole of order  $M \geq 1$  at  $z = a$ , then the coefficient  $h_{-1}$  is called the *residue* of  $h(z)$  at  $z = a$  and it is designated by

$$\text{Res}[h(z); z = a].$$

The important *residue theorem* relates global properties of a meromorphic function (its integral along curves) to its local properties at designated points, the poles.

**Theorem 1** (Cauchy's residue theorem). *Let  $\Gamma$  be a simple closed curve oriented positively and situated inside a simply connected region  $D$  (like a disk), and assume  $g(z)$  to be meromorphic in  $D$  and analytic on  $\Gamma$ . Then*

$$\frac{1}{2i\pi} \int_{\Gamma} g(z) dz = \sum_s \operatorname{Res}[g(z); z = s],$$

where the sum is extended to all poles of  $g(z)$  enclosed in  $\Gamma$ .

A direct application of the residue theorem concerns coefficients of analytic functions.

**Theorem 2** (Cauchy's coefficient formula). *Let  $f(z)$  be analytic in a simply connected region  $D$  and let  $\Gamma$  be a closed curve oriented positively and located inside  $D$  that simply encircles the origin. Then the coefficient  $[z^n]f(z)$  admits the integral representation*

$$f_n \equiv [z^n]f(z) = \frac{1}{2i\pi} \int_{\Gamma} f(z) \frac{dz}{z^{n+1}}.$$

**2.2. Singularities and exponential rate.** Most of the counting sequences encoded by functions have an asymptotic behavior that can be described by  $A_n \sim G^n \theta(n)$  where  $\theta(n)$  is a subexponential function: the real number  $G = \limsup_{n \rightarrow +\infty} |f_n|^{1/n}$  is called the *exponential rate of growth* of the counting sequence.

This parameter has a straightforward interpretation as regards the function which encodes the counting sequence. A *singularity* of such a function can be informally defined as a point where the function ceases to be analytic. Singularities of smallest modulus of a function analytic at 0 are called *dominant singularities*. The exponential rate of growth is linked to the modulus of dominant singularities by the following theorem.

**Theorem 3** (Exponential growth formula). *If  $f(z)$  is analytic at 0 and  $R$  is the modulus of a singularity of  $f(z)$  nearest to the origin, then the exponential rate of growth of the coefficients  $[z^n]f(z)$  is  $1/R$ .*

*Proof.* Cauchy's formula applied to a circle  $\Gamma$  of center 0 and radius  $R' < R$  gives

$$|f_n| = \left| \frac{1}{2i\pi} \int_{\Gamma} f(z) \frac{dz}{z^{n+1}} \right| \leq \frac{|2\pi R'|}{|2i\pi|} \sup \{ |f(z)| \mid |z| = R' \} R'^{-(n+1)} = \mathcal{O}(R'^{-n}),$$

so that  $G = \limsup_n |f_n|^{1/n} \leq \frac{1}{R'}$ , and  $G \leq \frac{1}{R}$  by letting  $R'$  approach  $R$ .

We now assume  $G < \frac{1}{R}$  and proceed to get a contradiction, proving  $G = \frac{1}{R}$  in this way. Fix  $R'$  such that  $G < \frac{1}{R'} < \frac{1}{R}$ . For some constant  $K$  and all sufficiently large  $n$ , we have  $|f_n| \leq \frac{K}{R'^n}$ . The series  $\sum_{n \geq 0} f_n z^n$  therefore converges normally on the set of all  $z$  of modulus  $R$ , since  $0 < \frac{R}{R'} < 1$ . This contradicts the existence of a singularity of modulus  $R$ .  $\square$

An additional property of functions defined by counting sequences is that their coefficients are non-negative. This situation allows to locate one dominant singularity more precisely.

**Theorem 4** (Pringsheim's theorem). *If a function has Taylor coefficients that are real non-negative, then one of its dominant singularities, if there is a singularity, is real positive.*

**2.3. Subexponential approximation.** If the location of the singularities of a function determines the exponential rate of growth of its coefficients, the nature of the singularities determines the way the dominant exponential term in coefficients is modulated by a subexponential factor.

For sake of simplicity, we assume that the singularities are isolated. By change of the variable, we can assume that all the dominant singularities are of modulus 1. Moreover we assume that there is a unique dominant singularity which is 1.

The notion of  $\Delta$ -analytic function is defined to describe the scope of the following transfer theorem which maps the local behavior of the function around its dominant singularity to the asymptotic form of its coefficients. Given two numbers  $\phi, R$ , with  $R > 1$  and  $0 < \phi < \frac{\pi}{2}$ , the open domain  $\Delta(\phi, R)$  is defined as

$$\Delta(\phi, R) = \left\{ z \mid |z| < R, z \neq 1, |\operatorname{Arg}(z - 1)| > \phi \right\}.$$

A domain is a  $\Delta$ -domain if it is a  $\Delta(\phi, R)$  for some  $R$  and some  $\phi$ . A function is  $\Delta$ -analytic if it is analytic in some  $\Delta$ -domain.

**Theorem 5** (Big-oh transfer [1]). *Let  $\alpha$  be a number not in  $\{0, -1, -2, \dots\}$ . Assume that  $f(z)$  is  $\Delta$ -analytic and that it satisfies in the intersection of a neighbourhood of 1 and of its  $\Delta$ -domain the condition*

$$f(z) = \mathcal{O} \left( (1 - z)^{-\alpha} \left( \log \frac{1}{1 - z} \right)^\beta \right).$$

Then

$$[z^n] f(z) = \mathcal{O} \left( n^{\alpha-1} (\log n)^\beta \right).$$

*Proof.* The starting point is Cauchy’s coefficient formula. We apply it to a particular loop around the origin which is internal to the  $\Delta$ -domain of  $f$ : we choose the positively oriented contour  $\gamma_n \equiv \gamma = \gamma_1 + \gamma_2 + \gamma_3 + \gamma_4$ , with

$$\begin{cases} \gamma_1 = \left\{ z \mid |z - 1| = \frac{1}{n}, |\operatorname{Arg}(z - 1)| \geq \theta \right\} \\ \gamma_2 = \left\{ z \mid \frac{1}{n} \leq |z - 1|, |z| \leq r, \operatorname{Arg}(z - 1) = \theta \right\} \\ \gamma_3 = \left\{ z \mid |z - 1| = r, |\operatorname{Arg}(z - 1)| \geq \theta \right\} \\ \gamma_4 = \left\{ z \mid \frac{1}{n} \leq |z - 1|, |z| \leq r, \operatorname{Arg}(z - 1) = -\theta \right\} \end{cases}$$

If the  $\Delta$ -domain of  $f$  is  $\Delta(\phi, R)$ , we assume that  $1 < r < R$ , and  $\phi < \theta < \frac{\pi}{2}$ , so that the contour  $\gamma$  lies entirely inside the domain of analyticity of  $f$ .

For  $j = 1, 2, 3, 4$ , let

$$f_n^{(j)} = \frac{1}{2i\pi} \int_{\gamma_j} f(z) \frac{dz}{z^{n+1}}.$$

The analysis proceeds by bounding the absolute value of the integral along each of the four parts. In order to keep notations simple, we detail the proof in the case where  $\beta = 0$ .

*Inner circle.* From trivial bounds, the contribution there is

$$\left| f_n^{(1)} \right| = \mathcal{O} \left( \left( \frac{1}{n} \right)^{1-\alpha} \right),$$

as the function  $f$  is  $\mathcal{O} \left( \left( \frac{1}{n} \right)^{-\alpha} \right)$ , the contour has length  $\mathcal{O} \left( \frac{1}{n} \right)$ , and  $z^{-n-1}$  is  $\mathcal{O}(1)$  there.

*Rectilinear parts.* Setting  $\omega = e^{i\theta}$  and performing the change of variable  $z = 1 + \frac{\omega t}{n}$ , we find

$$\left| f_n^{(2)} \right| < \frac{1}{2\pi} \int_1^\infty K \left( \frac{t}{n} \right)^{-\alpha} \left| 1 + \frac{\omega t}{n} \right|^{-n-1} dt,$$

for some constant  $K > 0$  such that  $|f(z)| < K(1 - z)^{-\alpha}$  “over the  $\Delta$ -domain.” In fact we have a constant for a small neighborhood  $V$  of 1 due to the asymptotic assumption and an other constant

Function $f(z)$	Asymptotic expansion of the coefficients $f_n$
1	0
$(1 - z)^{-1}$	1
$(1 - z)^{-2}$	$n + 1$
$(1 - z)^{-3}$	$\frac{1}{2}n^2 + \frac{3}{2}n + 1$
$(1 - z)^{1/2}$	$-\frac{1}{\sqrt{\pi n^3}} \left( \frac{1}{2} + \frac{3}{16n} + \frac{25}{256n^2} + \mathcal{O}\left(\frac{1}{n^3}\right) \right)$
$(1 - z)^{-1/2}$	$\frac{1}{\sqrt{\pi n}} \left( 1 - \frac{1}{8n} + \frac{1}{128n^2} + \frac{5}{1024n^3} + \mathcal{O}\left(\frac{1}{n^4}\right) \right)$
$\log(1 - z)^{-1}$	$\frac{1}{n}$
$(1 - z)^{-3/2} \log(1 - z)^{-1}$	$\sqrt{\frac{n}{\pi}} \left( 2 \log n + 2\gamma + 4 \log 2 - 2 + \frac{3 \log n}{4n} + \mathcal{O}\left(\frac{1}{n}\right) \right)$

TABLE 2. Examples of applications of the transfer theorem.

that comes from the compacity of a closed set  $C$  included in  $\Delta$  such that all the used loops are in  $C \cup V$ . From the relation

$$\left| 1 + \frac{\omega t}{n} \right| \geq 1 + \frac{t}{n} \cos \theta,$$

there results

$$\left| f_n^{(2)} \right| < \frac{K}{2\pi} J_n n^{\alpha-1}$$

where

$$J_n = \int_1^\infty t^{-\alpha} \left( 1 + \frac{t \cos \theta}{n} \right)^{-n} dt.$$

For a given  $\alpha$ , the integrals  $J_n$  are all bounded above by some constant since they admit a limit as  $n$  tends to infinity:

$$J_n \rightarrow \int_1^\infty t^{-\alpha} e^{-t \cos \theta} dt.$$

(The condition on  $\theta$  that  $0 < \theta < \frac{\pi}{2}$  precisely ensures convergence of the integral.) Thus, globally, on this part of the contour, we have

$$\left| f_n^{(2)} \right| = \mathcal{O}(n^{\alpha-1}),$$

and the same bound holds for  $f_n^{(4)}$  by symmetry.

*Outer circle.* There,  $f(z)$  is bounded while  $z^{-n}$  is of the order  $r^{-n}$ . Thus,  $f_n^{(3)}$  is exponentially small.

In summary, each of the four integrals of the split contour contributes  $\mathcal{O}(n^{\alpha-1})$ . The statement of the theorem thus follows. □

This theorem can be extended to equivalents giving a fairly mechanical process to translate asymptotic information on a function into information on its coefficients. These are simple functions that are used as a scale since any function equivalent to it around its dominant singularity as the same asymptotic expansion. See Table 2.

### 3. Examples

3.1. **Clouds.** Let us consider  $n$  lines in general position in the plane. A *cloud* is a subset of the set of the intersection points of the lines such that:

1. any three points of the cloud are not aligned;
2. any line has at least one of its points in the cloud;
3. the set is maximal for inclusion among the sets that satisfies points 1 and 2.

The size of a cloud is the number of points it contains.

There is a more combinatorial description of a cloud since they are in bijection with labeled 2-regular graphs (any vertex has degree 2, no loops, no multiple edges). In the bijection, the line labeled  $i$  is the vertex labeled  $i$  of the graph and the intersection between the line  $i$  and  $j$  is mapped to an edge between  $i$  and  $j$ . Indeed, point 1 in the definition exactly means that any vertex of the graph has degree at most 2 because three aligned intersections are necessarily on a common line since the picture is as general as possible. Point 2 translates the fact that any vertex has degree at least 1. Assume there are at least two vertices  $i, j$  of degree 1 in the cloud  $S$ . Then  $S \cup \{(ij)\}$  is a cloud and that is in contradiction with point 3. Finally, there cannot be only one vertex of degree 1 since the sum of the degree of vertices of a graph is even (each edge appears twice). As regards the size, since there are two intersections per line in a cloud and that an intersection is shared by two lines, the size of the cloud is the number of vertices. Thus instead of clouds we could equivalently consider the class of labeled 2-regular graphs where the size of an element is its number of vertices.

A labeled 2-regular graph is a set of non-oriented cycles of size at least 3 and we are interested in the exponential generating function of this structure. Oriented cycles of size at least 3 are the oriented cycles that do not contain 1 or 2 elements only so their generating function is

$$C_+^{>2}(z) = \log \frac{1}{1-z} - \left( \frac{1}{1!}z + \frac{1}{2!}z^2 \right).$$

A non-oriented cycle of at least 3 vertices admits exactly 2 distinct orientations, so that the generating function of non-oriented cycle of at least 3 vertices is

$$C^{>2}(z) = \frac{1}{2}C_+^{>2}(z).$$

Then the series of the sets of non-oriented cycles on at least 3 vertices and equivalently of the clouds is

$$\text{Clouds}(z) = \exp C^{>2}(z) = \frac{\exp\left(-\frac{1}{2}z - \frac{1}{4}z^2\right)}{\sqrt{1-z}}.$$

Thus,  $\text{Clouds}(z)$  is the product of  $1/\sqrt{1-z}$  which admits 1 as singularity of minimal modulus and is analytic in  $\mathbb{C} \setminus [1, +\infty)$ , and  $\exp\left(-\frac{1}{2}z - \frac{1}{4}z^2\right)$  that is entire. The behavior of  $\text{Clouds}(z)$  around 1 is the product of  $1/\sqrt{1-z}$  and  $\exp(-3/4)\left(1 + (1-z) + \frac{1}{4}(1-z)^2 - \frac{1}{12}(1-z)^3 + \mathcal{O}((1-z)^4)\right)$ , the standard Taylor expansion at 1 of  $\exp\left(-\frac{1}{2}z - \frac{1}{4}z^2\right)$ .

$$\text{Clouds}(z) = \frac{e^{-3/4}}{\sqrt{1-z}} + e^{-3/4}\sqrt{1-z} + \frac{e^{-3/4}(1-z)^{3/2}}{4} - \frac{e^{-3/4}(1-z)^{5/2}}{12} + \dots$$

This expansion is valid in a  $\Delta$ -domain so that by the principle of singularity analysis, the asymptotic determination of the coefficients  $c_n = [z^n] \text{Clouds}(z)$  results from a direct translation of the expansion

$$\text{Clouds}(z) = e^{-3/4} \frac{1}{\sqrt{1-z}} + e^{3/4} \sqrt{1-z} + \mathcal{O}\left((1-z)^{3/2}\right)$$



into

$$\begin{aligned} c_n &= e^{-3/4} \binom{n-1/2}{-1/2} + e^{-3/4} \binom{n-3/2}{-3/2} + \mathcal{O}\left(\frac{1}{n^{5/2}}\right) \\ &= \frac{e^{-3/4}}{\sqrt{\pi n}} \left(1 - \frac{1}{8n} + \frac{1}{128n^2} + \dots\right) - \frac{e^{-3/4}}{2\sqrt{\pi n^3}} \left(1 + \frac{3}{8n} + \dots\right) + \mathcal{O}\left(\frac{1}{n^{5/2}}\right). \end{aligned}$$

We finally have the asymptotic behavior of the counting sequence  $\{C_n\}_{n \geq 0}$  of clouds,

$$\frac{C_n}{n!} = c_n = \frac{e^{-3/4}}{\sqrt{\pi n}} + \frac{3e^{-3/4}}{8\sqrt{\pi n^3}} + \mathcal{O}\left(\frac{1}{n^{5/2}}\right) \quad \text{as } n \rightarrow +\infty.$$

**3.2.  $\Omega$ -trees.** A subset  $\Omega$  of  $\mathbb{N}$  is aperiodic if the greatest common divisor of its elements is 1. Given an aperiodic finite set  $\Omega$ , the class  $\mathcal{T}_\Omega$  of  $\Omega$ -trees is the set of rooted trees with a total order on the children of each node such that the degree of each node is in  $\Omega$ . For instance binary trees are  $\{0, 1, 2\}$ -trees. The size of an  $\Omega$ -tree is its number of nodes. This class is well defined if  $0 \in \Omega$  otherwise there are no finite  $\Omega$ -trees.

Since a  $\Omega$ -tree is made of a root and a sequence of length  $i \in \Omega$  of  $\Omega$ -trees, its ordinary generating function  $T$  satisfies

$$T(z) = z \cdot \sum_{\omega \in \Omega} (T(z))^\omega.$$

Let  $P(X)$  be the polynomial  $\sum_{\omega \in \Omega} X^\omega$ . The equation becomes  $T(z) = zP(T(z))$ . To check if the function  $T$  is analytic at  $z$  we rephrase the above equation as

$$z = T(z)/P(T(z)) = \psi(T(z))$$

so that it is a generic instance of the inversion problem for analytic functions ( $\psi(u) = \frac{u}{P(u)}$ ).

An important statement of the inversion theorem is that if  $\psi$  is analytic at  $t = t_0$ , then  $T(z)$  is analytic at  $z = \psi(t_0)$  if and only if  $\psi'(t_0) \neq 0$ . To have an intuition of this result, consider the analytic expansion of  $\psi$  near  $t_0$ :

$$\psi(t) = \psi(t_0) + (t - t_0)\psi'(t_0) + \frac{1}{2}(t - t_0)^2\psi''(t_0) + \dots.$$

If  $\psi'(t_0) \neq 0$ , solving formally for  $t$  suggests that  $t - t_0 \sim \frac{1}{\psi'(t_0)}(z - \psi(t_0))$  and a full expansion is obtained by repeated substitutions. If on the contrary  $\psi'(t_0) = 0$  and  $\psi''(t_0) \neq 0$ , solving formally now suggest that  $(t - t_0)^2 \sim \frac{2}{\psi''(t_0)}(z - \psi(t_0))$  so that the inversion problem should admit two solutions satisfying

$$t - t_0 \sim \pm \sqrt{\frac{2}{-\psi''(t_0)}}(\psi(t_0) - z)^{1/2}.$$

In this case the point  $\psi(t_0)$  is a branch point, so that  $T(z)$  cannot be analytic at this point. If the first nonzero derivative of  $\psi$  at  $t_0$  is of order  $r \geq 2$ , the same remark holds with a local behavior for  $t$  then of the form  $(\psi(t_0) - z)^{1/r}$ .

Because of Pringsheim's theorem, if  $T$  has a finite radius, then there is a dominant singularity in  $[0, +\infty)$ . Thus finding a dominant singularity of  $T$  results in searching the smallest positive zero of  $\psi'$ . Let  $\rho$  be this minimal zero of  $\psi'(x) = \frac{P(x) - xP'(x)}{P(x)^2}$ . This number satisfies

$$P(\rho) - \rho P'(\rho) = 0.$$

Now we have to check the number of distinct dominant singularities. By definition a dominant singularity can be written as  $\lambda = \rho e^{i\theta}$  and satisfies  $\psi'(\lambda) = 0$ . Assume there is an integer  $k \geq 2$

such that all  $\omega \in \Omega$  is divided by  $k$ . In this case  $P(x) - xP'(x) = \sum_{\omega \in \Omega} (1 - \omega)x^\omega$  can be rewritten in  $\sum_{\omega \in \Omega} (1 - \omega)(x^k)^{\omega/k}$ . Thus if  $\lambda^k = \rho^k$  then  $\lambda$  is an other dominant singularity so all complexes  $(\rho e^{2i\pi l/k})_{0 \leq l \leq k-1}$  are distinct dominant singularities. To apply the transfer theorems presented in the previous section safely we have to ensure that there is a unique dominant singularity,<sup>2</sup> therefore we made the assumption that the set  $\Omega$  is aperiodic. We admit that this condition is sufficient to have a unique dominant singularity  $\rho$  (there is a proof using the case of equality in the triangular inequality).

Since  $\rho$  satisfies  $P(\rho) - \rho P'(\rho) = 0$ , we have  $\psi''(\rho) = \frac{-\rho^2 P''(\rho)}{P(\rho)^3}$ . Thus if  $\Omega$  contains an element greater than 1,  $\Psi''(\rho) > 0$  and

$$T(z) = T(\rho) \pm \sqrt{\frac{2P(\rho)^3}{\rho^2 P''(\rho)}} \sqrt{\rho} \sqrt{1 - \frac{z}{\rho}} + \mathcal{O}\left(\left(1 - \frac{z}{\rho}\right)^{3/2}\right).$$

This expansion is valid on a  $\Delta$ -domain; thus using a transfer theorem, we obtain the asymptotic equivalent

$$[z^n] T(z) \sim \sqrt{\frac{2}{\rho^2 \Psi''(\rho)}} \sqrt{\rho} \frac{1}{2\sqrt{\pi n^3}} \rho^{-n} = C_\rho \rho^{-n} n^{-3/2}.$$

### Bibliography

- [1] Flajolet (Philippe) and Odlyzko (Andrew). – Singularity analysis of generating functions. *SIAM Journal on Discrete Mathematics*, vol. 3, n° 2, 1990, pp. 216–240.
- [2] Flajolet (Philippe) and Sedgewick (Robert). – *The Average Case Analysis of Algorithms: Complex Asymptotics and Generating Functions*. – Research Report n° 2026, Institut National de Recherche en Informatique et en Automatique, 1993. 100 pages.
- [3] Flajolet (Philippe) and Sedgewick (Robert). – *The Average Case Analysis of Algorithms: Counting and Generating Functions*. – Research Report n° 1888, Institut National de Recherche en Informatique et en Automatique, 1993. 116 pages.

---

<sup>2</sup>A generalisation of transfer theorems exists for the case of multiple of singularities; see [2, p. 85].

## Aléa discret et mouvement brownien<sup>†</sup>

*Philippe Chassaing*

Institut Élie Cartan, Université Nancy 1 (France)

March 26 and 27, 2001

*Summary by Philippe Chassaing*

Les liens entre le mouvement brownien et ses processus dérivés (méandre, pont, excursion) d'une part et d'autre part des objets combinatoires comme les mots de Dyck, les permutations bi-ordonnées, le tri *Shell's sort*, les arbres simples, les facteurs gauches, le hachage ou parking, les animaux dirigés, le graphe aléatoire, la marche aléatoire dans le plan fendu, . . . , rendent opportune une revue (forcément partielle) des innombrables propriétés du mouvement brownien.

En combinatoire et analyse d'algorithmes, beaucoup d'asymptotiques de statistiques intéressantes sont familières aux spécialistes du mouvement brownien : la hauteur ou la largeur des arbres simples normalisées convergent en loi vers une loi liée à la fonction  $\theta$  de Jacobi, connue pour être la loi du maximum de l'excursion brownienne. Dans l'asymptotique des nombres de Wright, dénombrant les graphes connexes à  $n$  sommets et  $k$  arêtes en excès [19], apparaissent les moments de la surface sous l'excursion brownienne, dont la distribution s'exprime à l'aide de la fonction d'Airy<sup>1</sup>. Le profil moyen d'un arbre simple suit asymptotiquement la loi de Rayleigh, qui est la loi du maximum du pont brownien. Le déplacement total dans une table de hachage pleine, est également asymptotiquement distribuée selon une loi d'Airy. Il est tentant de voir ces faits comme les fragments d'un même tableau : la convergence des chemins de Bernoulli (resp. de Dyck) et d'objets analogues vers le mouvement brownien (resp. l'excursion brownienne). Une version arbre en est donnée par Aldous avec sa convergence des arbres simples vers le *continuum random tree*.

À cette première explication de l'omniprésence de certaines lois vient s'ajouter le principe d'invariance [7]<sup>2</sup> selon lequel la loi limite de différentes fonctionnelles d'une marche aléatoire ne dépend que très peu (à un facteur multiplicatif près) de la loi d'un pas élémentaire : ce dernier principe se traduit, par exemple, en informatique fondamentale, par l'apparition de la même loi limite pour la hauteur de différents arbres simples [8, 16], ou encore de la même loi limite pour le cheminement total d'un arbre binaire ou pour le déplacement total d'une table de hachage pleine.

Pour beaucoup d'autres situations combinatoires (tailles de composantes connexes du graphes aléatoires, *minimum spanning tree*, *random mappings*, cartes planaires, etc.), l'existence d'un objet aléatoire limite est soupçonnée ou avérée, expliquant ainsi les lois limites déjà observées, fournissant éventuellement de nouveaux résultats asymptotiques en combinatoire et en analyse d'algorithmes, et posant de nouvelles questions sur l'omniprésent mouvement brownien. Il est sage pour un mini-cours de se limiter à la convergence d'objets combinatoires très basiques : chemins de Dyck (bilatères ou non) et facteurs gauches, tous étant plus généralement des chemins de Bernoulli, vers le mouvement brownien et ses avatars, excursion brownienne, méandre et pont. Le mouvement

---

<sup>†</sup>Notes de cours pour le cours donné pendant le groupe de travail ALÉA'01 à Luminy (France).

<sup>1</sup>Pour un aperçu agréable du lien entre mouvement brownien et fonctions spéciales, voir [1].

<sup>2</sup>cf. [5], lire l'introduction.

brownien, ses propriétés, et le théorème de Donsker requièrent une trentaine d'heures de cours pour un traitement rigoureux ; j'éviterai donc les démonstrations, et renverrai largement à la bibliographie abondante sur le sujet, en particulier à [17, 2, 12].

*Plan.*

1. Différents types de chemins aléatoires
2. Changement d'échelle brownien (*Brownian scaling*) et convergence faible
3. Convergence faible : définition et premières conséquences
4. Convergence faible : critères et autres caractérisations
5. Propriétés du mouvement brownien
6. Décompositions remarquables des trajectoires du mouvement brownien
7. Diverses propriétés de l'excursion brownienne normalisée, du pont et du méandre brownien
8. Conclusion

Les sections 6 à 8 seront rédigées dans un document ultérieur.

### 1. Différents types de chemins aléatoires

**Définition** (Chemins de Bernoulli). Un *chemin de Bernoulli* est un chemin sur le réseau engendré par  $NE = (1, 1)$  et  $SE = (1, -1)$ , partant de  $(0, 0)$ , admettant comme pas élémentaires précisément les pas  $NE$  et  $SE$ . Il y a  $2^n$  chemins de Bernoulli de longueur  $n$ .

**Définition** (Chemins de Dyck). Un *chemin de Dyck* de longueur  $2n$  est un chemin de Bernoulli de longueur  $2n$  qui se termine au point  $(2n, 0)$  et reste positif ou nul sur toute sa longueur. Il y a

$$C_n = \frac{\binom{2n+1}{n}}{2n+1}$$

chemins de Dyck de longueur  $2n$ . Un mot de Dyck est la description d'un chemin de Dyck par la suite de ses pas, *i. e.* un mot formé d'autant de caractères 'M' (pour « montées ») que de caractères 'D' (pour « descentes »), et dont n'importe quel préfixe contient au moins autant de 'M' que de 'D'. Il y a une bijection privilégiée (entre mots et chemins), alors notons indifféremment  $\mathcal{B}_{2n}^\oplus$  l'ensemble des  $C_n$  chemins de Dyck de longueur  $2n$  ou l'ensemble des  $C_n$  mots de Dyck de longueur  $2n$ .

**Définition** (Chemins de Dyck bilatères). Un *chemin de Dyck bilatère* de longueur  $2n$  est un chemin de Bernoulli de longueur  $2n$  qui se termine au point  $(2n, 0)$ . Il y a  $\binom{2n}{n}$  chemins de Dyck bilatères de longueur  $2n$ .

**Définition** (Facteurs gauches). Un *facteur gauche* de longueur  $n$  est un chemin de Bernoulli de longueur  $n$  qui reste positif ou nul tout au long de sa trajectoire. Il y a  $\binom{n}{\lfloor n/2 \rfloor}$  facteurs gauches de longueur  $n$ .

**Variables aléatoires correspondantes.** Quitte à identifier une fonction et son graphe, on peut voir l'ensemble  $\mathcal{B}_n$  des chemins de Bernoulli de longueur  $n$  et ses sous ensembles  $\mathcal{B}_n^\oplus$  (ensemble des chemins de Dyck<sup>3</sup>),  $\mathcal{B}_n^o$  (ensemble des chemins de Dyck bilatères) et  $\mathcal{B}_n^+$  (ensemble des facteurs gauches) comme des parties finies de l'espace  $\mathcal{C}[0, n]$  des fonctions continues. On notera  $\nu_n$  (resp.  $\nu_n^\oplus$ ,  $\nu_n^o$ ,  $\nu_n^+$ ) la mesure de probabilité sur  $\mathcal{C}[0, n]$  uniforme sur  $\mathcal{B}_n$  (resp.  $\mathcal{B}_n^\oplus$ ,  $\mathcal{B}_n^o$ ,  $\mathcal{B}_n^+$ ).

---

<sup>3</sup>Dans la suite, chaque fois que c'est nécessaire, pour les chemins de Dyck p. e., on sous entendra que  $n$  est pair, et dans ce cas on notera  $n = 2n'$ .

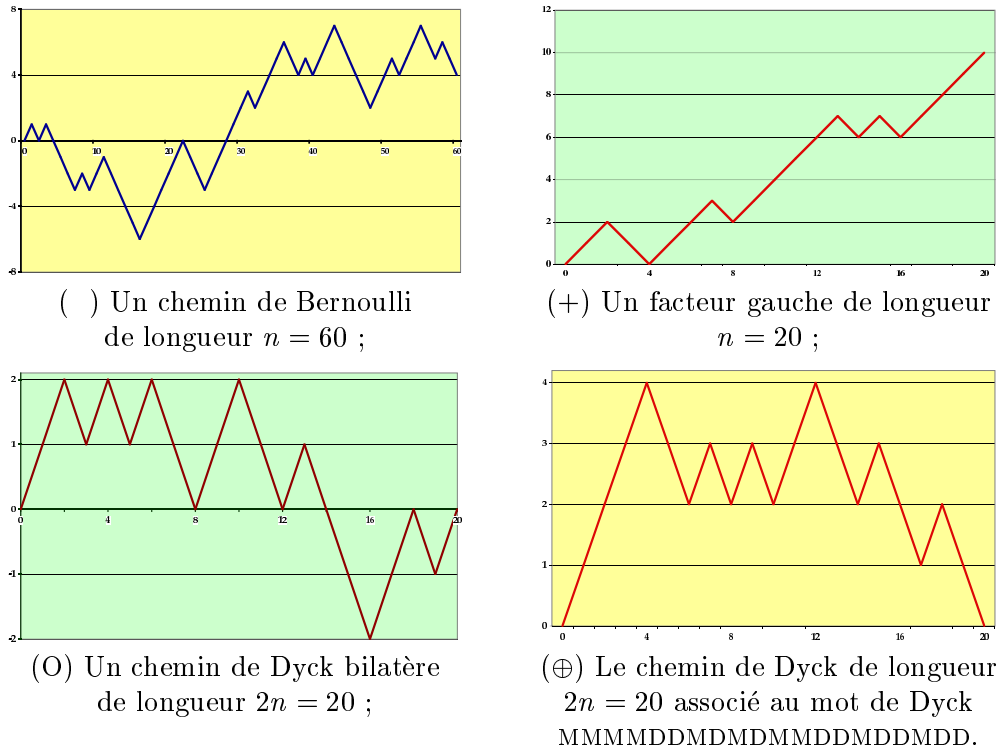


FIGURE 1. Différents types de chemins.

**Définition.** Dans la suite, une variable aléatoire de loi  $\nu_n$  (resp.  $\nu_n^\oplus, \nu_n^o, \nu_n^+$ ) sera appelée *marche aléatoire simple symétrique* (resp. *excursion de Bernoulli, pont de Bernoulli, méandre de Bernoulli*) de longueur  $n$ .

Une variable aléatoire  $X$  à valeur dans un espace de fonctions, p. e. dans  $\mathcal{C}[0, 1], \mathcal{C}[0, n]$  ou encore  $\mathcal{C}[0, +\infty)$ , est souvent appelée *processus stochastique*.

Seule l'appellation « marche aléatoire simple symétrique » est bien établie, les 3 autres étant inspirées d'un vocabulaire bien établi dans le cadre du mouvement brownien, où l'on parle d'*excursion brownienne*, de *pont brownien*, et de *méandre brownien*. Dans la suite, par un abus de langage sur lequel on ne s'attardera pas, on identifiera couramment une suite  $u = (u_k)_{0 \leq k \leq n}$  à son prolongement en une fonction  $f$  continue linéaire par morceaux sur  $[0, n]$ , ou encore au graphe de cette dernière fonction. En particulier, les fonctions de  $\mathcal{B}_n$  sont bien définies par leurs évaluations en  $0, 1, 2, \dots, n$ . La construction usuelle d'une marche aléatoire simple symétrique est plutôt celle de la suite des  $n + 1$  évaluations :

**Définition** (Marche aléatoire simple symétrique, définition équivalente). Notons  $(Y_i)_{i \geq 1}$  une suite de variables aléatoires indépendantes et de même loi (on abrègera « indépendantes et de même loi » en i. i. d. dans la suite), avec

$$\mathbf{P}(Y_k = 1) = \mathbf{P}(Y_k = -1) = 1/2,$$

et posons

$$S_0 = 0, \quad S_k = Y_1 + Y_2 + \dots + Y_k,$$

on dit que  $S = (S_k)_{0 \leq k \leq n}$  est la *marche aléatoire simple symétrique* de longueur  $n$ .

**Remarques.**

1. On verra dans la suite que cette construction révèle certaines propriétés cruciales des chemins de Bernoulli, dont le mouvement brownien va hériter par passage à la limite.
2. Il est naturel, dans ce contexte, de définir la marche simple symétrique *pour tout entier non négatif*, *i. e.* de définir un chemin de Bernoulli aléatoire de longueur infinie.
3. Plus généralement, une marche aléatoire  $S = (S_k)_{k \geq 0}$  est définie sur un groupe  $(G, \oplus)$ , *p. e.* ici  $(\mathbb{R}, +)$ , par

$$S_k = Y_1 \oplus Y_2 \oplus \cdots \oplus Y_k,$$

les  $Y_i$  étant *i. i. d.*, la loi de probabilité commune aux  $Y_i$  étant appelée « *pas* » de la marche. On peut par exemple associer aux arbres unaires-binaires aléatoires, ou aux arbres étiquetés aléatoires, une marche aléatoire dont le pas est différent du pas de la marche aléatoire simple symétrique, *i. e.* différent de  $\frac{1}{2}\delta_{-1} + \frac{1}{2}\delta_1$ .

Une fois la marche aléatoire simple symétrique ainsi définie, on peut voir  $\nu_n^\oplus$  (resp.  $\nu_n^0, \nu_n^+$ ) comme des lois conditionnelles de cette marche de longueur  $n$ , c'est-à-dire que, pour  $A \subset \mathcal{C}[0, n]$ ,

$$\begin{aligned} \nu_n(A) &= \frac{\#(A \cap \mathcal{B}_n)}{2^n} = \mathbf{P}(S \in A), \\ \nu_n^\oplus(A) &= \frac{\#(A \cap \mathcal{B}_n^\oplus)}{C_{n'}} = \mathbf{P}(S \in A \mid S_k \geq 0, 0 \leq k \leq n \text{ et } S_n = 0), \\ \nu_n^0(A) &= \frac{\#(A \cap \mathcal{B}_n^0)}{\binom{n}{n'}} = \mathbf{P}(S \in A \mid S_n = 0), \\ \nu_n^+(A) &= \frac{\#(A \cap \mathcal{B}_n^+)}{\binom{n}{\lfloor n/2 \rfloor}} = \mathbf{P}(S \in A \mid S_k \geq 0, 0 \leq k \leq n). \end{aligned}$$

Ces définitions de  $\nu_n$  (resp.  $\nu_n^\oplus, \nu_n^0, \nu_n^+$ ) fournissent un algorithme efficace pour la génération d'un chemin de Bernoulli aléatoire, et des algorithmes de rejet parfaitement inefficaces pour la génération des chemins de Dyck (bilatères ou non) ou encore des facteurs gauches.

## 2. Changement d'échelle brownien (*Brownian scaling*) et convergence faible

**Définition** (Changement d'échelle brownien (*Brownian scaling*)). Étant donné une fonction  $f$  définie sur un intervalle  $[a, b]$  borné, on note  $f^{[a,b]}$  la fonction définie sur  $[0, 1]$  par

$$f^{[a,b]}(t) = \frac{1}{\sqrt{b-a}} f(a + t(b-a)).$$

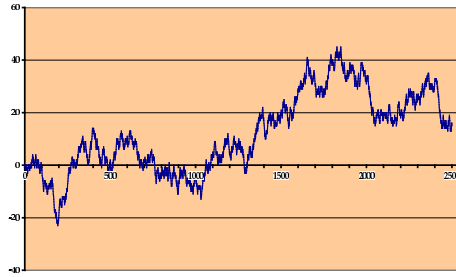
En particulier cette opération envoie bijectivement  $\mathcal{C}[a, b]$  sur  $\mathcal{C}[0, 1]$ .

Le graphe de  $f^{[a,b]}$  est ainsi obtenu, à partir de celui de  $f$ , en multipliant la largeur par un facteur  $\frac{1}{b-a}$  et la hauteur par un facteur  $\frac{1}{\sqrt{b-a}}$ . Bachelier en 1900, ou Einstein en 1905 (dans leur étude respectivement du cours des actions en bourse, et du mouvement, observé par Brown en 1826, de certaines particules en suspension dans un liquide) utilisent explicitement ou implicitement, une propriété remarquable : *le changement d'échelle brownien d'un chemin de Bernoulli de longueur  $n$  converge vers un objet limite, quand  $n$  tend vers  $+\infty$ .*

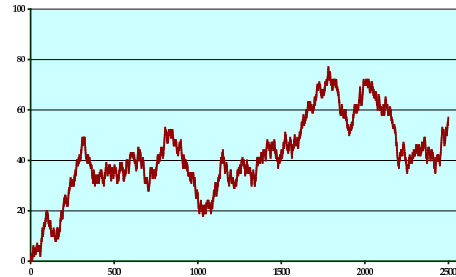
Notons  $\mu_n$  (resp.  $\mu_n^\oplus, \mu_n^0, \mu_n^+$ ) l'image de  $\nu_n$  (resp.  $\nu_n^\oplus, \nu_n^0, \nu_n^+$ ) par le changement d'échelle brownien. Le résultat clé de ce mini-cours est le

**Théorème.** *La suite de mesures de probabilités  $\mu_n$  (resp.  $\mu_n^\oplus, \mu_n^0, \mu_n^+$ ) sur l'espace  $\mathcal{C}[0, 1]$  possède, au sens de la convergence faible, une mesure de probabilité limite,  $\mu$  (resp.  $\mu^\oplus, \mu^0, \mu^+$ ).*

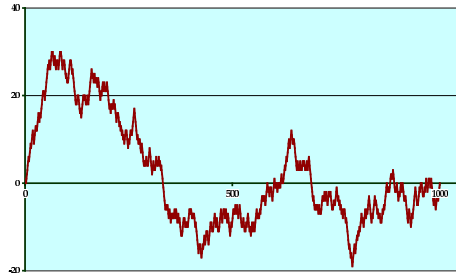
La notion de convergence faible est développée Sections 3 et 4. Fixons le vocabulaire.



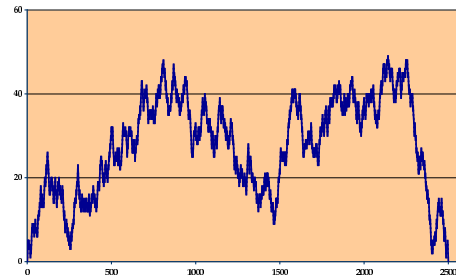
( ) Un chemin de Bernoulli de longueur  $n = 2500$  ;



(+) Un méandre de Bernoulli au hasard de longueur  $n = 2500$  ;



(O) Un chemin de Dyck bilatère au hasard de longueur  $n = 1000$  ;



(⊕) Un chemin de Dyck au hasard de longueur  $n = 2500$ .

FIGURE 2. Chemins de Bernoulli au hasard de longueur 1000 à 2500 : ils possèdent en général des fluctuations d'ordre de grandeur quelques dizaines.

**Définition** (Mouvement brownien). La mesure de probabilité  $\mu$ , définie sur  $\mathcal{C}[0, 1]$  muni de sa tribu de boréliens, est appelée *mesure de Wiener*. Une variable aléatoire  $B$  à valeur dans  $\mathcal{C}[0, 1]$ , ayant pour loi la mesure de Wiener, est appelée *mouvement brownien* (linéaire) (standard).

**Définition** (Excursion, pont et méandre browniens). Une variable aléatoire  $e$  (resp.  $b$ ,  $m$ ) à valeur dans  $\mathcal{C}[0, 1]$ , ayant pour loi la mesure  $\mu^\oplus$  (resp.  $\mu^o$ ,  $\mu^+$ ), est appelée *excursion brownienne* (normalisée) (resp. *pont brownien*, *méandre brownien*).

Les chemins de Bernoulli de la Figure 2 donnent une idée de l'allure typique du mouvement brownien ( ), resp. du méandre (+), du pont (O), de l'excursion brownienne (⊕). On peut résumer les définitions précédentes en un tableau  $2 \times 2$ , suivant la présence ou l'absence des deux contraintes (de positivité et de retour en 0 à la fin) :

### Remarques.

- Le théorème ci-dessus rassemble en fait quatre théorèmes et possède quatre auteurs : la convergence des marches aléatoires vers la mesure de Wiener  $\mu$ , ou vers le mouvement brownien, a été démontrée par Donsker [6], la convergence vers l'excursion brownienne par Kaigh [11], la convergence vers le méandre brownien par Iglehart [9], et celle vers le pont brownien par Liggett [13].
- Les résultats de Donsker, Iglehart et autres portent en fait sur la convergence de marches aléatoires, *conditionnées ou non*, de pas plus généraux que ceux de la marche aléatoire simple symétrique : les pas  $Y_i$  sont toujours i. i. d., mais de loi commune quasiment quelconque

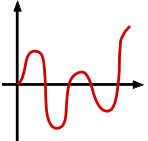
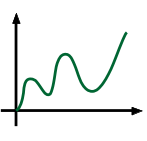
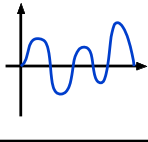
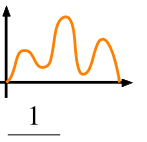
		contrainte de positivité:	
		sans	avec
contrainte de retour en 0:	sans	<p><b>mouvement Brownien B</b> chemin de Bernoulli</p> <p><math>v_n, \mu_n \xrightarrow{\text{faiblement}} \mu</math></p> <p>mesure de Wiener <math>\mu</math></p> <p><math>\# \mathcal{B}_n = 2^n</math></p> 	<p><b>+ méandre Brownien m</b> facteur gauche</p> <p><math>v_n^+, \mu_n^+ \xrightarrow{\text{faiblement}} \mu^+</math></p> <p><math>\# \mathcal{B}_n^+ = \binom{n}{\lfloor n/2 \rfloor}</math></p> 
	avec	<p><b>○ pont Brownien b</b> chemin de Dyck bilatère</p> <p><math>v_n^\circ, \mu_n^\circ \xrightarrow{\text{faiblement}} \mu^\circ</math></p> <p><math>\# \mathcal{B}_n^\circ = \binom{n}{n/2}</math></p> 	<p><b>⊕ excursion Brownienne e</b> chemin de Dyck</p> <p><math>v_n^\oplus, \mu_n^\oplus \xrightarrow{\text{faiblement}} \mu^\oplus</math></p> <p><math>\# \mathcal{B}_n^\oplus = \binom{n+1}{n/2} \frac{1}{n+1}</math></p> 

FIGURE 3. Les différents types de chemins et leurs analogues browniens.

(parfois  $Y_i$  doit être à valeurs entières, il doit toujours être centré ( $\mathbf{E}[Y_i] = 0$ ) et à variance finie ( $0 < \mathbf{E}[Y_i^2] < +\infty$ )). Cette généralité est bien sûr intéressante, mais plus particulièrement en combinatoire, ou en analyse d'algorithme. Par exemple, il est expliqué dans Aldous ou dans [19] que le mot de Lukasiewicz associé à un arbre général (resp. étiqueté) de taille  $n$  est aussi associé à une marche aléatoire de longueur  $n$ , conditionnée  $\oplus$ , de pas  $p = (p_k)_{k \geq -1}$  géométrique (donné par  $p_k = 2^{-k-2}$ ) (resp. de pas Poisson, donné par  $p_k = \frac{1}{k+1!e}$ ). Ainsi le résultat de Kaigh permet d'expliquer un faisceau de comportements asymptotiques de statistiques liées aux arbres « généraux » (resp. aux arbres étiquetés, graphes connexes et, par exemple, constantes de Wright, hachage linéaire, etc.).

- Une multitude de caractérisations et constructions différentes du mouvement brownien, du pont, de l'excursion et du méandre brownien, souvent découlant de propriétés combinatoires des chemins de Dyck ou de Bernoulli, seront données aux Sections 5 et 7.

### 3. Convergence faible : définition et premières conséquences

J'abrège encore ici ce qui est expliqué de manière très claire et assez économique dans le livre fondamental de Billingsley. On se placera dans un espace métrique  $(S, \mathcal{S})$ , qui, pour nous, sera exclusivement  $\mathbb{R}^d$  ou  $\mathcal{C}[0, 1]$ , muni de la distance usuelle dans le premier cas, de la distance de la convergence uniforme dans le second cas ;  $\mathcal{S}$  désignera la tribu engendrée par (la plus petite tribu contenant les) ouverts de la topologie induite. Les mesures considérées seront des mesures de probabilité sur  $\mathcal{S}$ . Les résultats ci-dessous s'appliquent à des espaces métriques plus généraux, dont on exige parfois qu'il soient complets et séparables (voir [2, 14]).

**Définition** (Convergence faible). On dit que la suite de mesures de probabilité  $(\mu_n)_{n \geq 0}$  converge faiblement vers la mesure de probabilité  $\mu$ , si et seulement si, pour toute fonction continue bornée  $f$



de  $S$  dans  $\mathbb{R}$ ,

$$\lim_n \int f d\mu_n = \int f d\mu.$$

On dit qu'une suite de variables aléatoires  $X_n$  à valeurs dans  $S$  converge faiblement vers la variable aléatoire  $X$  si et seulement si la suite  $(\mu_n)_{n \geq 0}$  des lois des v. a.  $X_n$  converge faiblement vers la loi  $\mu$  de  $X$ . La CNS de la définition se traduit alors ainsi : pour toute fonction *continue bornée*  $f$  de  $S$  dans  $\mathbb{R}$ ,

$$\lim_n \mathbf{E}[f(X_n)] = \mathbf{E}[f(X)].$$

Il en découle immédiatement que

**Propriété** (Corollaire fondamental). *Si  $X_n$  converge faiblement vers  $X$ , et si  $\Phi$  est une fonction continue de  $S$  dans  $T$  (deux espaces métriques), alors  $\Phi(X_n)$  converge faiblement vers  $\Phi(X)$ .*

*Démonstration.* Pour toute fonction  $f$  continue bornée de  $T$  dans  $\mathbb{R}$ ,  $f \circ \Phi$  est continue bornée de  $S$  dans  $\mathbb{R}$ , donc

$$\lim_n \mathbf{E}[f(\Phi(X_n))] = \mathbf{E}[f(\Phi(X))].$$

□

*Quelques exemples de fonctions continues sur  $S = \mathcal{C}[0, 1]$ .*

1. Pour  $T = \mathbb{R}$  ou  $\mathbb{R}^d$ , et pour des nombres réels  $t, t_1, \dots, t_d$  fixés dans  $[0, 1]$ , les applications  $f \mapsto \Phi_t(f) = f(t)$  et  $f \mapsto \Phi_t(f) = (f(t_1), \dots, f(t_d))$  sont continues, donc  $X_n(t) \xrightarrow{\text{faiblement}} X(t)$  et

$$(X_n(t_1), \dots, X_n(t_d)) \xrightarrow{\text{faiblement}} (X(t_1), \dots, X(t_d)).$$

Cette conséquence de la convergence faible est appelée convergence des *distributions fini-dimensionnelles* de  $X_n$  vers celles de  $X$ . La convergence des distributions fini-dimensionnelles ne suffit pas à assurer la convergence faible, elle implique seulement que s'il y a convergence, alors  $X$  est la limite. Pour un exemple simple où il n'y a pas convergence faible, alors qu'il y a convergence des distributions fini-dimensionnelles, voir la section suivante.

2.  $f \mapsto (\max f, \min f, \int_0^1 f(t) dt)$  est continue. Dans le cas du maximum, la convergence en loi de la hauteur des arbres généraux apparaît alors comme une conséquence du théorème clé, version Kaigh. Dans le même goût, la convergence en loi de la largeur des arbres simples est une conséquence de la convergence du profil, démontrée par Drmota et Gittenberger.
3.  $f \mapsto \operatorname{argmax} f$  n'est pas continue sur  $\mathcal{C}[0, 1]$ , non plus que la suite des longueurs des intervalles séparant les zéros de  $f$  (on parle de longueurs des « excursions » de  $f$ ).

En particulier, la convergence jointe de deux statistiques intéressantes ne coûte pas plus cher que la convergence d'une seule. Les derniers contre-exemples frustrants appellent un théorème relaxant l'hypothèse de continuité sur  $\Phi$ . Notons  $\mathcal{D}_\Phi$  l'ensemble des discontinuités de  $\Phi$ .

**Théorème** (Voir [2, Th. 5.1, p. 30]). *Si  $X_n \xrightarrow{\text{faiblement}} X$ , et si  $\mathbf{P}(X \in \mathcal{D}_\Phi) = 0$ , alors  $\Phi(X_n)$  converge faiblement vers  $\Phi(X)$ .*

La démonstration utilise le théorème « porte-manteau », qu'on verra un peu plus tard. Donnons deux exemples d'application :

- posons  $\theta(f) = \sup \{ x \in [0, 1] \mid f(x) = \max_{[0,1]} f \}$ . Alors  $\theta$  n'est pas continue sur  $\mathcal{C}[0, 1]$ ,  $\mathcal{D}_\theta$  étant l'ensemble des fonctions continues qui atteignent leur maximum en plus d'un point.

Il se trouve que le mouvement brownien standard  $B$ , avec probabilité 1, atteint son maximum en un seul point de  $[0, 1]$ , donc

$$(1) \quad \mathbf{P}(B \in \mathcal{D}_\theta) = 0.$$

– posons  $T_a(f) = \inf \{x \geq 0 \mid f(x) \geq a\}$ , l'infimum de l'ensemble vide étant par convention pris égal à  $+\infty$ ;  $\mathcal{D}_{T_a}$  est l'ensemble des fonctions  $f$  satisfaisant  $f \leq a$  sur un intervalle  $[T_a(f), T_a(f) + h]$ ,  $h > 0$ . Il se trouve qu'avec probabilité 1,  $T_a(B)$  est un point d'accumulation de  $\{t \mid B_t > a\}$ , entraînant que

$$(2) \quad \mathbf{P}(B \in \mathcal{D}_{T_a}) = 0.$$

Les propriétés 3. et 4. du mouvement brownien sont des conséquences plus ou moins directes de la propriété de Markov forte<sup>4</sup>. De nombreux processus stochastiques héritent<sup>5</sup> des propriétés (1) et (2) du mouvement brownien.

Les théorèmes de cette section permettent d'exploiter les résultats de Donsker *et al.*, mais réciproquement, joints avec des considérations combinatoires, ils permettent de trouver ou de retrouver les lois de fonctionnelles intéressantes du mouvement brownien et de ses avatars.

### Exercices.

1. Posons  $M_n = \max_{0 \leq k \leq n} S_k$ . Montrer que pour  $k \geq 0$

$$\mathbf{P}(M_n \geq k) = \mathbf{P}(S_n \geq k + 1) + \mathbf{P}(S_n \geq k).$$

Utiliser le Théorème central limite (version de Moivre<sup>6</sup>) pour en déduire que

$$\max \{B_s \mid 0 \leq s \leq 1\} \stackrel{\text{loi}}{=} |B_1|.$$

Une étape possible est de calculer

$$\mathbf{P}(M_n \geq k \text{ et } S_n \geq \ell),$$

ce qui permet en prime d'obtenir la densité jointe de  $(B_1, \max \{B_s \mid 0 \leq s \leq 1\})$ .

2. Notons  $\theta$  le lieu où le mouvement brownien atteint son maximum. Montrer que  $\theta$  suit la loi de l'arcsinus, *i. e.* pour  $0 \leq a \leq b \leq 1$ ,

$$\mathbf{P}(\theta \in [a, b]) = \int_a^b \frac{dx}{\pi \sqrt{x(1-x)}} = \frac{1}{\pi} (\arcsin(2b-1) - \arcsin(2a-1)).$$

Pour cela, on pourra montrer que le lieu  $\theta_n$  du premier maximum d'un chemin de Bernoulli de longueur  $n$  satisfait, pour  $1 \leq k \leq n-1$ ,

$$\mathbf{P}(\theta_n = k) = \binom{k-1}{\lfloor \frac{k-1}{2} \rfloor} \binom{n-k}{\lfloor \frac{n-k}{2} \rfloor} 2^{-n},$$

et établir une convergence locale à l'aide de bornes sur le deuxième terme dans la formule de Stirling (si on veut être complètement rigoureux). On voit que le maximum est atteint avec une forte probabilité hors des intervalles  $[a, 1-a]$ , la densité de probabilité de  $\theta$  ayant des pôles en 0 et 1.

<sup>4</sup>Pour (1), voir [12, preuve du Th. 2.9.12 p. 107]. Pour (2), voir la Section 5 de ce mini-cours.

<sup>5</sup>en vertu du théorème de Cameron–Martin–Girsanov, cf. [17, Ch. 8].

<sup>6</sup>Voir Section 5.

3. Montrer que la valeur terminale du méandre brownien,  $m(1)$ , suit la loi de Rayleigh, à savoir, pour  $0 \leq a \leq b$ ,

$$\mathbf{P}(m(1) \in [a, b]) = \int_a^b x \exp\left(-\frac{x^2}{2}\right) dx = \exp\left(-\frac{a^2}{2}\right) - \exp\left(-\frac{b^2}{2}\right).$$

4. Montrer que le lieu du maximum du pont brownien est uniformément distribué sur  $[0, 1]$ <sup>7</sup>. Y a-t-il une démonstration combinatoire du fait que la valeur maximale du pont brownien suit la loi de Rayleigh<sup>8</sup>?
5. Démontrer la formule (11.5) page 78 de [2]. En déduire la loi du maximum de l'excursion brownienne<sup>9</sup>.

#### 4. Convergence faible : critères et autres caractérisations

**Théorème** (Théorème « porte-manteau », voir [2, Th. 2.1, p. 11]).  $X_n$  converge faiblement vers  $X$  si et seulement si une des conditions suivantes est remplie :

1.  $\lim_n \mathbf{E}[f(X_n)] = \mathbf{E}[f(X)]$  pour toute fonction  $f$  continue bornée de  $S$  dans  $\mathbb{R}$  ;
2.  $\lim_n \mathbf{E}[f(X_n)] = \mathbf{E}[f(X)]$  pour toute fonction  $f$  bornée, uniformément continue, de  $S$  dans  $\mathbb{R}$  ;
3.  $\limsup_n \mathbf{P}(X_n \in F) \leq \mathbf{P}(X \in F)$  pour tout fermé  $F$  de  $S$  ;
4.  $\liminf_n \mathbf{P}(X_n \in G) \geq \mathbf{P}(X \in G)$  pour tout ouvert  $G$  de  $S$  ;
5.  $\lim_n \mathbf{P}(X_n \in A) = \mathbf{P}(X \in A)$  pour tout  $A$  de  $S$  qui vérifie  $\mathbf{P}(X \in \partial A) = 0$ .

Ici encore on pourra se reporter à [2] pour les développements. Une classe  $\mathcal{A}$  de fonctions de  $S$  caractérise une loi de probabilité si pour tout choix de deux variables  $X$  et  $Y$  à valeurs dans  $S$ , on a

$$\forall f \in \mathcal{A}, \quad \mathbf{E}[f(X)] = \mathbf{E}[f(Y)] \quad \Rightarrow \quad X \stackrel{\text{loi}}{=} Y.$$

*Exemples.*

1. Pour  $S = \mathbb{R}$ , la fonction de répartition caractérise une loi de probabilité, ce qui revient à dire que la classe  $\mathcal{A} = \{\mathbf{1}_{(-\infty, x]} \mid x \in \mathbb{R}\}$  est caractérisante.
2. Pour  $S = \mathbb{R}^d$ , la classe  $\mathcal{A} = \{\Phi_{\vec{t}} \mid \vec{t} \in \mathbb{R}^d\}$ , où  $\Phi_{\vec{t}}$  est défini par

$$\Phi_{\vec{t}}(\vec{x}) = e^{i\vec{t} \cdot \vec{x}}$$

est caractérisante,  $\vec{t} \mapsto \mathbf{E}[e^{i\vec{t} \cdot X}]$  étant appelée *fonction caractéristique* de  $X$ .

3. Pour  $S = \mathcal{C}[0, 1]$ ,  $\mathcal{C}[a, b]$  ou  $\mathcal{C}[0, +\infty)$  la classe  $\mathcal{A} = \{\Phi_{\vec{t}} \mid d \geq 1, \vec{t} \in \mathbb{R}^d\}$ , où  $\Phi_{\vec{t}}$  est défini par

$$\Phi_{\vec{t}}(f) = (f(t_1), \dots, f(t_d))$$

est caractérisante.

4. La classe  $\mathcal{A}$  des fonctions bornées et uniformément continues de  $S$  dans  $\mathbb{R}$  est caractérisante.

La convergence de  $\mathbf{E}[\Phi(X_n)]$  vers  $\mathbf{E}[\Phi(X)]$  pour toutes les fonctions  $\Phi$  d'une classe caractérisante  $\mathcal{A}$  suffit-elle à assurer la convergence faible de  $X_n$  vers  $X$  ? La réponse est différente pour chacun des exemples ci-dessus : pour 2. c'est oui, en vertu du Théorème de continuité de Paul Lévy [2, Théorème 7.6, p. 46], et il s'agit d'une CNS. Pour 1. c'est aussi oui, mais la condition est largement trop restrictive : il s'agit d'une condition nécessaire seulement si la loi limite est *diffuse* (i. e.

<sup>7</sup>Utiliser le lemme cyclique attribué parfois à Raney, parfois à Dvoretzki ou à Motzkin.

<sup>8</sup>cf. [2, Section 11], mais on peut sûrement trouver un raccourci (je n'ai pas eu le temps de m'en assurer).

<sup>9</sup>C'est, en particulier, la loi limite pour la hauteur ou la largeur des arbres simples [8, 16].

$\mathbf{P}(X = a) = 0$  pour tout  $a$  dans  $\mathbb{R}$ ), en vertu du 5. du Théorème « porte-manteau », puisque  $\partial(-\infty, a] = \{a\}$  ! Enfin, la réponse est non pour l'exemple 3., comme le montre l'exemple suivant tiré de [2] : prenons  $X$  et  $X_n$  non aléatoires à valeur dans  $\mathcal{C}[0, 1]$ ,  $X \equiv 0$  et  $X_n \equiv f_n$ ,  $f_n$  ayant le graphe ci-dessous : les distributions fini-dimensionnelles de  $X_n$  convergent bien faiblement vers

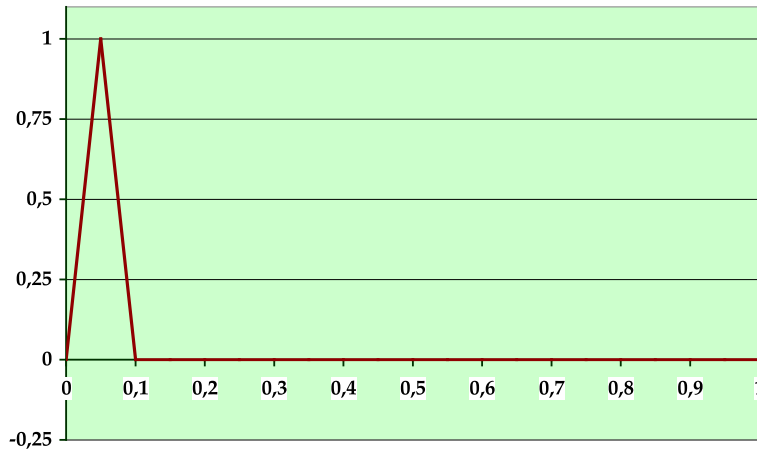


FIGURE 4.  $f_n$  est continue et affine par morceau, avec ici pour  $n = 10$ ,  $(f(0), f(1/2n), f(1/n), f(1)) = (0, 1, 0, 0)$ .

les probabilités concentrées sur  $0 \in \mathbb{R}^d$ , *i. e.* vers les distributions fini-dimensionnelles de  $X$ , mais  $\Phi(X_n) = \max X_n \equiv 1$  ne converge pas faiblement vers  $\Phi(X) = \max X \equiv 0$ .

Il faut donc une condition supplémentaire à la convergence des distributions fini-dimensionnelles pour obtenir la convergence faible des variables aléatoires à valeur dans  $\mathcal{C}[0, 1]$  : c'est la condition de *tension*.

**Définition.** La suite de variables aléatoires  $X_n$  est *tendue* (ou *équitendue*) si et seulement si pour tout  $\varepsilon > 0$  il existe un compact  $K_\varepsilon$  de  $S$  tel que

$$\forall n, \quad \mathbf{P}(X_n \notin K_\varepsilon) \leq \varepsilon.$$

Le Théorème de Prohorov [2, Section 6] assure que la tension est une CS (et une CNS si  $S = \mathcal{C}[0, 1]$ ) pour la *relative compacité* d'une suite de mesures de probabilité (ici les lois des v. a.  $X_n$ ). Il suit que cette suite de variables  $(X_n)_{n \geq 0}$  possède au moins une valeur d'adhérence pour la convergence faible. On connaît les distributions fini-dimensionnelles de cette valeur d'adhérence, ce sont les limites des distributions fini-dimensionnelles de  $X_n$ , donc ce sont les distributions fini-dimensionnelles de  $X$ , donc  $X$  est la seule valeur d'adhérence de  $X_n$ , or une suite relativement compacte ayant une seule valeur d'adhérence est convergente. Finalement :

**Théorème.** *Si une suite de variables aléatoires  $X_n$  variables aléatoires à valeurs dans  $\mathcal{C}[0, 1]$  est tendue, et si ses distributions fini-dimensionnelles convergent vers celles de  $X$ , alors  $X_n$  converge faiblement vers  $X$ .*

Le chapitre 2 de [2] donne une foule de critères de tension dans  $\mathcal{C}[0, 1]$ , basées sur la caractérisation d'Arzelà–Ascoli des compacts de  $\mathcal{C}[0, 1]$ . Par exemple, les démonstrations de Donsker, Iglehart et Kaigh sont basées sur de tels critères, ainsi que la démonstration par Drmota et Gitteberger de la convergence du profil des arbres simples. Il existe des traitements plus modernes [10, 14, 15], mais [2] est déjà très lisible et complet.

Il faut aussi parler du lien entre convergence presque sûre, en probabilité, et dans  $\mathcal{L}^p$  d'une part, convergence faible d'autre part. Les premières citées exigent que les variables  $X_n$  et  $X$ , à valeurs dans le même espace  $S$  à l'arrivée, soient aussi définie sur *le même triplet probabiliste*  $(\Omega, \mathcal{A}, \mathbb{P})$  au départ, alors que la convergence faible, étant en fait uniquement la convergence de la mesure image par  $X_n$  vers la mesure image par  $X$  des mesures de probabilité des espaces de départ, exige seulement que  $X_n$  et  $X$  aient *le même espace d'arrivée*  $S$ . Notons  $d(\cdot, \cdot)$  la distance sur  $S$ .

**Définition.** Une suite  $(X_n)_{n \geq 0}$  converge :

1. *presque sûrement* vers  $X$  si et seulement si

$$\mathbf{P}\left(\left\{\omega \in \Omega \mid \lim_n d(X_n(\omega), X(\omega)) = 0\right\}\right) = 1 ;$$

2. *en probabilité* vers  $X$  si et seulement si

$$\forall \varepsilon > 0, \quad \lim_n \mathbf{P}\left(\left\{\omega \in \Omega \mid d(X_n(\omega), X(\omega)) \geq \varepsilon\right\}\right) = 0 ;$$

3. vers  $X$  dans  $\mathcal{L}^p$  si et seulement si

$$\lim_n \mathbf{E}\left[d(X_n, X)^p\right] = 0.$$

**Théorème.** *Les trois convergences ci-dessus entraînent la convergence faible.*

*Démonstration.* Seulement pour le 1., pour une fonction continue  $\Phi$ ,  $\Phi(X_n)$  converge presque sûrement vers  $\Phi(X)$ , et si de plus  $\Phi$  est bornée, le Théorème de convergence dominée entraîne bien que  $\lim_n \mathbf{E}[\Phi(X_n)] = \mathbf{E}[\Phi(X)]$ . Par ailleurs, 3. entraîne 2. en vertu de l'inégalité de Markov. Pour montrer que 2. entraîne la convergence faible, il faut utiliser la caractérisation 2. du Théorème « porte-manteau » et travailler à peine un peu plus.  $\square$

Finalement il y a une quasi-réciproque utile au théorème précédent, c'est le

**Théorème** (Théorème de représentation de Skorohod, voir [18, II.86.1, p. 215]). *Si  $S$  est un espace de Lusin (en particulier pour  $S = \mathcal{C}[0, 1]$ ) et si la suite de variables aléatoires  $(X_n)_{n \geq 0}$ , à valeurs dans  $S$ , converge faiblement vers  $X$ , alors il existe un triplet probabiliste  $(\Omega, \mathcal{A}, \mathbb{P})$ , et, définies sur ce triplet, des copies  $(\hat{X}_n)_{n \geq 0}$  et  $\hat{X}$  de  $(X_n)_{n \geq 0}$  et de  $X$ , telles que  $(\hat{X}_n)_{n \geq 0}$  converge presque sûrement vers  $\hat{X}$ .*

Par « copie », on entend que  $X_n$  et  $\hat{X}_n$ , ou encore  $X$  et  $\hat{X}$ , ont même loi. Par exemple, il n'est pas toujours naturel de construire des arbres simples aléatoires, ou des graphes aléatoires, de tailles différentes, sur le même espace de probabilité : il est beaucoup plus fréquent de considérer, par exemple, l'ensemble  $\mathcal{T}_n$  des arbres étiquetés de taille  $n$  comme un espace de probabilité à lui tout seul, muni de la probabilité uniforme. Plonger tous les  $\mathcal{T}_n$  dans un même triplet probabiliste évite pourtant parfois certains calculs de lois fini-dimensionnelles : ils sont remplacés par des estimations plus faciles conduisant à une convergence presque sûre<sup>10</sup>. Par ailleurs, le Théorème de représentation de Skorohod est un outil très commode pour les démonstrations de la Section 6.

## 5. Propriétés du mouvement brownien

Le but ici n'est certainement pas de donner de démonstration, mais, à titre mnémotechnique, de montrer comment le mouvement brownien imite les (ou hérite des) propriétés de la marche aléatoire simple symétrique.

<sup>10</sup>Voir par exemple [4].

**Accroissements indépendants et stationnaires.** La marche aléatoire simple symétrique possède des accroissements indépendants : sous  $\mu_n$ , pour  $1 \leq k_1 \leq k_2 \leq \dots \leq k_i \leq n$ ,

$$(Y_1 + \dots + Y_{k_1}) \perp (Y_{k_1+1} + \dots + Y_{k_2}) \perp \dots \perp (Y_{k_{i-1}+1} + \dots + Y_{k_i})$$

*i. e.*

$$S_{k_1} \perp (S_{k_2} - S_{k_1}) \perp \dots \perp (S_{k_i} - S_{k_{i-1}})$$

et stationnaires

$$(Y_{k+1} + \dots + Y_{k+\ell}) \stackrel{\text{loi}}{=} (Y_1 + \dots + Y_\ell)$$

*i. e.*

$$S_{k+\ell} - S_k \stackrel{\text{loi}}{=} S_\ell.$$

Le mouvement brownien aussi ! C'est-à-dire sous  $\mu$ , pour  $0 \leq t_1 \leq t_2 \leq \dots \leq t_i \leq 1$ ,

$$B_{t_1} \perp (B_{t_2} - B_{t_1}) \perp \dots \perp (B_{t_i} - B_{t_{i-1}})$$

et pour  $t \geq 0$ ,  $s \geq 0$ ,

$$B_{t+s} - B_t \stackrel{\text{loi}}{=} B_s.$$

Cela entraîne la propriété de Markov faible.

**Propriété** (Propriété de Markov faible). *Le nouveau processus  $W = (W_s)_{0 \leq s \leq h}$ , défini par*

$$W_s = B_{t+s} - B_t$$

*est indépendant de  $(B_s)_{0 \leq s \leq t}$ . De plus  $W$  a même loi que  $(B_s)_{0 \leq s \leq h}$ .*

La démonstration requiert seulement de vérifier que pour chaque  $k, \ell$ , et pour chaque suite de nombres réels  $0 < t_1 < t_2 < \dots < t_k \leq t$  et  $0 < s_1 < s_2 < \dots < s_\ell \leq h$ ,

$$(B_{t_i})_{1 \leq i \leq k} \perp (W_{s_i})_{1 \leq i \leq \ell} \quad \text{et} \quad (W_{s_i})_{1 \leq i \leq \ell} \stackrel{\text{loi}}{=} (B_{s_i})_{1 \leq i \leq \ell}.$$

Pour l'indépendance, il suffit de remarquer que  $(B_{t_i})_{1 \leq i \leq k} \perp (W_{s_i})_{1 \leq i \leq \ell}$  est équivalent à

$$(B_{t_1}, B_{t_2} - B_{t_1}, \dots, B_{t_k} - B_{t_{k-1}}) \perp (W_{s_1}, W_{s_2} - W_{s_1}, \dots, W_{s_\ell} - W_{s_{\ell-1}})$$

et de remarquer que

$$(W_{s_1}, W_{s_2} - W_{s_1}, \dots, W_{s_\ell} - W_{s_{\ell-1}}) = (B_{t+s_1} - B_t, B_{t+s_2} - B_{t+s_1}, \dots, B_{t+s_\ell} - B_{t+s_{\ell-1}}).$$

Cette dernière égalité plus la stationnarité des accroissements donne aussi l'égalité en loi.

**Remarque.** On a bien sûr  $t > 0$ ,  $h > 0$ , et on doit pour le moment imposer  $t + h \leq 1$ , mais cette dernière inégalité est en fait superflue car il est naturel de définir le mouvement brownien sur  $[0, +\infty)$  (comme de définir la marche aléatoire simple symétrique  $(S_k)_{k \geq 0}$  pour chaque entier positif).

**Une construction possible du mouvement brownien sur la demi-droite des entiers positifs.** Considérons par exemple une suite  $(B^{(n)})_{n \geq 0}$ ,  $B^{(n)} = (B_s^{(n)})_{0 \leq s \leq 1}$  de mouvements browniens mutuellement indépendants<sup>11</sup>. Définissons alors  $B = (B_t)_{t \geq 0}$  comme un élément aléatoire de  $\mathcal{C}[0, +\infty)$ , tel que pour  $n \leq s \leq t \leq n+1$ ,

$$B_t - B_s = B_t^{(n)} - B_s^{(n)},$$

c'est à dire qu'on recolle les graphes (trajectoires) des  $B^{(n)}$  pour former le graphe de  $B$ . Il est alors facile de voir que  $B$  hérite des  $B^{(n)}$  l'indépendance des accroissements. Il en hérite aussi la stationnarité des accroissements, mais, pour le voir, il faut parler un peu de la loi de ces accroissements.

**Lois des accroissements du mouvement brownien.** La formule de Stirling, fondamentale en combinatoire, est née des travaux de de Moivre qui sont en quelque sorte un premier pas vers le mouvement brownien<sup>12</sup>. Posons

$$S_{k+\ell} - S_k = -\ell + 2Z.$$

Alors  $Z$  suit la loi binomiale  $(\ell, \frac{1}{2})$ , i. e. pour  $0 \leq i \leq \ell$ ,

$$\mathbf{P}(Z = i) = \binom{\ell}{i} \frac{1}{2^\ell}.$$

On sait, depuis que de Moivre<sup>13</sup> a démontré la formule de Stirling<sup>14</sup>, et l'approximation « gaussienne » de la loi binomiale<sup>15</sup>, que l'on peut écrire, pour  $\ell = 2\lfloor sn/2 \rfloor \sim ns$ ,

$$\begin{aligned} \mathbf{P}(S_{k+\ell} - S_k = 2\lfloor x\sqrt{n}/2 \rfloor) &= \mathbf{P}\left(\frac{S_{k+\ell} - S_k}{\sqrt{n}} \in \left[\frac{2\lfloor x\sqrt{n}/2 \rfloor - 1}{\sqrt{n}}, \frac{2\lfloor x\sqrt{n}/2 \rfloor + 1}{\sqrt{n}}\right]\right) \\ &\sim \frac{2}{\sqrt{n}} \frac{1}{\sqrt{2\pi s}} e^{-x^2/2s} \sim \mathbf{P}\left(N\sqrt{s} \in \left[x - \frac{1}{\sqrt{n}}, x + \frac{1}{\sqrt{n}}\right]\right), \end{aligned}$$

où  $N$  est une variable aléatoire suivant la loi normale (ou gaussienne) centrée réduite, souvent notée  $\mathcal{N}(0, 1)$ , à savoir

$$\mathbf{P}(N \in [a, b]) = \int_a^b \frac{1}{\sqrt{2\pi}} e^{-x^2/2} dx.$$

En d'autres termes,  $\frac{S_{k+\ell} - S_k}{\sqrt{n}}$  a approximativement la même loi que  $\sqrt{s} N$ , à savoir, la loi normale (ou gaussienne) centrée de variance  $s$ , notée traditionnellement  $\mathcal{N}(0, s)$ . D'autre part,  $\frac{S_{k+\ell} - S_k}{\sqrt{n}}$  est l'accroissement, entre les points  $\frac{k}{n}$  et, approximativement,  $\frac{k}{n} + s$ , de la fonction obtenue, à partir de la marche aléatoire simple symétrique, par changement d'échelle brownien. Par passage à la limite, on en déduit que

**Propriété (Accroissements gaussiens).** *Indépendamment de  $t$ ,*

$$B_{t+s} - B_t \stackrel{\text{loi}}{=} \sqrt{s} N.$$

<sup>11</sup>On peut par exemple définir une telle suite comme un élément au hasard de  $\mathcal{C}[0, 1]^{\mathbb{N}}$  muni du produit infini de mesures de Wiener  $\mu^{\otimes \mathbb{N}}$ .

<sup>12</sup>un peu forcé, le rapprochement, non ?

<sup>13</sup>Voir [http://www-groups.dcs.st-andrews.ac.uk/~history/Mathematicians/De\\_Moivre.html](http://www-groups.dcs.st-andrews.ac.uk/~history/Mathematicians/De_Moivre.html).

<sup>14</sup>dans *Miscellanea Analytica*, 1730.

<sup>15</sup>dans *Approximatio ad Summam Terminorum Binomii  $(a+b)^n$  in Seriem expansi*, 1733.

Notons

$$p_s(x, y) = \frac{1}{\sqrt{2\pi s}} e^{-\frac{(y-x)^2}{2s}}.$$

On peut voir  $p_s(x, y)$  comme la densité de probabilité de  $x + \sqrt{s}N$ , *i. e.*, en vertu de la propriété d'accroissements gaussiens indépendants, comme la densité conditionnelle de  $B_{t+s}$ , sachant que  $B_t = x$ . On en déduit que

**Propriété** (Distributions fini-dimensionnelles du mouvement brownien). *La densité de probabilité  $f$  de  $(B_{t_1}, B_{t_2}, \dots, B_{t_d})$  est donnée par la formule*

$$f(x_1, x_2, \dots, x_d) = p_{t_1}(0, x_1) p_{t_2-t_1}(x_1, x_2) \dots p_{t_d-t_{d-1}}(x_{d-1}, x_d).$$

*Une autre manière de caractériser les distributions fini-dimensionnelles du mouvement brownien est de remarquer que  $(B_{t_1}, B_{t_2}, \dots, B_{t_d})$  est un vecteur gaussien centré, dont la loi est donc caractérisée par sa matrice de covariance. On calcule facilement le terme général :*

$$\mathbf{Cov}(B_{t_i}, B_{t_j}) = \min(t_i, t_j).$$

En effet, pour  $s \leq t$ ,

$$\mathbf{Cov}(B_s, B_t) = \mathbf{Cov}(B_s, B_s) + \mathbf{Cov}(B_s, B_t - B_s) = \mathbf{Var}(B_s) = \mathbf{Var}(\sqrt{s}N) = s \mathbf{Var}(N) = s,$$

la deuxième égalité découlant de  $B_s \perp B_t - B_s$ .

Rappelons qu'une variable aléatoire  $X = (X_1, X_2, \dots, X_d)$  à valeurs dans  $\mathbb{R}^d$  est un *vecteur gaussien* si et seulement si toutes les combinaisons linéaires de ses composantes sont gaussiennes (ont même loi que  $m + \sigma N$ , pour un choix approprié de  $m$  et  $\sigma$ ), ou encore, si et seulement si  $X$  est image par une transformation affine (disons,  $X = \vec{m} + A\vec{N}$ ) d'un vecteur  $\vec{N} = (N_1, N_2, \dots, N_k)$  dont les composantes  $N_i$  sont i. i. d. et de loi  $\mathcal{N}(0, 1)$ . Dans le cas des distributions fini-dimensionnelles du mouvement brownien, on a  $\vec{m} = 0$ , et on peut exhiber  $A$  et  $\vec{N}$ , en posant

$$N_1 = \frac{B_{t_1}}{\sqrt{t_1}}, \quad N_i = \frac{B_{t_i} - B_{t_{i-1}}}{\sqrt{t_i - t_{i-1}}}.$$

La loi d'un vecteur gaussien est caractérisée par l'espérance de chacune de ses composantes et par sa matrice de covariance. Dans la représentation affine ci-dessus,  $\vec{m}$  est le vecteur des espérances des composantes, et la matrice de covariance est  $\Gamma = {}^t A A$ .

**Définition.** Un processus  $X$  dont les distributions fini-dimensionnelles sont gaussiennes est appelé *processus gaussien*. La loi du processus est alors caractérisée par sa *fonction moyenne*  $m(t) = \mathbf{E}[X_t]$  et sa *fonction covariance*  $\Gamma(s, t) = \mathbf{Cov}(X_s, X_t)$ .

Le mouvement brownien et, comme on le verra en Section 7, le pont brownien, sont deux exemples de processus gaussiens centrés ( $m(t) \equiv 0$ ). La fonction covariance du mouvement brownien est

$$\Gamma(s, t) = \min(s, t).$$

**Théorème** (Transformations des trajectoires du mouvement brownien). *Le mouvement brownien est préservé par les transformations suivantes :*

- Symétrie :  $W^{(1)} = (-B_t)_{t \geq 0}$  est un mouvement brownien.
- Décalage : Pour  $t_0 \geq 0$ ,  $W^{(2)} = (B_{t_0+t} - B_{t_0})_{t \geq 0}$  est un mouvement brownien.
- Changement d'échelle : Pour  $c > 0$ ,  $W^{(3)} = (\frac{1}{\sqrt{c}} B_{ct})_{t \geq 0}$  est un mouvement brownien.
- Inversion du temps :  $W^{(4)} = (W_t^{(4)})_{t \geq 0}$  défini par  $W_t^{(4)} = tB_{1/t}$ , pour  $t > 0$ , et par  $W_0^{(4)} = 0$ , est un mouvement brownien.



*Démonstration.* Chacun de ces processus est gaussien centré : il suffit de calculer sa fonction covariance. Dans les quatre cas, on trouve  $\Gamma^{(i)}(s, t) = \min(s, t)$ . Reste un petit problème : la continuité de  $W^{(4)}$  en 0, qui n'est pas automatique. La loi forte des grands nombres<sup>16</sup> pour le mouvement brownien, stipule que

$$\mathbf{P} \left( \lim_{+\infty} \frac{B_t}{t} = 0 \right) = 1.$$

En conséquence

$$\mathbf{P} \left( \left\{ \omega \in \Omega \mid t \rightarrow W_t^{(4)}(\omega) \text{ est continue en } 0 \right\} \right) = 1.$$

Le processus  $W^{(4)}$  est donc *presque* sûrement continu en 0, alors que le mouvement brownien, tel qu'on l'a défini, est à valeurs dans  $\mathcal{C}[0, 1]$ , c'est-à-dire que  $t \rightarrow B_t(\omega)$  est continu en 0 pour *tout*  $\omega$ . Régler ce genre de problème rigoureusement est justement ce que je veux éviter dans une introduction au mouvement brownien prévue pour être succincte<sup>17</sup>.  $\square$

**Temps d'atteinte.** Le temps d'atteinte de la hauteur  $a > 0$ , noté  $T_a$ , est défini par

$$T_a = \begin{cases} \inf \{ t \geq 0 \mid B_t \geq a \} & \text{si l'ensemble n'est pas vide,} \\ +\infty & \text{si l'ensemble est vide.} \end{cases}$$

**Théorème.**  $T_a$  a même loi que  $\frac{a^2}{N^2}$ , en particulier  $\mathbf{P}(T_a = +\infty) = 0$ .

*Démonstration.* On a

$$\begin{aligned} \mathbf{P}(T_a > t) &= \mathbf{P}(\max \{ B_s \mid 0 \leq s \leq t \} < a) = \mathbf{P} \left( \max \left\{ \frac{1}{\sqrt{t}} B_{ts} \mid 0 \leq s \leq 1 \right\} < \frac{a}{\sqrt{t}} \right) \\ &= \mathbf{P} \left( \max \{ B_s \mid 0 \leq s \leq 1 \} < \frac{a}{\sqrt{t}} \right) = \mathbf{P} \left( |B_1| < \frac{a}{\sqrt{t}} \right) = \mathbf{P} \left( \frac{a^2}{B_1^2} > t \right), \end{aligned}$$

la troisième égalité par changement d'échelle, la quatrième comme conséquence de l'exercice 1, Section 3.  $\square$

**Définition.** Une v. a.  $T$  à valeurs dans  $[0, +\infty]$  est un *temps d'arrêt* du mouvement brownien si et seulement si  $\{ \omega \mid T(\omega) \leq t \}$  est dans la tribu engendrée par  $(B_s)_{0 \leq s \leq t}$ , en d'autres termes, si on peut décider de la véracité de l'affirmation «  $T(\omega) \leq t$  » en observant la trajectoire du mouvement brownien seulement jusqu'à l'instant  $t$  (inclus).

En particulier, les temps d'atteinte  $T_a$  sont des temps d'arrêts.

**Propriété** (Propriété de Markov forte, cf. [12, Section 2.5]).  $T$  étant un temps d'arrêt, le nouveau processus  $W^T = (W_s^T)_{0 \leq s \leq T}$ , défini par

$$W_s^T = B_{T+s} - B_T$$

est indépendant de  $(B_s)_{0 \leq s \leq T}$ . De plus  $W^T$  a même loi que le mouvement brownien.

<sup>16</sup>Pour une démonstration simple, voir [12, Problème 9.3, p. 104 et Remarque 3.10, p. 15]. On peut être plus précis sur le comportement du mouvement brownien en  $+\infty$  : voir, [12, p. 112], la loi du logarithme itérée due à Khintchine, 1933.

<sup>17</sup>Il se trouve que  $W^{(4)}$  est indistinguable d'un processus à valeurs dans  $\mathcal{C}[0, 1]$ , voir [12, Section 1.1].

**Quelques conséquences.**

- $T_{a+b} - T_a$  est le temps d'atteinte de  $b$  par le processus  $W^{T_a}$ , il est donc indépendant de  $T_a$  et a la même loi que  $T_b$ . En d'autres termes, le processus  $(T_a)_{a \geq 0}$  est à accroissements indépendants et stationnaires<sup>18</sup>.
- Presque sûrement,  $+\infty$  est un point d'accumulation de l'ensemble des zéros du mouvement brownien : posons  $T$  le premier zéro du mouvement brownien après l'instant 1 (*i. e.*  $T = \inf \{ t \geq 1 \mid B_t = 0 \}$ ). La loi conditionnelle de  $T$  sachant que  $B_1 = a$  est la loi de  $T_a$ , donc  $T$  est presque sûrement fini ;  $T$  est un temps d'arrêt donc  $W^T$  est lui-même un mouvement brownien et possède lui aussi un zéro après son instant 1 (donc  $B$  possède un zéro après l'instant 2, etc.).
- De la même manière on voit que, presque sûrement,  $+\infty$  est un point d'accumulation de l'ensemble  $\{ t \geq 0 \mid B_t > 0 \}$ , ou de l'ensemble  $\{ t \geq 0 \mid B_t < 0 \}$ . Ainsi, par inversion du temps, 0 est un point d'accumulation des ensembles  $\{ t > 0 \mid B_t = 0 \}$ ,  $\{ t > 0 \mid B_t > 0 \}$  et  $\{ t > 0 \mid B_t < 0 \}$ .
- Ainsi  $T_a$  est un point d'accumulation des ensembles  $\{ t > T_a \mid B_t = a \}$ ,  $\{ t > T_a \mid B_t < a \}$  et  $\{ t > T_a \mid B_t > a \}$ . Cette toute dernière assertion implique la relation (2).

Ce ne sont que quelques exemples d'application de la propriété de Markov forte, mais en fait on l'applique comme on respire, sans s'en rendre compte. On a commencé à aborder la structure de l'ensemble des zéros du mouvement brownien, alors mentionnons que

**Théorème** (Structure de l'ensemble des zéros du mouvement brownien). *Presque sûrement, l'ensemble des zéros du mouvement brownien est fermé, non borné, sans point isolé, de mesure de Lebesgue nulle, et possède 0 comme point d'accumulation*<sup>19</sup>.

Finalement, mentionnons

**Quelques propriétés locales du mouvement brownien.** Pour un chemin de Bernoulli  $f$  quelconque dans  $\mathcal{C}[0, n]$ , on a

$$\sum_{k=a}^{b-1} |f(k+1) - f(k)|^2 = b - a,$$

pour  $a$  et  $b$  entiers,  $0 \leq a < b \leq n$ . Par scaling brownien, on obtient que presque sûrement pour la mesure de probabilité  $\mu_n$ ,

$$\sum_{k=0}^{n(b-a)-1} \left| f\left(a + \frac{k+1}{n}\right) - f\left(a + \frac{k}{n}\right) \right|^2 = b - a,$$

si  $a$  et  $b$  sont dans  $[0, 1]$  et de la forme  $\frac{\ell}{n}$ ,  $\ell$  entier. Cela se traduit par le fait que le mouvement brownien possède une variation quadratique égale à  $t$  (toute fonction continument dérivable, *p. e.*, possède une variation quadratique nulle). Plus précisément, pour une subdivision  $\Pi = \{t_0, t_1, \dots, t_m\}$  de  $[0, t]$  (*i. e.*  $0 = t_0 \leq t_1 \leq \dots \leq t_m = t$ ), notons

$$V_t^{(2)}(\Pi) = \sum_{k=1}^m |B_{t_k} - B_{t_{k-1}}|^2$$

la *variation quadratique* du mouvement brownien sur la subdivision  $\Pi$ , et notons

$$\|\Pi\| = \max_{1 \leq k \leq m} |t_k - t_{k-1}|$$

<sup>18</sup>mais ses trajectoires ne sont pas continues, cf. [12, Section 6.2.A].

<sup>19</sup>cf. [12, Th. 2.9.6].

le pas de la subdivision  $\Pi$ . On a alors

**Propriété** (Variation quadratique, cf. [12, Th. 1.5.8 et Problème 2.5.5]). *En probabilité,  $V_t^{(2)}(\Pi)$  converge vers  $t$  quand  $\|\Pi\|$  tend vers 0, i. e. pour chaque  $\varepsilon, \eta > 0$ , on peut trouver  $\delta > 0$  tel que  $\|\Pi\| < \delta$  entraîne*

$$\mathbf{P}\left(\left|V_t^{(2)}(\Pi) - t\right| > \varepsilon\right) < \eta.$$

Ceci, avec le fait que presque sûrement sous  $\mu_n$  une fonction possède une dérivée dont la valeur absolue en tout point (sauf en  $\frac{k}{n}$ ) est  $\sqrt{n}$ , laisse à penser que le mouvement brownien a peu de chances d'être dérivable en un point donné. En fait on a un résultat beaucoup plus précis :

**Théorème** (Paley, Wiener & Zygmund, 1933, cf. [12, Th. 2.9.18]).

$$\mathbf{P}\left(\left\{\omega \in \Omega \mid \text{la fonction } t \rightarrow B_t(\omega) \text{ n'est dérivable nulle part}\right\}\right) = 1.$$

Une autre propriété, que l'on peut aussi pressentir en générant des chemins de Bernoulli aléatoires, illustre bien le comportement erratique du mouvement brownien :

**Théorème** (Dvoretzky, Erdős & Kakutani, 1961, cf. [12, Th. 2.9.13]).

$$\mathbf{P}\left(\left\{\omega \in \Omega \mid \text{la fonction } t \rightarrow B_t(\omega) \text{ n'a aucun point de croissance}\right\}\right) = 1.$$

Un point  $t$  est un point de croissance de  $f$  si on peut trouver  $\delta > 0$  tel que pour tout  $y \in [t - \delta, t]$  et tout  $z \in [t, t + \delta]$ ,  $f(y) \leq f(t) \leq f(z)$ .

Cet aperçu des propriétés du mouvement brownien est à la fois très incomplet et assez désordonné. Heureusement la littérature sur le sujet est riche, et on pourra s'y reporter.

### Bibliography

- [1] Biane (Philippe), Pitman (Jim), and Yor (Marc). – Probability laws related to the Jacobi theta and Riemann zeta functions, and Brownian excursions. *Bulletin of the American Mathematical Society (New Series)*, vol. 38, n° 4, 2001, pp. 435–465.
- [2] Billingsley (Patrick). – *Convergence of probability measures*. – John Wiley & Sons Inc., New York, 1968, xii+253p.
- [3] Borodin (Andrei N.) and Salminen (Paavo). – *Handbook of Brownian motion—facts and formulae*. – Birkhäuser Verlag, Basel, 1996, xiv+462p.
- [4] Chassaing (Philippe) and Marckert (Jean-François). – Parking functions, empirical processes, and the width of rooted labeled trees. *Electronic Journal of Combinatorics*, vol. 8, n° 1, 2001, p. Research Paper 14. 19 pages.
- [5] Csörgő (M.) and Révész (P.). – *Strong approximations in probability and statistics*. – Academic Press Inc. [Harcourt Brace Jovanovich Publishers], New York, 1981, 284p.
- [6] Donsker (Monroe D.). – An invariance principle for certain probability limit theorems. *Memoirs of the American Mathematical Society*, vol. 1951, n° 6, 1951, p. 12.
- [7] Erdős (P.) and Kac (M.). – On certain limit theorems of the theory of probability. *Bulletin of the American Mathematical Society*, vol. 52, 1946, pp. 292–302.
- [8] Flajolet (Philippe) and Odlyzko (Andrew). – The average height of binary trees and other simple trees. *Journal of Computer and System Sciences*, vol. 25, n° 2, 1982, pp. 171–213.
- [9] Iglehart (Donald L.). – Functional central limit theorems for random walks conditioned to stay positive. *Annals of Probability*, vol. 2, 1974, pp. 608–619.
- [10] Jacod (Jean) and Shiryaev (Albert N.). – *Limit theorems for stochastic processes*. – Springer-Verlag, Berlin, 1987, xviii+601p.
- [11] Kaigh (W. D.). – An invariance principle for random walk conditioned by a late return to zero. *Annals of Probability*, vol. 4, n° 1, 1976, pp. 115–121.
- [12] Karatzas (Ioannis) and Shreve (Steven E.). – *Brownian motion and stochastic calculus*. – Springer-Verlag, New York, 1991, second edition, xxiv+470p.

- [13] Liggett (Thomas M.). – An invariance principle for conditioned sums of independent random variables. *J. Math. Mech.*, vol. 18, 1968, pp. 559–570.
- [14] Parthasarathy (K. R.). – *Probability measures on metric spaces*. – Academic Press Inc., New York, 1967, xi+276p.
- [15] Pollard (David). – *Convergence of stochastic processes*. – Springer-Verlag, New York, 1984, xiv+215p.
- [16] Rényi (A.) and Szekeres (G.). – On the height of trees. *Journal of the Australian Mathematical Society*, vol. 7, 1967, pp. 497–507.
- [17] Revuz (Daniel) and Yor (Marc). – *Continuous martingales and Brownian motion*. – Springer-Verlag, Berlin, 1999, third edition, xiv+602p.
- [18] Rogers (L. C. G.) and Williams (David). – *Diffusions, Markov processes, and martingales. Vol. 1*. – John Wiley & Sons Ltd., Chichester, 1994, second edition, xx+386p. Foundations.
- [19] Spencer (Joel). – Enumerating graphs and Brownian motion. *Communications on Pure and Applied Mathematics*, vol. 50, n° 3, 1997, pp. 291–294.
- [20] Yor (M.). – *Local times and excursions for brownian motion: a concise introduction*. – Postgrado de Matemáticas, Facultad de Ciencias, Universidad Central de Venezuela, Caracas, 1995, *Lecciones en Matemáticas*, vol. 1.

## CONTENTS

### Part I. Combinatorics

Enumeration of Sand Piles. <i>Talk by S. Corteel, summary by M. Nguyễn-Thé</i> .....	3
On the Group of a Sandpile. <i>Talk by D. Rossin, summary by D. Gouyou-Beauchamps</i> .....	9
The Tennis Ball Problem. <i>Talk by D. Merlini, summary by C. Banderier</i> .....	15
Hyperharmonic Numbers and the Phratry of the Coupon Collector. <i>Talk by D. Foata, summary by C. Banderier</i> .....	19
Mac Mahon's Partition Analysis Revisited. <i>Talk by P. Paule, summary by S. Corteel</i> .....	23
Engel Expansions of $q$ -Series. <i>Talk by P. Paule, summary by B. Salvy</i> .....	27
Eulerian Calculus: a Technology for Computer Algebra and Combinatorics. <i>Talk by D. Foata, summary by D. Gouyou-Beauchamps</i> .....	31

### Part II. Analysis of Algorithms and Combinatorial Structures

Asymptotics for Random Combinatorial Structures. <i>Talk by A. Dembo, summary by P. Flajolet</i> .....	39
Random Walks and Heaps of Cycles. <i>Talk by Ph. Marchal, summary by C. Banderier</i> .....	45
Tail Bounds for Occupancy Problems. <i>Talk by P. Spirakis, summary by S. Boucheron</i> .....	49
Patricia Tries in the Context of Dynamical Systems. <i>Talk by J. Bourdon, summary by M. Nguyễn-Thé</i> .....	53
New and Old Problems in Pattern Matching. <i>Talk by W. Szpankowski, summary by M. Régnier</i> .....	59
Genome Analysis and Sequences with Random Letter Distribution. <i>Talk by M. Termier, summary by M. Vandenbogaert</i> .....	63
Random Sequences and Genomic Analysis. <i>Talk by A. Denise</i> .....	67
The Primal-Dual Schema for Approximation Algorithms: Where Does It Stand, and Where Can It Go? <i>Talk by V. Vazirani, summary by C. Kenyon</i> .....	69

Distributed Decision Making: The Case of No Communication. <i>Talk by P. Spirakis</i> .....	73
--	----

### Part III. Computer Algebra and Applications

Thirty Years of Integer Factorization. <i>Talk by F. Morain, summary by M. Durand</i> .....	77
Variations on Computing Reciprocals of Power Series. <i>Talk by A. Schönhage, summary by L. Meunier</i> .....	81
Fast Multivariate Power Series Multiplication in Characteristic Zero. <i>Talk by G. Lecerf, summary by L. Meunier</i> .....	85
A Tutorial on Closed Difference Forms. <i>Talk by B. Zimmermann, summary by F. Chyzak</i> .....	89
Transformations Exhibiting the Rank for Skew Laurent Polynomial Matrices. <i>Talk by M. Bronstein, summary by A. Bostan</i> .....	95
A Criterion for Non-Complete Integrability of Hamiltonian Systems. <i>Talk by D. Boucher, summary by Ph. Dumas and B. Salvy</i> .....	99
Effective Algebraic Analysis in Linear Control Theory. <i>Talk by A. Quadrat, summary by F. Chyzak</i> .....	105
Effective Test of Local Algebraic Observability — Applications to Systems and Control Theory. <i>Talk by A. Sedoglavic</i> .....	113

### Part IV. Probabilistic Methods

Reflected Brownian Bridge Area Conditioned on its Local Time at the Origin. <i>Talk by G. Louchard, summary by M. Nguyễn-Thé</i> .....	117
Cover Time and Favourite Points for Planar Random Walks. <i>Talk by A. Dembo, summary by Ch. Fricker and P. Nicodème</i> .....	121
Introduction to Random Walks on Groups. <i>Talk by Y. Guivarc’h, summary by Ph. Robert</i> .....	127
Random Matrices and Queues in Series. <i>Talk by Y. Baryshnikov, summary by M. Durand</i> .....	131
Information Theory by Analytic Methods: The Precise Minimax Redundancy. <i>Talk by W. Szpankowski, summary by T. Klausner</i> .....	135

### Part V. Asymptotics and Analysis

On Jackson’s $q$ -Bessel Functions. <i>Talk by C. Zhang, summary by B. Salvy</i> .....	141
On the Convergence of Borel Approximants. <i>Talk by D. Lutz, summary by M. Durand</i> .....	145

**Part VI. ALEA'2001 Lecture Notes**

Enumerative Combinatorics: Combinatorial Decompositions and Functional Equations.  
*Talk by M. Bousquet-Mélou, summary by M. Vandenbogaert* ..... 151

Symbolic Enumerative Combinatorics and Complex Asymptotic Analysis.  
*Talk by Ph. Flajolet, summary by Y. Le Borgne*..... 161

Aléa discret et mouvement brownien (*Discrete Randomness and Brownian Motion*).  
*Talk by Ph. Chassaing, summary by Ph. Chassaing*..... 171









---

Unité de recherche INRIA Lorraine, Technopôle de Nancy-Brabois, Campus scientifique,  
615 rue du Jardin Botanique, BP 101, 54600 VILLERS LÈS NANCY  
Unité de recherche INRIA Rennes, Irisa, Campus universitaire de Beaulieu, 35042 RENNES Cedex  
Unité de recherche INRIA Rhône-Alpes, 655, avenue de l'Europe, 38330 MONTBONNOT ST MARTIN  
Unité de recherche INRIA Rocquencourt, Domaine de Voluceau, Rocquencourt, BP 105,  
78153 LE CHESNAY Cedex  
Unité de recherche INRIA Sophia-Antipolis, 2004 route des Lucioles, BP 93, 06902 SOPHIA-ANTIPOLIS  
Cedex

---

Éditeur  
INRIA, Domaine de Voluceau, Rocquencourt, BP 105, 78153 LE CHESNAY Cedex  
(France)  
<http://www.inria.fr>  
ISSN 0249-6399