

And/Or trees revisited

B. Chauvin ^{*}, P. Flajolet [†], D. Gardy [‡], B. Gittenberger [§]

Abstract

We consider boolean functions over n variables. Any such function can be represented (and computed) by a complete binary tree with and or or in the internal nodes and a literal in the external nodes, and many different trees can represent the same function, so that a fundamental question is related to the so-called *complexity* of a boolean function: $L(f) :=$ minimal size of a tree computing f .

The existence of a limiting probability distribution $P(\cdot)$ on the set of and/or trees was shown by Lefmann and Savicky [8]. We give here an alternative proof, which leads to effective computation in simple cases. We also consider the relationship between the probability $P(f)$ and the complexity $L(f)$ of a boolean function f . A detailed analysis of the functions enumerating some sub-families of trees, and of their radius of convergence, allows us to improve on the upper bound of $P(f)$, established by Lefmann and Savicky.

1 Introduction

Random And/Or boolean formulas and functions play an important rôle in the literature of theoretical computer science, and one of the fundamental questions they raise is that of their representation by a data structure such as a tree or a circuit.

^{*}LAMA, CNRS UMR 8100, Université de Versailles Saint-Quentin, 78035 Versailles Cedex, France.

[†]INRIA Rocquencourt, Domaine de Voluceau, 78153 Le Chesnay, France.

[‡]PRISM, CNRS FRE 2510, Université de Versailles Saint-Quentin, 78035 Versailles Cedex, France.

[§]Department of Geometry, Technische Universität Wien, Wiedner Hauptstraße 8-10/113, A-1040 Wien, Austria.

Some properties of a representation, such as its size, are actually properties of the associated boolean function, and estimating the usefulness of a given representation, e.g. its average size, quickly leads to investigating probability distributions on some space of boolean functions.

Several recent works have attempted to define probability distributions for suitable models of such functions. For instance, Friedman [6] investigated distributions involved in the random k -SAT problem, where iterated conjunctions of small disjunctions appear. A sequence of probability distributions defined on formulas of the same size and with the same tree structure appear in his model. Paris *et al.* [11] show that the proportion $pr_{n,k}$ of boolean formulas f on n variables which take the value **true** for k affectations of the variables (and the value **false** for the remaining $2^n - k$ affectations) has a well defined limit, when the number n of boolean variables and the formula size both tend to infinity; this amounts to defining the limiting probability distribution of the $pr_{n,k}$. Closely related to the present work, and indeed at its origin, are some papers by Savicky *et al.* Savicky and Lefmann [8] obtain a relation between the probability of a boolean function and its complexity. Further papers [15, 16] establish relationships between functions of given limited complexity and some probability distributions or some enumeration results. By contrast, earlier results of Savicky [13, 14] present a way of iteratively building boolean formulas, which leads to a limiting uniform distribution on the set of boolean functions on a fixed number of variables, and is closest in spirit to the work of Friedman [6]. Woods [17] presents, among other results relative to logical sentences, a model of boolean formulas which is closely related to our approach, both in the tools used (generating functions) and in the final result (existence of a limiting distribution on the set of boolean functions).

We consider in the present paper the model for boolean functions presented for example by Lefmann and Savicky in [8], where the probability of a boolean function is proportional to the number of boolean formulas of a given type which compute it. The boolean functions are defined on n variables x_1, x_2, \dots, x_n . Since the literals are then $x_1, \bar{x}_1, x_2, \bar{x}_2, \dots, x_n, \bar{x}_n$, there are 2^{2^n} such boolean functions. All along the paper, n is fixed, and we consider some special values in sections 2.2 to 2.4.

In such a context, formulas of size m with n variables are represented by labeled rooted binary trees where the m internal nodes are labeled by **and** and **or** and the $m + 1$ external nodes by a literal, i.e., a variable or its negation. Each of the m inner nodes is labeled by **and** or **or** with equal probability $1/2$ and independently of the other nodes; each leaf is labeled by a literal, chosen according to the uniform distribution on the $2n$ literals and independently of the labeling of all the other nodes. Many different trees can compute the same function, so that a fundamental

question is to evaluate the so-called *complexity* of a boolean function, which we define as

$$L(f) := \text{minimal size of a tree computing } f.$$

In this paper, we define the *size* of a binary tree as the number of its internal nodes¹. We should also mention that several different complexity measures for boolean functions have been proposed in the literature; see e.g. [2] for a recent survey.

The aim of this paper is dual: a better understanding of the limiting probability distribution on the space of boolean functions, and a study of the relationship between the probability of a given boolean function and its complexity. For the first topic, we need to make precise what we mean by *the probability of a given function*. As suggested by Woods [17] and further argued by Lefmann and Savicky [8], a natural definition of the limiting distribution is as the limit of the uniform distribution on **and/or** trees of finite size approaching infinity. More precisely, for any fixed m there is a uniform distribution P_m on the set of **and/or** trees with n variables and m internal nodes. Let f be some boolean function and define $t(f, m)$ as the number of trees of size m that compute f . Then the probability of this function is by definition

$$P_m(f) = \frac{t(f, m)}{T_m},$$

where T_m is the total number of **and/or** trees with size m (we give later on a simple formula (3) for T_m). For simplicity's sake, we use the same notation P_m for the distributions on trees and on functions, although the last one is actually the image probability of P_m by the canonical application which associates to a tree the boolean function it computes.

Lefmann and Savicky proved [8, Theorem 2.3] that these distributions P_m have a limit P when m goes to $+\infty$, which they describe as a biased tree distribution. Section 2 of this paper is devoted to an alternative description of this limiting distribution P and to some explicit computations for the cases $n = 1, 2$ or 3 using generating functions. The case of general n relies on results for systems of algebraic equations due to Drmota [3], Lalley [7] or Woods [17].

The combined approach by generating functions and branching processes also allows us to define a second probability distribution on boolean functions: starting from a critical branching process, we label at random its internal and external nodes

¹We can also choose the number of external nodes or the total number of nodes. In Lefmann and Savicky's paper [8], the size is the number of external nodes, so that there is a $+1$ shift when comparing our results to theirs.

to obtain a random **and/or** tree, i.e. a random boolean function. We present this approach in Section 2.6, together with some numerical computations and comparisons with our first probability distribution P .

The second topic, namely the study of the relationships between $P(f)$ and $L(f)$, improves on the relation proved by Lefmann and Savicky:

$$\frac{1}{4} \cdot \left(\frac{1}{8n}\right)^{L(f)} \leq P(f) \leq (1 + O(1/n)) \exp\left(-c \frac{L(f)}{n^3}\right). \quad (1)$$

The lower bound seems to be tight, but the upper bound can be improved, to yield an order n^{-2} instead of n^{-3} . Following Lefmann and Savicky, our method is to start from Markov inequality: for any function f and for $\varepsilon > 0$, the definition of the complexity gives (again using the same notation for both distributions P on trees and on functions, as we did for P_m)

$$\begin{aligned} P(f) = P(\text{the tree } \tau \text{ computes } f) &\leq P((1 + \varepsilon)^{\|\tau\|} \geq (1 + \varepsilon)^{L(f)}) \\ &\leq \frac{E[(1 + \varepsilon)^{\|\tau\|}]}{(1 + \varepsilon)^{L(f)}}. \end{aligned}$$

The upper bound in (1) can be improved as ε becomes larger. That means a better control of $E[(1 + \varepsilon)^{\|\tau\|}]$, in other words a fine evaluation of the radius of convergence of the generating function of the size of a tree. This is achieved in Section 3 where the following theorem is proved.

Theorem 1 *Almost surely,*

$$\frac{1}{4} \cdot \left(\frac{1}{8n}\right)^{L(f)} \leq P(f) \leq (1 + O(1/n)) \exp\left(-c \frac{L(f)}{n^2}\right). \quad (2)$$

Finally we discuss our results, both on the improvement on Lefmann and Savicky's bound for the complexity and on the probability distributions, in Section 4, where we also consider possible extensions to other models of boolean formulae, which would take into account the commutativity and associativity of the boolean operators.

2 Enumerating functions and limit distribution for **and/or** trees

We recall that the generating function for binary trees, counted by the number of (internal and external) nodes, satisfies the equation

$$b(z) = 1 + z b(z)^2,$$

which gives

$$b(z) = \frac{1 - \sqrt{1 - 4z}}{2z}.$$

Now define the set \mathcal{T} of **and/or** trees, assuming that the number of variables is n :

$$\mathcal{T} = \oplus_{1 \leq i \leq n} (\{x_i\} + \{\bar{x}_i\}) \oplus (\wedge, \mathcal{T}, \mathcal{T}) \oplus (\vee, \mathcal{T}, \mathcal{T});$$

hence the equation on the generating function for these trees enumerated by number of *internal* nodes:

$$T(z) = 2n + 2zT(z)^2,$$

which gives

$$T(z) = \frac{1 - \sqrt{1 - 16nz}}{4z}.$$

This gives readily the number of **and/or** trees with m internal nodes (a formula already given in former papers [8, 11]):

$$T_m := [z^m]T(z) = 2^m(2n)^{m+1}C_m, \tag{3}$$

with C_m the Catalan number: $C_m = (2m)!/m!(m+1)!$.

2.1 From the distribution on trees to the distribution on boolean functions

We assume in this section that the probability distribution over **and/or** trees on n variables and of size m (i.e. number of internal nodes, or number of leaves minus 1) is uniform. Let us stress again that this distribution depends on m , and that throughout this section n is a fixed parameter. This induces a probability distribution P_m over the boolean functions on n variables: Let $t(f, m)$ be the number of trees of size m that compute a given boolean function f^2 ; then the probability of this function is

$$P_m(f) = t(f, m)/T_m.$$

Assume that we know the generating function $t_f(z) = \sum_m t(f, m)z^m$ enumerating the trees that compute the function f ; then the probability $P_m(f)$ is simply

$$P_m(f) = \frac{[z^m]t_f(z)}{[z^m]T(z)}.$$

²Both the numbers of trees T_m and $t(f, m)$ and the distribution $P_m(f)$ depend on the parameter n , which we do not mention explicitly unless necessary.

Hence knowing the asymptotic behaviour, as m tends to infinity, of the coefficients of the functions t_f (and of T) will give us the existence of a limiting distribution on the boolean functions and possibly a way of computing it. Now Lefmann and Savicky [8, Thm. 2.3] simply assert the existence of this limiting distribution:

$$P(f) = \lim_{m \rightarrow +\infty} P_m(f).$$

In the sequel, we first examine how we can explicitly compute the limiting distribution for $n = 1..3$, before turning to the case of general n .

Remark: One might wish to study a different probability distribution $P_{\leq m}(f)$, defined as the ratio of the number of trees of size *smaller than or equal to* m that compute the boolean function f , over the total number of trees of size smaller than or equal to m . Such an approach has the advantage that the supports of the probability distributions $P_{\leq m}$, for increasing m , are increasing subsets of the set of all (finite or infinite) binary trees labelled by \wedge , \vee and literals. By the Kolmogorov existence theorem [1, Sect. 36], we know the existence of a limiting probability distribution on the set of *finite and/or* trees. Using generating functions, we can write

$$P_{\leq m}(f) = \frac{[z^m]\{t_f(z)/(1-z)\}}{[z^m]\{T(z)/(1-z)\}}.$$

Assume that the function $t_f(z)$ has a radius of convergence $1/16n$ (see below); as the function $T(z)$ has the same radius of convergence $1/16n$, we see that dividing by $1-z$ introduces, in both cases, a singularity at 1, larger than the radius of convergence; hence the asymptotic behaviour is determined by the singularities at $1/16n$ and the asymptotic limit of $P_{\leq m}(f)$ is exactly the limit of $P_m(f)$ (although the values for finite m differ).

2.2 Case of a single variable

In this part, we consider what happens when there is a single variable x , and two literals x and \bar{x} . There are four functions:

$$f_1 = False; \quad f_2 = \bar{x}; \quad f_3 = x; \quad f_4 = True.$$

Let us denote by A_f the set of trees computing the boolean function f . We have that

$$\begin{aligned} A_{True} &= (\wedge, A_{True}, A_{True}) \oplus (\vee, A_x, A_{\bar{x}}) \oplus (\vee, A_{\bar{x}}, A_x) \\ &\quad \oplus (\vee, A_{True}, A) \oplus (\vee, A, A_{True}) \setminus (\vee, A_{True}, A_{True}). \end{aligned}$$

The subtraction of the last term comes from the fact that the two preceding terms both contain the trees $(\vee, A_{True}, A_{True})$, which we must count only once. We get an equation on the generating functions, where t_f is the enumerating g.f. for the trees that compute the boolean function f , and where T is, as above, the enumerating function for all trees:

$$t_{True}(z) = 2zt_x(z)t_{\bar{x}}(z) + 2zt_{True}(z)T(z). \quad (4)$$

By symmetry (exchange \vee with \wedge and $True$ with $False$, in the equation defining A_{True} , to get the equation on A_{False}),

$$t_{False}(z) = 2zt_x(z)t_{\bar{x}}(z) + 2zt_{False}(z)T(z). \quad (5)$$

Now consider the set of trees that compute the function $f_3 = x$:

$$\begin{aligned} A_x = & \{x\} \oplus (\wedge, A_x, A_x) \oplus (\wedge, A_x, A_{True}) \oplus (\wedge, A_{True}, A_x) \\ & \oplus (\vee, A_x, A_{False}) \oplus (\vee, A_{False}, A_x) \oplus (\vee, A_x, A_x) \end{aligned}$$

Hence the equation on the generating functions:

$$t_x(z) = 1 + 2zt_x(z)^2 + 2zt_x(z)t_{True}(z) + 2zt_x(z)t_{False}(z). \quad (6)$$

By symmetry, we get a similar equation for the function $f_2 = \bar{x}$:

$$t_{\bar{x}}(z) = 1 + 2zt_{\bar{x}}(z)^2 + 2zt_{\bar{x}}(z)t_{True}(z) + 2zt_{\bar{x}}(z)t_{False}(z). \quad (7)$$

Now we solve this system of four equations in four variables, t_{True} , t_{False} , t_x and $t_{\bar{x}}$, to get first (obvious) that $t_{True} = t_{False}$ and $t_x = t_{\bar{x}}$, then that

$$t_{True}(z) = \frac{2zt_x^2(z)}{1 - 2zT(z)},$$

and finally a polynomial equation on the function $t_x(z)$:

$$1 + 2zy^2 - y + 8\frac{z^2y^3}{1 - 2zT(z)} = 0,$$

with $T(z) = (1 - \sqrt{1 - 16z})/4z$. Solving, we get three solutions for this equation:

$$\frac{2}{1 - \sqrt{1 - 16z}}; \quad \frac{-1}{8z} \left(1 + \sqrt{1 - 16z} \pm \sqrt{2 + 16z + 2\sqrt{1 - 16z}} \right).$$

The solution $t_x(z)$ is the (unique) function such that its value at $z = 0$ exists and is equal to 1; hence

$$t_x(z) = \frac{-1}{8z} \left(1 + \sqrt{1 - 16z} - \sqrt{2 + 16z + 2\sqrt{1 - 16z}} \right).$$

Expanding this expression around the singularity $z_0 = 1/16$, we get

$$\begin{aligned} t_x(z) &= 2(\sqrt{3} - 1) + 2 \left(\frac{1}{\sqrt{3}} - 1 \right) \sqrt{1 - 16z} \\ &\quad + 2 \left(\frac{7\sqrt{3}}{9} - 1 \right) (1 - 16z) + O((1 - 16z)^{3/2}), \end{aligned}$$

and a transfer lemma [4] gives readily

$$t(x, m) = [z^m]t_x(z) \sim 2^{2m+2} \frac{\sqrt{3} - 1}{\sqrt{3}} C_{m-1}.$$

We finally obtain the asymptotic probability of the function x by dividing the number $t(x, m)$ of trees computing this function by the total number of trees T_m :

$$P_m(x) \sim \frac{\sqrt{3} - 1}{\sqrt{3}} \cdot \frac{m + 1}{2m - 1} \rightarrow \frac{\sqrt{3} - 1}{2\sqrt{3}} = 0.21132486....$$

Now the function enumerating the trees that compute the function *True* is

$$t_{True}(z) = \frac{1}{8z} \left(2 - \sqrt{2 + 16z + 2\sqrt{1 - 16z}} \right),$$

which gives for $P_m(True)$ an asymptotic probability equal to $1/2\sqrt{3} = 0.28867513...$

In view of future generalization, let us look again at the initial system of four equations: We can rewrite each equation in a standard form

$$t_f = 1_{f \text{ literal}} + z^t F(A_{f,\vee} + B_{f,\wedge})F,$$

where F is the vector $(t_{False}, t_x, t_{\bar{x}}, t_{True})$, and the matrices $A_{f,\vee}$ and $B_{f,\wedge}$ are obtained by a process to be described in Section 2.5, and are given below for the

functions $True$ and x (the other cases are symmetrical):

$$\begin{aligned} B_{True, \wedge} &= \begin{pmatrix} 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 \end{pmatrix} & A_{True, \vee} &= \begin{pmatrix} 0 & 0 & 0 & 1 \\ 0 & 0 & 1 & 1 \\ 0 & 1 & 0 & 1 \\ 1 & 1 & 1 & 1 \end{pmatrix} \\ B_{x, \wedge} &= \begin{pmatrix} 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 1 \\ 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \end{pmatrix} & A_{x, \vee} &= \begin{pmatrix} 0 & 1 & 0 & 0 \\ 1 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{pmatrix} \end{aligned}$$

2.3 Functions on two variables

We have now two variables x_1 and x_2 , four literals and sixteen boolean functions. Writing down the recurrence relations on the sets A_f and translating them into generating functions gives us a system of sixteen algebraic equations, each of degree two. Symmetries (the generating function for a boolean function f is equal to the generating function for $\neg f$; the variables x_1 and x_2 can be exchanged; the generating functions for the boolean functions $l_1 \wedge l_2$ and $l_1 \vee l_2$, where the l_i are the literals on x_1 and x_2 , are the same) reduce it to a system of order 4, where a, b, c and d are the generating functions respectively for the boolean functions $True$, x_1 , $x_1 \wedge x_2$ and $x_1 \text{ xor } x_2$ ³:

$$\begin{cases} a = 2zaT + 4zb^2 + 20zc^2 + 2zd^2 + 16zbc + 8zcd; \\ b = 1 + 2zb^2 + 4zc^2 + 4zab + 8zbc; \\ c = 2zb^2 + 8zc^2 + 4zac + 8zbc + 4zbd + 4zcd; \\ d = 4zc^2 + 2zd^2 + 4zad + zcd. \end{cases}$$

Furthermore, we have that

$$T = 2a + 4b + 8c + 2d. \quad (8)$$

The solution of this system will give us closed-form expressions of the functions, which can be readily expanded around the singularity $z = 1/32$, and finally exact and approximate expressions for the probabilities. We leave the details to the reader and give for example the function $d(z)$:

$$d(z) = \frac{1}{8z} \left(-1 - \tau_0 + 2\tau_1 - \sqrt{\tau_2 + \sqrt{8\tau_3}} \right),$$

³We denote by $x_1 \text{ xor } x_2$ the boolean function $(x_1 \wedge \bar{x}_2) \vee (\bar{x}_1 \wedge x_2)$. This is simply the $+$ operator on $\{0, 1\}$.

with $\tau_0 = \sqrt{1-32z}$, $\tau_1 = \sqrt{2+2\tau_0}$ and

$$\begin{aligned}\tau_2 &= 2(1+\tau_0)(2-\tau_1) - 32z; \\ \tau_3 &= (1+\tau_0)(5-2\tau_1) + 16z(-3+2\tau_0+\tau_1-\tau_0\tau_1) - 128z^2.\end{aligned}$$

Now, if we have an expansion $t_f = \alpha_f - \beta_f \sqrt{1-z/\rho} + O(z-\rho)$ with $\beta_f > 0$ at the singularity $\rho = 1/32$, then the probability of the function f is (the details are given in Section 2.5)

$$P(f) = \frac{\beta_f}{T(\rho)} = \beta_f/8.$$

Define

$$\begin{aligned}\alpha &= (2\sqrt{2}-1)^2 = 9-4\sqrt{2}; \\ \beta &= -129+90\sqrt{2}+61\sqrt{3}-38\sqrt{6}; \\ \gamma^2 &= (\sqrt{3}-1)(2\sqrt{2}+\sqrt{3}) = 3-2\sqrt{2}-\sqrt{3}+2\sqrt{6}; \\ \delta &= \sqrt{6-2\sqrt{2}+\sqrt{3}-2\sqrt{6}} = \frac{1}{\sqrt{2}}(2\sqrt{2}-1-\sqrt{3}); \\ \nu &= 153-117\sqrt{2}+61\sqrt{3}-38\sqrt{6}.\end{aligned}$$

We obtain the following asymptotic probabilities:

$$\begin{aligned}P(True) &= \frac{\beta}{6\alpha\gamma\sqrt{2}}; \\ P(x_1) &= 1 - \frac{3}{2\sqrt{2}} + P(True) - \frac{\nu}{6\alpha\delta} = P(True) - \frac{(\sqrt{2}-1)^2}{2\sqrt{2}} - \frac{\nu}{6\alpha\delta}; \\ P(x_1 \wedge x_2) &= \frac{\sqrt{2}-1}{2} - P(True) + \frac{\nu}{12\alpha\delta} = \frac{\nu}{12\alpha\delta} - P(x_1 \text{ xor } x_2); \\ P(x_1 \text{ xor } x_2) &= P(True) - \frac{\sqrt{2}-1}{2}.\end{aligned}$$

Floating values are easy to compute:

$$\begin{aligned}P(True) &= .20940201\dots; & P(x_1) &= .06717345\dots; \\ P(x_1 \wedge x_2) &= .03848896\dots; & P(x_1 \text{ xor } x_2) &= .00229522\dots\end{aligned}$$

In other terms, a random boolean function is one of the constant functions (**True** or **False**) almost 42% of the time, a literal 27% of the time, a function of the kind $l_1 \wedge l_2$ in 30% of the cases, and either $x_1 \text{ xor } x_2$ or its negation less than .5% of the times. The average complexity of a random boolean function under this probability distribution is $2/\sqrt{3} - \sqrt{2} + 1 = 0.740486\dots$

2.4 Functions on three variables

We consider here the case where $n = 3$. There are fourteen different classes of boolean functions; in each class the same generating function enumerates the binary trees associated to the boolean functions. Let us denote these fourteen generating functions by the column vector $\mathbf{t}(z) = {}^t(t_1(z), \dots, t_{14}(z))$. Then $\mathbf{t}(z)$ satisfies a functional equation of the form

$$\mathbf{t}(z) = \mathbf{Q}(\mathbf{t}(z)), \quad (9)$$

where each component of the vector-valued function \mathbf{Q} is quadratic in each of the t_i (for details see section 2.5). By the Drmota-Lalley-Woods theorem (see again section 2.5) we know that each of the functions t_i admits a representation of the form

$$t_i(z) \sim \alpha_i - \beta_i \sqrt{1 - 48z}, \quad z \rightarrow \frac{1}{48}$$

Thus for $z = 1/48$ the system (9) has a unique solution $(\alpha_1, \dots, \alpha_{14})$. Hence the fixed point can be obtained by iteration, starting from a vector whose coordinates are all equal to zero.

In order to compute the values β_i , $i = 1, \dots, 14$, observe that Drmota [3] showed that the vector $(\beta_i)_{i=1, \dots, 14}$ is an eigenvector with eigenvalue 1 of the Jacobian

$$\frac{\partial \mathbf{Q}}{\partial \mathbf{t}} = \left(\frac{\partial Q_i}{\partial t_j} \right)_{i,j=1 \dots 14}$$

evaluated at $z = 1/48$. Since 1 is an eigenvalue of multiplicity 1 at $z = 1/48$, we can easily compute the eigenvector and normalize it to obtain the results presented below. We give, for each class, the generic form of boolean functions belonging to it ($l_i \in \{x_i, \bar{x}_i\}$ for $i = 1..3$; other functions of the class are obtained by permuting literals or exchanging \vee and \wedge), its cardinality (number of boolean functions), the complexity and probability common to all the functions of the class, and finally the cumulated probability of the class. The classes are in decreasing order of individual probability. The values given below were obtained with 30.000 iteration steps, and rounded to three digits.

Boolean Function	Card.	Compl.	Indiv. Prob.	Cumul. Prob.
True	2	1	0.165	0.330
l_1	6	0	0.0314	0.188
$l_1 \wedge l_2$	24	1	0.00995	0.239
$l_1 \wedge l_2 \wedge l_3$	16	2	0.00768	0.123
$(l_1 \wedge l_2) \vee l_3$	48	2	0.00211	0.101
$(l_1 \wedge l_3) \vee (\bar{l}_1 \wedge l_2)$	24	3	0.28710^{-3}	0.00689
$l_1 \text{ xor } l_2$	6	3	0.19210^{-3}	0.00115
$(l_1 \text{ xor } l_2) \vee l_3$	24	4	0.15710^{-3}	0.00377
$(l_1 \wedge (l_2 \vee l_3)) \vee (l_1 \wedge l_2)$	8	4	0.14910^{-3}	0.00119
$(l_1 \wedge l_2 \wedge l_3) \vee (\bar{l}_1 \wedge \bar{l}_2)$	48	4	0.96210^{-4}	0.00462
$(l_1 \wedge l_2 \wedge l_3) \vee (\bar{l}_1 \wedge \bar{l}_2 \wedge \bar{l}_3)$	8	5	0.56010^{-4}	0.44810^{-3}
$(l_1 \wedge (l_2 \vee l_3)) \vee (\bar{l}_1 \wedge \bar{l}_2 \wedge \bar{l}_3)$	24	5	0.21710^{-4}	0.52110^{-3}
$(l_1 \wedge (l_2 \text{ xor } \bar{l}_3)) \vee (\bar{l}_1 \wedge (l_2 \vee \bar{l}_3))$	16	7	0.27910^{-5}	0.44610^{-4}
$(l_1 \text{ xor } l_2) \text{ xor } l_3$	2	9	0.81410^{-7}	0.16310^{-6}

Again, we notice that this model gives a predominant place to constants and to very simple functions: a function has a probability 0.330 to be constant and 0.757 to have complexity 0 or 1; the functions of complexity 2 have a global 0.224 probability; the probability that the complexity is equal to 3 drops to 0.008, and the cumulated probability of functions with complexity 4 or larger is 0.0106. The average complexity of a random boolean function under this model is equal to 1.08.

2.5 General case

For each boolean function f on n variables we can write an equation on the set A_f of trees computing it:

$$A_f = 1_{\{f \text{ literal}\}} \oplus \sum_{g,h:f=g \vee h} (\vee, A_g, A_h) \oplus \sum_{g,h:f=g \wedge h} (\wedge, A_g, A_h).$$

This equation on sets of trees translates into an equation on the generating functions enumerating these sets. We obtain

$$t_f(z) = 1_{\{f \text{ literal}\}} + z \sum_{g,h:g \vee h=f} t_g(z) t_h(z) + z \sum_{g,h:g \wedge h=f} t_g(z) t_h(z).. \quad (10)$$

The 2^{2^n} boolean functions on n variables can be defined by their truth table: We associate to each function f a word $(f[1], \dots, f[p])$ of length $p = 2^n$ over the alphabet

$\{0, 1\}$, representing its value for each of the p assignments of values 0 or 1 to the n variables x_i , i.e. a column of the truth table. The alphabetical order of the words also gives an ordering on the functions of the set \mathcal{F} of boolean functions, which we denote then by f_1, \dots, f_{2^p} : $f_1 = (0, \dots, 0) = \text{False}$, $f_2 = (0, \dots, 0, 1) = \bar{x}_1 \wedge \bar{x}_2 \wedge \dots \wedge \bar{x}_n$, and $\neg f_i = f_{2^p-i}$.

Now the relations $f = g \vee h$ and $f = g \wedge h$ translate into relations on the corresponding words: For each j , we have that $f[j] = g[j] + h[j]$ (for \vee , with $1+1 = 1$) or $f[j] = g[j].h[j]$ (for \wedge). For example, for $n = 3$, a boolean function is defined by a word of $\{0, 1\}^8$; the function $f = x_1 \vee x_3$ can equivalently be defined by the word $(0, 1, 0, 1, 1, 1, 1, 1)$; let $g = (0, 0, 0, 1, 0, 0, 0, 1)$ ($g = x_2 \wedge x_3$); then the functions h such that $f = g \vee h$ are all the functions $(0, 1, 0, *, 1, 1, 1, *)$, where $*$ stands for 0 or 1, which gives us four possible functions.

Define the vector $\mathbf{f} = (t_{f_1}, \dots, t_{f_{2^p}})$ of the generating functions; furthermore define, for each boolean function f_k , two matrices

$$A_{f_k} = (a_{i,j}(f_k)), \quad B_{f_k} = (b_{i,j}(f_k)),$$

with $a_{i,j}(f_k) = 1$ if $f_k = f_i \vee f_j$ and 0 otherwise, and with $b_{i,j}(f_k) = 1$ if $f_k = f_i \wedge f_j$ and 0 otherwise. Then we can write the equation (10) as

$$t_f(z) = 1_{\{f \text{ literal}\}} + z {}^t\mathbf{f} \cdot (A_f + B_f) \cdot \mathbf{f} \quad (11)$$

$$= 1_{\{f \text{ literal}\}} + z \sum_{1 \leq i, j \leq 2^p} (a_{i,j} + b_{i,j}) t_{f_i}(z) t_{f_j}(z), \quad (12)$$

where ${}^t\mathbf{f}$ is the vector obtained from \mathbf{f} by transposition. Now we have such an equation for each of the 2^p functions f , which gives a system of 2^p algebraic equations on 2^p unknown functions f_i .

Theorem 2 *The limiting probability distribution of $P_m(f)$, $m \rightarrow +\infty$, exists and can be computed.*

Proof

If we instantiate the equation (11) for each of the f_i , we obtain a nonlinear polynomial system $\vec{y} = \Phi(\vec{y})$, where each component has nonnegative coefficients, and such that the dependency graph is connected and that the system is a-proper (i.e. a certain Lipschitz condition is satisfied). Then, by results on systems of algebraic equations (see Drmota [3], Lalley [7] or Woods [17]), all component solutions are algebraic with a common singularity, and we can expand the functions around their singularity to get the asymptotic behaviour of the probabilities $P_m(f)$ for large m . We give below the version due to Flajolet and Sedgewick [5, Th. 8.13, p. 71].

Positive polynomial systems. Consider a nonlinear polynomial system $\vec{y} = \Phi(\vec{y})$ that is *a-proper*, *a-positive* and *a-irreducible*. In that case, all component solutions y_j have the same radius of convergence $\rho < \infty$. Then, there exist functions h_j analytic at the origin such that

$$y_j = h_j \left(\sqrt{1 - z/\rho} \right) \quad (z \rightarrow \rho^-) \quad (1 \leq j \leq 2^p).$$

In addition, all other dominant singularities are of the form $\rho\omega$ with ω a root of unity. If furthermore the system is *a-aperiodic*, all y_j have ρ as unique dominant singularity. In that case, the coefficients admit a complete asymptotic expansion of the form

$$[z^n]y_j(z) \sim \rho^{-n} \left(\sum_{k \geq 1} d_k n^{-1-k/2} \right).$$

In our case, it is easy to check that the system is non linear, *a-proper*, *a-positive*, *a-irreducible* and *a-periodic* (the definitions come again from [5, Th. 8.13, p. 71]). We can then apply the theorem: There exists a solution $(t_{f_1}, \dots, t_{f_{2^p}})$ to the algebraic system; the t_f have a common, strictly positive, radius of convergence ρ and a unique dominant algebraic singularity at $\rho < +\infty$, with an expansion around ρ

$$t_f(z) = \alpha_f - \beta_f \sqrt{1 - z/\rho} + O(1 - z/\rho), \quad (13)$$

which gives by a transfer lemma [4]

$$t(f, m) = [z^m]t_f(z) = -\beta_f [z^m] \sqrt{1 - z/\rho} (1 + O(1/m)).$$

The radius of convergence is also the radius for the function $T(z) = \sum_f t_f(z)$; hence $\rho = 1/16n$ and

$$t(f, m) = \beta_f C_{m-1} \rho^{-m} 2^{1-2m} (1 + O(1/m)).$$

Now we have that $P_m(f) = t(f, m)/T_m$; plugging in the value of T_m from equation (3), we get

$$P_m(f) = \frac{\beta_f}{4n} (1 + O(1/m)),$$

which provides another proof of the existence of the asymptotic distribution P . \square

Numerical computation of the asymptotic probabilities? The functions t_f are defined at ρ ; plugging the values (13) for each f into the equations (10) and identifying in each equation the coefficient of $\sqrt{1 - z/\rho}$, we get a system of size 2^{p+1} on the α_f and β_f , whose generic equations are

$$\begin{cases} \alpha_f &= 1_{\{f \text{ literal}\}} + \rho \sum_{g \vee h = f} \alpha_g \alpha_h + \rho \sum_{g \wedge h = f} \alpha_g \alpha_h; \\ \beta_f &= \rho \sum_{g \vee h = f} (\alpha_g \beta_h + \alpha_h \beta_g) + \rho \sum_{g \wedge h = f} (\alpha_g \beta_h + \alpha_h \beta_g). \end{cases}$$

Now this system can be solved numerically with the help of a Computer Algebra System such as Maple; see the indications given in Section 2.4 for the special case $n = 3$.

2.6 Yet another probability distribution

The terms α_f that appear in the expansion (13) can be interpreted to give rise to a different probability distribution π on the space of boolean functions.

Proposition 1 *The probability of computing the function f , when we start from a simple critical branching process (the probabilities that a node has 0 or 2 sons are both equal to $1/2$) and label the nodes randomly and independantly, to obtain a random and/or tree, is*

$$\pi(f) = \frac{t_f(\rho)}{T(\rho)} = \frac{\alpha_f}{4n}.$$

Proof: Let τ be a random and/or tree obtained by considering a simple critical Galton Watson process and labelling at random the nodes: We label the internal nodes by \wedge and \vee with equal probability, and the leaves by the $2n$ literals, again with equal probability. Let $\tilde{\tau}$ be the underlying unlabelled tree; almost surely $\tilde{\tau}$ is finite. Denote by $|\tilde{\tau}|$ or $|\tau|$ its size (total number of nodes), and by $||\tau||$ the number of its internal nodes: $|\tau| = 2||\tau|| + 1$. In this critical Galton Watson process, the probability that the process stops at one node and the probability that the node has two sons are both equal to $1/2$; hence the probability that we obtain an unlabelled tree $\tilde{\tau}$ of size $|\tilde{\tau}|$ is

$$Proba(\tilde{\tau}) = \frac{1}{2^{|\tilde{\tau}|}}.$$

Now the probability of a given labelled tree τ is obtained by multiplying the probability of $\tilde{\tau}$ by the probability of the labelling, which is itself equal to the probability $2^{-||\tau||}$ of labelling the internal nodes as in τ , times the probability $(1/2n)^{||\tau||+1}$ of labelling the leaves:

$$Proba(\tau) = \frac{1}{2^{|\tau|}} \cdot \frac{1}{2^{||\tau||}} \cdot \frac{1}{(2n)^{||\tau||+1}}.$$

Define $\pi(f) := \sum_{\tau \text{ computes } f} \text{Proba}(\tau)$; we have that

$$\begin{aligned} \pi(f) &= \frac{1}{4n} \sum_{\tau \text{ computes } f} \left(\frac{1}{16n} \right)^{||\tau||} \\ &= \frac{1}{4n} \sum_m t(f, m) \left(\frac{1}{16n} \right)^m \\ &= \frac{t_f(1/16n)}{T(1/16n)} = \frac{\alpha_f}{4n}. \end{aligned}$$

It is interesting to compute the distribution $\{\pi(f)\}$ for small n and to compare it to the distribution $\{P(f)\}$. Numerical data suggest that the distribution $\pi(f)$ is even more strongly biased towards functions of low complexity, mostly literals, than the distribution P , and that the average complexity of a random boolean function under $\pi(f)$ is less than half its average complexity under P . For example, $n = 1$ leads to

$$\begin{aligned} \pi(True) &= \frac{2-\sqrt{3}}{2} = 0.1339745960... \\ \pi(x) &= \frac{\sqrt{3}-1}{2} = 0.3660254040... \end{aligned}$$

The average complexity under this distribution is $2 - \sqrt{3} = 0.268$, versus $1/\sqrt{3} = 0.577$ under the distribution P .

For $n = 2$, we obtain (as in Section 2.3, we use $\gamma = \sqrt{3 - 2\sqrt{2} + 2\sqrt{6} - \sqrt{3}}$ and $\delta = (2\sqrt{2} - 1 - \sqrt{3})/\sqrt{2}$):

$$\begin{aligned} \pi(True) &= 1 - \frac{\gamma}{2} = 0.0864216570... \\ \pi(x_1) &= \frac{3}{\sqrt{2}} - 1 - \frac{\gamma}{2} - \frac{\delta}{\sqrt{2}} = 0.1595538420... \\ \pi(x_1 \wedge x_2) &= \frac{1+\gamma}{2} - \sqrt{2} + \frac{\delta}{2\sqrt{2}} = 0.0234588600... \\ \pi(x_1 \text{ xor } x_2) &= \sqrt{2} - \frac{1+\gamma}{2} = 0.000635219... \end{aligned}$$

The average complexity is equal now to $1 + 2(\sqrt{2} + \sqrt{3}) = 0.364$, to be contrasted with 0.740, the average complexity under $P(\cdot)$.

Finally, the case $n = 3$ gives the following results, where for comparison purposes we give both distributions $P(f)$ and $\pi(f)$, and where the weighted probabilities (for π) are relative to the cumulative values of boolean functions in the fourteen classes:

Boolean Function	$P(f)$	$\pi(f)$	Cumul. prob. for $\pi(f)$
True	0.165	0.0642	0.128
l_1	0.0314	0.0994	0.596
$l_1 \wedge l_2$	0.00995	0.00776	0.186
$l_1 \wedge l_2 \wedge l_3$	0.00768	0.00282	0.0451
$(l_1 \wedge l_2) \vee l_3$	0.00211	0.81710^{-3}	0.0392
$(l_1 \wedge l_2) \vee (\bar{l}_1 \wedge l_3)$	0.28710^{-3}	0.88010^{-4}	0.00211
$l_1 \text{ xor } l_2$	0.19210^{-3}	0.67310^{-4}	0.40410^{-3}
$(l_1 \text{ xor } l_2) \vee l_3$	0.15710^{-3}	0.31410^{-4}	0.75410^{-3}
$(l_1 \wedge (l_2 \vee l_3)) \vee (l_1 \wedge l_2)$	0.14910^{-3}	0.32110^{-4}	0.25710^{-3}
$(l_1 \wedge l_2 \wedge l_3) \vee (\bar{l}_1 \wedge \bar{l}_2)$	0.96210^{-4}	0.22010^{-4}	0.00106
$(l_1 \wedge l_2 \wedge l_3) \vee (\bar{l}_1 \wedge \bar{l}_2 \wedge \bar{l}_3)$	0.56010^{-4}	0.99910^{-5}	0.79910^{-4}
$(l_1 \wedge (l_2 \vee l_3)) \vee (\bar{l}_1 \wedge \bar{l}_2 \wedge \bar{l}_3)$	0.21710^{-4}	0.37010^{-5}	0.88810^{-4}
$(l_1 \wedge (l_2 \text{ xor } \bar{l}_3)) \vee (\bar{l}_1 \wedge (l_2 \vee \bar{l}_3))$	0.27910^{-5}	0.35410^{-6}	0.56610^{-5}
$(l_1 \text{ xor } l_2) \text{ xor } l_3$	0.81410^{-7}	0.76710^{-7}	0.15310^{-6}

We see that the distribution $\pi(f)$ leads to literals almost 60% of the time, to functions of complexity 1 in more than 31% cases, of complexity 2 less than 9%, and that functions of complexity 3 or larger are less than 0.5%. Accordingly, the average complexity under this distribution is 0.4998, again less than half the average complexity 1.086 under the distribution $P(.)$.

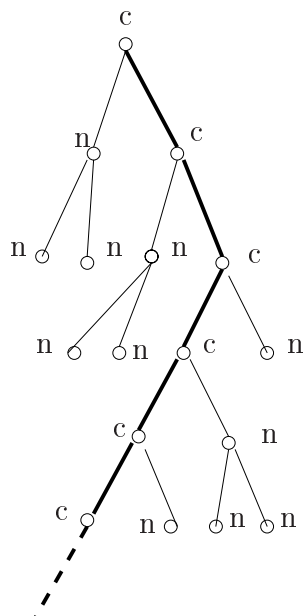
Comparing the two sets of values $P(f)$ and $\pi(f)$ for $n = 1...3$, we see that, apart from the interversion of literals and constants, the relative order of boolean functions, sorted in decreasing probability, is the same for the two distributions. Moreover, with the exception of literals, we always have that $\pi(f) < P(f)$. It would be desirable to give a rigorous proof of this fact for all n , which is equivalent to $\alpha_f < \beta_f$, with the exception of literals where $\alpha_l > \beta_l$.

3 Improving the bounds of Lefmann and Savicky

3.1 Lefmann and Savicky's model

In [8, th 2.3], Lefmann and Savicky obtain a limiting distribution P , which is both the limiting distribution of the uniform probability on finite trees of given size P_m , when m goes to infinity, and the limiting distribution when k goes to infinity, of some probabilities p_k constructed by the machinery of segments. This second representation leads to the following description of the limiting distribution P , as a

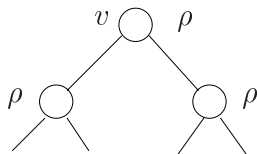
pruned infinite binary biased tree (the terminology comes from the branching literature, see [9] for instance). Start from a binary biased tree, in which there are two types of nodes: c -nodes and n -nodes. All the nodes on the spine of the biased tree (including the root), are c -nodes which reproduce always with two children, a c -node and a n -node. The nodes that don't belong to the spine are n -nodes, and they split having no descendant with probability $1/2$ and having two n -nodes children with probability $1/2$ (note that these critical branching subtrees are a.s. finite and are those considered in Section 2.6).



We then define the pruning as follows: For a given **and/or** tree, let us mark the internal nodes by conditions, with the following inductive procedure:

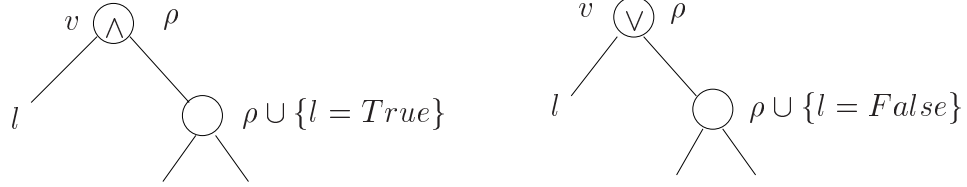
* if node v is the root, it is marked by a set ρ_0 of conditions; ρ_0 can be the empty set \emptyset , as it happens further.

* if node v , marked by a set of conditions ρ , has two internal nodes as children, then ρ is inherited by the children;

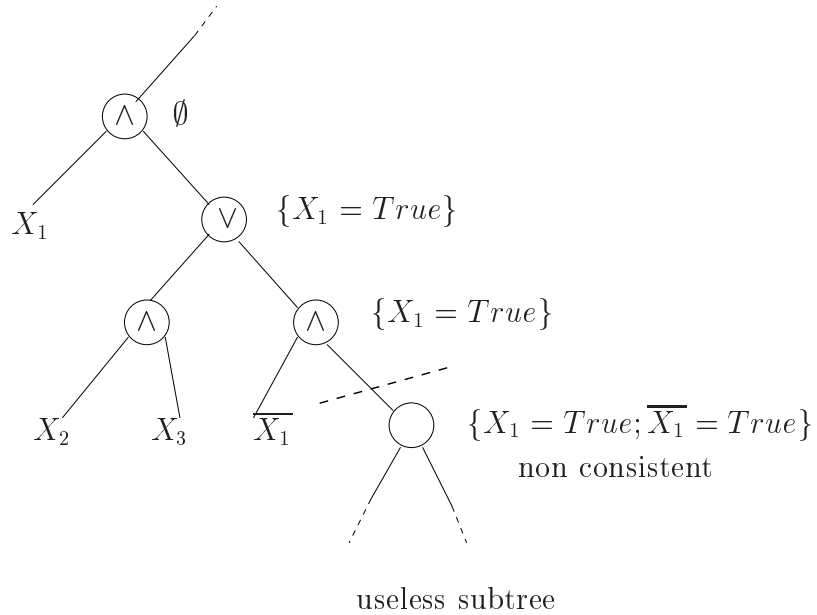


* if node v , marked by a set of conditions ρ , has for children an internal node and an external node containing a literal l , then

- if node v is labelled by a “and”, then the internal child node is marked by $\rho \cup \{l = True\}$;
- if node v is labelled by a “or”, then the internal child node is marked by $\rho \cup \{l = False\}$.



When applying this procedure to the whole tree, maybe some sets of conditions associated to some nodes are not consistent. In this case, the subtree beginning at such a node does not influence the boolean function and thus it can be replaced by any constant. By pruning the tree, such a subtree is deleted.



After pruning, standard simplifications rules (e.g. $f \wedge f = f$) can also be applied to get a smaller marked tree. Finally, a tree τ gives a pruned tree $\hat{\tau}$, then a simplified tree $\tilde{\tau}$, and these three trees compute the same function $f \in \mathcal{F}$, so that there is an upper bound for the complexity of f :

$$L(f) \leq \|\tilde{\tau}\| \leq \|\hat{\tau}\| .$$

In the following, for simplicity, we denote by τ (instead of $\tilde{\tau}$) a pruned, simplified tree and $\|\tau\|$ is the number of internal nodes of the tree τ associated with a consistent set of conditions. The Markov inequality then gives a way of estimating the complexity: for any function f ,

$$\begin{aligned} P(f) = P(\text{tree } \tau \text{ computes } f) &\leq P((1 + \varepsilon)^{\|\tau\|} \geq (1 + \varepsilon)^{L(f)}) \\ &\leq \frac{E[(1 + \varepsilon)^{\|\tau\|}]}{(1 + \varepsilon)^{L(f)}} . \end{aligned}$$

This inequality holds as soon as $1 + \varepsilon$ is less than the radius of convergence of the generating function of the size of a tree. That is why, in the following, we are looking for a good evaluation of this radius. We get it by successive approximations, both a truncation at height d , and a set of k conditions at the root.

3.2 Basic relations

For a tree τ and an integer d , let $\|\tau\|_{d,k}$ be the number of internal nodes in tree τ at height at most d , when τ is a tree obtained by pruning according to a set of k conditions at the root. The exact values of the conditions do not matter, only the cardinality of the set of conditions does. Notice that $0 \leq k \leq n$.

Define the two generating functions

$$F_{d,k}(z) := \mathbb{E}(z^{\|\tau\|_{d,k}} / \text{the root of } \tau \text{ is a } n - \text{node});$$

$$H_{d,k}(z) := \mathbb{E}(z^{\|\tau\|_{d,k}} / \text{the root of } \tau \text{ is a } c - \text{node}) .$$

We need to investigate these generating functions and their recurrence relations, especially the function H for $k = 0$ and for $d \rightarrow +\infty$ to get finally $\mathbb{E}(z^{\|\tau\|})$. We then take $z = 1 + \varepsilon$ for ε as large as possible.

The following relations are almost obvious, they are summarized in the following lemmas.

Lemma 1 *Starting from more conditions at the root gives a smaller pruned tree, so that*

$$\|\tau\|_{d,k+1} \leq \|\tau\|_{d,k}$$

and consequently, the following inequalities on the generating functions $F_{d,k}$ and $H_{d,k}$ hold for $z \geq 1$

$$0 \leq F_{d,k+1}(z) \leq F_{d,k}(z); \tag{14}$$

$$0 \leq H_{d,k+1}(z) \leq H_{d,k}(z) . \tag{15}$$

Lemma 2 *Cutting the tree at height $d + 1$ gives a larger tree than cutting at height d :*

$$\|\tau\|_{d,k} \leq \|\tau\|_{d+1,k}$$

so obviously for $z \geq 1$

$$F_{d,k}(z) \leq F_{d+1,k}(z);$$

$$H_{d,k}(z) \leq H_{d+1,k}(z) .$$

The following recursions come from the structure of pruning the infinite tree. They are given by Lefmann and Savicky ([8]) in their lemmas 3.3 and 3.4.

Lemma 3

$$F_{0,k}(z) = H_{0,k}(z) = z;$$

$$F_{d+1,k}(z) = \frac{z}{4} \left(F_{d,k}^2(z) + \frac{k}{n} F_{d,k}(z) + 2 \left(1 - \frac{k}{n} \right) F_{d,k+1}(z) + \frac{k}{n} + 1 \right);$$

$$H_{d+1,k}(z) = \frac{z}{4} \left(2 F_{d,k}(z) H_{d,k}(z) + \frac{k}{n} H_{d,k}(z) + 2 \left(1 - \frac{k}{n} \right) H_{d,k+1}(z) + \frac{k}{n} \right) .$$

3.3 Comparing with a branching Galton-Watson process

Let us gather up the threads of the proof of theorem 1: recall that the aim is to get a fine evaluation of the radius of convergence of the generating function $\mathbb{E}(z^{\|\tau\|})$ of the size of a tree τ . For a positive real argument, this generating function is the increasing limit when $d \rightarrow +\infty$ of the generating functions $H_{d,k}$ for trees truncated at height d with $k = 0$ conditions. We are going to control the radius of convergence of $H_{d,k}$ with a uniform (in d) upper bound of the functions $H_{d,k}$ (lemma 5). Because of the intricate relations between $F_{d,k}$ and $H_{d,k}$ (which appear in lemma 3), we need the same kind of uniform (in d) upper bound of the functions $F_{d,k}$ (lemma 4). Nicely, it comes from comparing the generating function $F_{d,k}$ to the generating function of a Galton-Watson process (depending on k and not on d). Finally the study of the special case $k = 0$ is achieved in lemma 6.

We begin with lemma 3 together with inequality (14) which give for every fixed k :

$$F_{d+1,k}(z) \leq \frac{z}{4} \left(F_{d,k}^2(z) + \left(2 - \frac{k}{n} \right) F_{d,k}(z) + \frac{k}{n} + 1 \right)$$

so that defining

$$\tilde{\varphi}_k(u) := \frac{1}{4} u^2 + \left(\frac{1}{2} - \frac{k}{4n} \right) u + \frac{1}{4} + \frac{k}{4n}$$

it reads

$$F_{d+1,k}(z) \leq z\tilde{\varphi}_k(F_{d,k}(z)) . \quad (16)$$

This is the key equation for comparing our generating function to the generating function of a Galton-Watson one: indeed, $\tilde{\varphi}_k$ is the reproduction function associated to a Galton-Watson process, with generating function G_k solution of

$$G_k(z) = z\tilde{\varphi}_k(G_k(z)) . \quad (17)$$

In this branching process,

$$\begin{aligned} p_2 &:= \mathbb{P}(\text{ 2 children }) = \frac{1}{4}, \\ p_1 &:= \mathbb{P}(\text{ 1 child }) = \frac{1}{2} - \frac{k}{4n}, \\ p_0 &:= \mathbb{P}(\text{ 0 child }) = \frac{1}{4} + \frac{k}{4n}. \end{aligned}$$

Let $\alpha := k/n$ and notice that $\alpha \in [0, 1]$. Equation (17) is quadratic and can be explicitly solved, giving

$$G_k(z) = \frac{1}{2z} [4 - z(2 - \alpha) - \sqrt{16 - 8z(2 - \alpha) - (8 - \alpha)\alpha z^2}] . \quad (18)$$

This expression allows to compute explicitly the radius of convergence $\rho(G_k)$ of the generating function G_k (directly or by a derivation standard in branching processes):

$$\rho(G_k) = \frac{-2 + \alpha + 2\sqrt{1 + \alpha}}{2\alpha(1 - \alpha/8)}$$

and its expansion around $\alpha = 0$ is

$$\rho(G_k) = 1 + \frac{\alpha^2}{16} + O(\alpha^3) = 1 + \frac{k^2}{16n^2} + O\left(\frac{k^3}{n^3}\right)$$

so that for any constant $C < 1/16$, for $k > 0$ and n large enough,

$$\rho(G_k) \geq 1 + C \frac{k^2}{n^2} .$$

Now, the key equation (16) provides a uniform (in d) upper bound for the generating functions $F_{d,k}$, and the following lemma is straightforward.

Lemma 4 *For any fixed $k > 0$ and for n large enough, the sequence $(F_{d,k})_d$ increases to a limit F_k when $d \rightarrow +\infty$. Moreover, at the limit:*

$$\forall z, 1 \leq z \leq \rho(G_k), \quad F_k(z) \leq G_k(z) .$$

This implies that the radius of convergence $\rho(F_k)$ of the generating function F_k , is greater than $\rho(G_k)$, so that for any constant $C < 1/16$,

$$\rho(F_k) \geq 1 + C \frac{k^2}{n^2} .$$

Now we use the same kind of comparison for the generating functions $H_{d,k}$. For the same reasons as $F_{d,k}$, lemma 3 with inequality (15) gives for every fixed k

$$H_{d+1,k} \leq \frac{z}{4} \left[2F_{d,k}H_{d,k} + \left(2 - \frac{k}{n}\right)H_{d,k} + \frac{k}{n} \right]$$

and by the previous study of the $F_{d,k}$,

$$H_{d+1,k} \leq \frac{z}{4} \left[\left(2G_k + 2 - \frac{k}{n}\right)H_{d,k} + \frac{k}{n} \right] .$$

Define a candidate upper bound function H_k by the equation

$$H_k = \frac{z}{4} \left[\left(2G_k + 2 - \frac{k}{n}\right)H_k + \frac{k}{n} \right] ;$$

the explicit expression (18) of $G_k(z)$ gives the following form for H_k :

$$H_k(z) = \frac{\alpha z}{\sqrt{16 - 8(2 - \alpha)z - \alpha(8 - \alpha)z^2}} \quad (19)$$

and it appears that H_k has the same radius of convergence as G_k :

$$\rho(H_k) = \rho(G_k) \geq 1 + C \frac{k^2}{n^2} ,$$

for any constant $C < 1/16$. An obvious recurrence shows that for every d and k , the functions $H_{d,k}$ are upperbounded by H_k . By dominated convergence (recall that the sequence $(H_{d,k})$ is increasing in d), the sequence $(H_{d,k})$ converges to a function h_k , when d goes to infinity and $\rho(h_k) \geq \rho(H_k)$. We have proved:

Lemma 5 *For any fixed $k > 0$, and for n large enough, the sequence $(H_{d,k})_d$ increases to a limit h_k when $d \rightarrow +\infty$. There exists a function H_k , given by (19), which dominates h_k . This implies that the radius of convergence $\rho(h_k)$ of h_k , is greater than $\rho(H_k)$, so that for any constant $C < \frac{1}{16}$, and n large enough,*

$$\rho(h_k) \geq 1 + C \frac{k^2}{n^2} .$$

Study for $k = 0$:

The final step now comes from the *direct study* of $F_{d,0}$ and $H_{d,0}$ whose evolution is given by lemma 3:

$$F_{d+1,0}(z) = \frac{z}{4} [F_{d,0}(z)^2 + 2F_{d,1}(z) + 1];$$

$$H_{d+1,0}(z) = \frac{z}{4} [2H_{d,0}(z)F_{d,0}(z) + 2H_{d,1}(z)] .$$

As before, we begin with the study of $F_{d,0}$ to deduce the radius of $H_{d,0}$. By lemma 4,

$$F_{d+1,0} \leq \frac{z}{4} [F_{d,0}(z)^2 + 2G_1(z) + 1]$$

and we define a fixed point F by the equation

$$F = \frac{z}{4} [F^2 + 2G_1 + 1] .$$

This gives

$$F(z) = \frac{2 - \sqrt{4 - z^2(2G_1(z) + 1)}}{z} .$$

In a small area around 1 (recall that we are interested in z smaller than $1 + C/n^2$), $\sqrt{4 - z^2(2G_1(z) + 1)} = 1 + 1/n + O(1 - z)$, so that the radical does not introduce a singularity and $\rho(F) \geq \rho(G_1) \geq 1 + C/n^2$. By the usual recursion, it is clear that for every $d \geq 0$, $F_{d,0} \leq F$. Then

$$\rho(F_{d,0}) \geq \rho(F) \geq 1 + \frac{C}{n^2} ,$$

for any constant $C < 1/16$.

From the previous uniform upper bound on $F_{d,0}$ and from lemma 5, the functions $H_{d,0}$ satisfy

$$H_{d+1,0} \leq \frac{z}{4} [2H_{d,0}F + 2H_1] = \frac{z}{2} [H_{d,0}F + H_1]$$

and H_1 is a solution of

$$H_1 = \frac{z}{4} \left[(2G_1 + 2 - \frac{1}{n})H_1 + \frac{1}{n} \right] .$$

Define a fixed point H by

$$H = \frac{z}{2} [HF + H_1] ,$$

so that $\rho(H) \geq \rho(H_1) \geq 1 + C/n^2$ and by recursion $H_{d,0} \leq H$, for every $d \geq 0$. We have proved

Lemma 6 *For every $d \geq 0$,*

$$\rho(H_{d,0}) \geq \rho(H) \geq 1 + C \frac{1}{n^2} .$$

When d goes to infinity, by lemma 5, we finally obtain that $H_{d,0}$ increases to the generating function of the size of a pruned, simplified and/or tree, and the radius of convergence is at least $1 + C/n^2$, thus giving the theorem.

4 Concluding remarks

We begin by commenting on the quality of the bounds in Theorem 1: The lower bound is tight, but the upper bound is not, and can possibly be improved. A natural approach would be to get a better lower bound of the radius of convergence, possibly using results by Nguyễn Thế [10]. However, it is possible that Markov's inequality is not strong enough to give a really tight upper bound, and that a different approach may have to be sought.

Numerical computations suggest a point worth mentioning (accordingly for a small number n of variables, but the situation is probably even more marked for larger n). Assuming all trees equally likely gives a very high probability $P(f)$ to functions of complexity 0 or 1 (mainly the constants), and functions of higher complexity quickly become negligible. The alternative distribution $\pi(f)$ behaves in a similar way, and is even more biased towards literals.

We should investigate further the relationship between the probability distributions $P(\cdot)$ (defined by equiprobable and/or trees) and $\pi(f)$ (defined by labelling critical branching processes). We conjecture that, except for literals, we always have that $\pi(f) < P(f)$ for binary planar trees.

A desirable extension of the model for boolean formulae takes into account the commutativity or associativity of the boolean operators when representing a function by a tree (i.e. by a boolean formula). Preliminary investigations show that the

natural model becomes that of non planar (for commutativity) general (for associativity) trees, for which we can write generating functions using Polya's theory of tree enumeration [12]. Such tree models are also related to those of Woods [17], where he proved a general theorem on the existence of a limiting distribution, although no explicit computations were given. This should allow us to prove, in a manner similar to that of Section 2, the existence of a limiting probability distribution, and to compute numerical distributions for small values of n .

It is then natural to try and compare these different distributions, which are all defined on the same set of boolean functions but stem from different underlying tree models. We might extend the distribution $\pi(f)$ to non planar general trees, and examine the conjecture $\pi(f) < P(f)$ in this new context. The relationship between complexity and probability also deserves further investigations: Does a modification of Theorem 1 still hold if we substitute different probability distributions for P ? What about other complexity measures? We hope to study all these points in a forthcoming paper.

References

- [1] P. Billingsley. *Probability and Measure, Third Edition*. Wiley, 1995.
- [2] H. Buhrman and R. de Wolf. Complexity measures and decision tree complexity : a survey. *Theoretical Computer Science.*, 288:21–43, 2002.
- [3] M. Drmota. Systems of functional equations. *Random Structures and Algorithms*, 10:103–124, 1997.
- [4] P. Flajolet and A. M. Odlyzko. Singularity analysis of generating functions. *SIAM J. on Discrete Math.*, 3(2):216–240, 1990.
- [5] P. Flajolet and R. Sedgewick. Analytic combinatorics: Functional equations, rational and algebraic functions. Technical Report 4103, INRIA, January 2001.
- [6] J. Friedman. Probabilistic spaces of boolean functions of a given complexity: generalities and random k -sat coefficients. Technical Report CS-TR-387-92, Princeton University, Princeton, NJ, 1992.
- [7] S.P. Lalley. Finite range random walk on free groups and homogeneous trees. *Ann. Probab.*, 21(4):2087–2130, 1993.

- [8] H. Lefmann and P. Savický. Some typical properties of large and/or boolean formulas. *Random Structures and Algorithms*, 10:337–351, 1997.
- [9] R. Lyons, R. Pemantle, and Y. Peres. Conceptual proofs of $l \log l$ criteria for mean behavior of branching processes. *Ann. Probab.*, 23:1125–1138, 1995.
- [10] M. Nguyễn Thế. Distribution of the size of simplified or reduced trees. In *Mathematics and Computer Science II*. Birkhauser, 2002.
- [11] J. B. Paris, A. Vencovská, and G. M. Wilmers. A natural prior probability distribution derived from the propositional calculus. *Annals of Pure and Applied Logic*, 70:243–285, 1994.
- [12] G. Pólya and R.C. Read. *Combinatorial enumeration of Groups, Graphs and Chemical Compounds*. Springer Verlag, New York, 1987.
- [13] P. Savický. Random boolean formulas representing any boolean function with asymptotically equal probability. *Discrete Mathematics*, 83:95–103, 1990.
- [14] P. Savický. Bent functions and random boolean formulas. *Discrete Mathematics*, 147:211–234, 1995.
- [15] P. Savický. Complexity and probability of some boolean formulas. *Combinatorics, Probability and Computing*, 7:451–463, 1998.
- [16] P. Savický and A. Woods. The number of boolean functions computed by formulas of a given size. *Random Structures and Algorithms*, 13:349–382, 1998.
- [17] A. Woods. Coloring rules for finite trees, and probabilities of monadic second order sentences. *Random Structures and Algorithms*, 10:453–485, 1997.