The Average Height of Binary Trees and Other Simple Trees

PHILIPPE FLAJOLET

INRIA, 78150 Rocquencourt, France

AND

ANDREW ODLYZKO

Bell Laboratories, Murray Hill, New Jersey 07974 Received January 5, 1981; revised April 14, 1982

The average height of a binary tree with *n* internal nodes is shown to be asymptotic to $2\sqrt{\pi n}$. This represents the average stack height of the simplest recursive tree traversal algorithm. The method used in this estimation is also applicable to the analysis of traversal algorithms of unary-binary trees, unbalanced 2-3 trees, t-ary trees for any *t*, and other families of trees. It yields the two previously known estimates about average heights of trees, namely for labeled nonplanar trees (a result due to Renyi and Szekeres) and for planar trees (a result of De Bruijn, Knuth, and Rice). The method developed here, which relies on a singularity analysis of generating functions, is new and widely applicable.

0. INTRODUCTION

We consider the problem of the relation between *height* and *size* in *trees*, for various types of trees. Given a family F of trees with F_n the subset of those trees formed with n nodes, the problem is to determine the average height defined by

$$\overline{H}_n(F) = \frac{1}{\operatorname{card} F_n} \sum_{t \in F_n} \operatorname{height}(t).$$

In this paper we solve this problem for the family B of binary trees.

THEOREM B. The average height of binary trees with n internal nodes satisfies

$$\overline{H}_n(B) \sim 2 \sqrt{\pi n} \qquad as \quad n \to \infty.$$

So far the only result available about average heights of planar trees dealt with the family G of general trees, i.e., planar trees with unrestricted node degrees [3].

THEOREM G (De Bruijn et al.). The average height of general planar trees (of arbitrary node specification) with n nodes satisfies

$$H_n(G) \sim \sqrt{\pi n}$$
 as $n \to \infty$.

The similarities in the forms of Theorems G and B might induce the reader to believe that Theorem B is only a simple modification of Theorem G. The methods differ, however, in an essential way.

Theorem G is proved by first giving exact enumerations for the number of trees of fixed height and fixed size; these are expressed as certain sums of binomial coefficients. The asymptotics are then performed by appealing to properties of the *Mellin integral transform*. This method is an important starting point of a number of analyses [12] amongst which we mention those of radix exchange sort, digital search, Patricia trees, sorting networks, and register allocation. Many other enumeration results, such as those in [1], also are obtained by starting with explicit formulae for generating functions.

The problem we encounter with binary trees is that exact enumeration formulae are no longer available for the number of trees of fixed size and height and we only have recursive formulae. The path we follow relies on the principle that the coefficients of a generating function are largely determined by the location and nature of its *singularities*. It is also the only recourse we know of when one has at one's disposal nothing but functional equations over generating functions.

The power of the method is due to the fact that many enumeration problems have generating functions satisfying functional equations of some sort. Singularities are located by applying approximations and obtaining asymptotic expansions in the complex plane. Coefficients of generating functions are then estimated using contour integration.

Despite its power this method has only rarely been used in algorithmic analyses. The work closest to ours is the determination by Odlyzko of the number of balanced 2–3 trees [15]. We demonstrate the generality of our approach by showing

THEOREM S. For each simple family of trees S there exists an effectively computable constant c(S) such that the average height of a tree in S with n nodes is

$$\overline{H}_n(S) \sim C(S) \sqrt{\pi n} \quad as \quad n \to \infty.$$

A family of trees is said to be *simple* if, essentially, for each r there is a finite set of allowable labels for nodes of degree r. Theorem S contains as subcases the result by De Bruijn *et al.* on the average height of planar trees, and (though it does not immediately fit into our framework) a result by Renyi and Szekeres about nonplanar labeled trees.

Since the height of a tree represents the stack size needed in recursively traversing the tree, Theorem S also yields the analysis of the simplest recursive tree traversal algorithm in a diversity of contexts. The reader should, however, be warned that statistics on binary search trees represent a different problem to be briefly discussed later.

To conclude this introduction, we should like to emphasize that the interest of this paper is largely methodological. Almost all *classical* analyses of algorithms follow a chain starting with exact enumeration formulae derived by *direct* counting arguments continued by real approximations (usually approximating discrete sums by integrals). There is a very clear stage at which this approach fails to apply: either the nature of the problem leads to a combinatorial expression whose estimation proves intractable, or even more plainly—as in the case here—no combinatorial expression is available at all. In both cases, studying the analytical properties of the corresponding generating functions—especially their singularities—leads to solution of problems not tractable by more elementary methods.

The plan of the paper is as follows: In the binary case, a certain generating function of the \overline{H}_{n} , H(z), is shown to be the sum of quantities defined by a quadratic recurrence (Section 2). Recovering the \overline{H}_{n} from H(z) requires a detailed analytical investigation of the behavior of H(z). A detailed outline of the method is given at the beginning of Section 3. This method is then developed fully in Sections 3-5.

We shall indicate how to extend the method to any simple family of trees (Section 6). This includes all previously known results about the heights of trees and provides the very general result stated in Thereom S. Last (Section 7), we shall discuss the limits of the present approach and some of its extensions to estimates of higher moments and limit distributions.

A priliminary version of this paper [5] was presented at the 21st Symposium on Foundation of Computer Science, Syracuse, New York, October 13–15, 1980. Similar results have been obtained by a somewhat different analytic method by G. B. Brown and B. O. Shubert, "On Random Binary Trees" (preprint).

1. TREE TRAVERSAL

We shall limit ourselves here to a short algorithmic discussion of tree traversal, referring the reader to [11] for more details.

Perhaps one of the simplest recursive algorithms is the algorithm for *visiting*—one also says *traversing* or *exploring*—nodes of a planar *tree*. The algorithm occurs in a number of contexts in compiling, program transformation, term rewriting systems, optimization, and related areas. Loosely described, this simple algorithm looks like

> procedure VISIT(T: tree) do-something-with(root(T)); for U subtree-of-root-of T do VISIT(U) rof erudecorp.

In specific applications, the trees input to the algorithm usually obey some particular format. For instance, one may encounter: expression trees involving nullary symbols (variables), unary symbols (log, sin, $\sqrt{}$) and binary symbols (+, -, \times , \div); syntax trees of various types with nodes of possibly unbounded degrees (as in list-of-instruction nodes); trees to represent terms in formal manipulation systems; and others.

We are interested here in the behavior of the tree exploration procedures in such contexts. The running time analysis of the VISIT procedure is not difficult since the complexity is clearly linear in the size of the input tree. The main problem is to evaluate *storage utilization*, i.e., to determine the average *stack size* (equivalently recursion depth) required for exploring a tree, as a function of the size of the tree. For a given tree, the stack size required by the visit is equal to the height of the tree. Average case analysis of the algorithm applied to a family F of input trees thus reduces to determining *average heights* of trees in F.

The results of this paper completely solve the average cases analysis of tree traversal applied to any simple family of inputs. In particular, Theorem B can be rephrased as

THEOREM B. The recursive traversal procedure applied to binary trees of size n has average storage complexity

$$\overline{H}_n(B) \sim 2\sqrt{\pi n} \qquad as \quad n \to \infty.$$

It should be mentioned here that the result by De Bruijn et al. relative to the family S of general planar trees, namely, that

$$\overline{H}_n(G) = \sqrt{\pi n} - \frac{1}{2} + o(1) \quad \text{as} \quad n \to \infty,$$

gives some information on the height of binary trees, as well as on binary tree traversal. Indeed the rotation correspondence ([11], Sect. 2.3.2) transforms a general tree with n nodes into a binary tree containing (n-1) internal (binary) nodes, hence n external (nullary) nodes. Let ρ be this correspondence exemplified by Fig. 1. The reader can convince himself easily that

$$height(t) = height^*(\rho(t)) + 1,$$

where height* denotes the one-sided height of binary trees, defined as the maximum number of (internal) left branching nodes on any branch of the tree. Since for any binary tree

$$height(u) \ge height^*(u) + 1$$
,

it follows for the family B of binary trees that

$$\overline{H}_{n-1}(B) \geqslant \overline{H}_n(G).$$



FIG. 1. The Rotation Correspondence transforms a general tree into a binary tree: the leftmost-son relation becomes the left-son relation and the right-brother relation becomes the right-son relation; the root of the general tree is dropped. External nodes of the binary tree are not represented.

Thus the estimation of the average height of general planar trees shows $\overline{H}_n(B)$ to be at least of order $\sqrt{\pi n}$. Theorem B shows that $\overline{H}_n(B)$ is essentially twice as large; i.e., we obtain the surprising result that the average height of binary trees is practically the sum of the average right and left heights.

The result about heights of general trees is also of interest in another context. It is possible [11, 12] to optimize the recursive visit procedure in the case of binary trees by eliminating endrecursion. The resulting iterative algorithm keeps at each stage a list of right subtrees that still remain to be explored; the storage complexity of this optimized iterative algorithm is easily seen to correspond exactly to one-sided height. Hence, Theorem G can be expressed as

THEOREM G'. The iterative traversal procedure for binary trees of size n has average storage complexity

$$\overline{H}_{n+1}(G) \sim \sqrt{\pi n} \quad as \quad n \to \infty.$$

Thus the expected memory complexity of the optimized iterative exploration algorithm is asymptotically (for large sizes of trees) half the expected complexity of recursive exploration.

To conclude this brief algorithmic discussion, let us mention that if the left-ro-right order in the exploration need not be kept, then exploration can be reduced to a pebbling game on trees which is equivalent to register allocation. The analysis of optimal register allocation applies there, and rephrasing results of [6, 8, 14] one gets the following result:

THEOREM O (Optimal exploration of binary trees). The minimal stack size for exploring binary trees with n internal nodes when the left-to-right order is irrelevant has average value

$$\bar{O}_n = \log_4 n + P(\log_4 n) + o(1),$$

where P is a continuous function with period 1.

FLAJOLET AND ODLYZKO

This estimation applies, e.g., in the context of preprocessing (allowing one bit per node).

Some comments are now in order about the relevance of our statistics: we perform analyses of tree traversal by averaging over all possible trees. The results are thus significant only when inputs do not satisfy any further conditions. Basically our analyses apply to input trees with an independent labelling of nodes; such is the case at least for expression trees in compiling, or term trees in formal manipulation systems.

As a first approximation, our treatment can also be applied to term trees in heterogeneous algebra. In this context several types of objects are present and operators have type restrictions. This involves syntax trees of various sorts. Counting of such trees then leads to similar statistics with generating functions that are still algebraic, and an exact treatment along our lines should be feasible (for the particular case of syntax trees of linear grammars, see [9]).

An analysis of our type does not apply when trees occur as components of more complex structures, as appears in binary search trees or tournament trees. For instance, binary search trees have monotonic labellings, and the probability distribution induced on shapes of trees by random insertion is known [12] and far from uniform. Indeed for binary search trees, the average height for size n is $O(\log n)$ corresponding to a logarithmic search, and Robson [19] has obtained the following bounds:

THEOREM BST. Let \overline{K}_n be the average height of binary search trees generated by n independent random insertions. Then

$$c_1 \log n + o(\log n) \leqslant \overline{K}_n \leqslant c_2 \log n + o(\log n),$$

with $c_1 > 3.6$ and $c_2 = 4.31170...$

The precise asymptotic behavior of $\overline{K}_n/\log n$ is not yet known, although it is known to tend to a limit [20].

To conclude this presentation of alternative statistics, let us mention the result of Flajolet [4] relative to the height of index trees in dynamic hashing, which also applies to digital search trees (tries):

THEOREM D. Let \overline{L}_n be the average height of a digital search tree constructed over n keys uniformly drawn on [0, 1]. Then

$$\overline{L}_n \sim 2 \log_2(n)$$
 as $n \to \infty$.

Some considerations about heights in combinatorial structures are developed in our final section. We have not addressed in this paper the somewhat different problem of path lengths in trees, (see [11, 12]) and the related question of levels of nodes in trees (which can be used to derive upper bounds on heights). For this last problem the reader is referred to the excellent paper of Meir and Moon [13].

2. THE HEIGHT OF BINARY TREES: BASIC RECURRENCES

We consider the set *B* of *binary trees* in the sense of Knuth [11]: every node has either 0 or 2 successors and left and right successors are distinguished. The *size* of a tree in *B* is the number of its internal binary nodes, i.e., the number of nodes with two successors. We let |t| denote the size of *t*. We also define

$$B_n = \operatorname{card} \{t \in B \colon |t| = n\}.$$

The *height* of a binary tree is the number of nodes along the longest branch from the root and is given inductively by

$$\begin{aligned} \text{height}(\Box) &= 1\\ \text{height}(t) &= 1 + \max\{\text{height}(t_2), \text{height}(t_2)\},\\ \text{where} \quad t_1 &= \text{left}(t) \quad \text{and} \quad t_2 &= \text{right}(t). \end{aligned}$$

Figure 2 shows the distribution of height on trees of size 4. We introduce the quantities

$$B_n^{[h]} = \operatorname{card} \{t \in B : |t| = n \text{ and } \operatorname{height}(t) \leq h\},\$$

and \overline{H}_n , the average height of all trees of size n, is

$$\overline{H}_n = \frac{H_n}{B_n}$$
 with $H_n = \sum_{h \ge 1} h(B_n^{[h]} - B_n^{[h-1]}).$ (1a)

From the definition, we clearly have that $B_n^{[h]} = B_n$ if h > n. Rearranging the sum in (1a), we thus get

$$H_n = \sum_{h \ge 0} (B_n - B_n^{[h]}).$$
(1b)

The first values of these quantities are displayed in Table I.



FIG. 2. Amongst the 14 trees of size 4, there are 8 trees of height 5(a), and 6 trees of height 4(b). Here \bigcirc denotes internal nodes.

FLAJOLET AND ODLYZKO

ΤA	BL	LΕ	
----	----	----	--

The Distribution of Height in Trees of Size ≤ 7 with $A_{n,h} = B_n^{[h]} - B_n^{[h-1]}$

n	B _n	A _{n.2}	A _{n,3}	A _{n.4}	$A_{n,5}$	A ,, 6	A _{n,7}	A _{n,8}	\overline{H}_n
1	1	1							2.0
2	2	0	2						3.0
3	5	0	1	4					3.8
4	14	0	0	6	8				4.57
5	42	0	0	6	20	16			5.24
6	132	0	0	4	40	56	32		5.88
7	429	0	0	1	68	152	144	64	6.47

We now introduce the generating functions relative to the B_n , $B_n^{[h]}$, and H_n :

$$B(z) = \sum_{n \ge 0} B_n z^n,$$
$$B^{[h]}(z) = \sum_{n \ge 0} B^{[h]}_n z^n,$$
$$H(z) = \sum_{n \ge 0} H_n z^n.$$

The inductive definition of binary trees shows that the B_n satisfy the recurrence

$$B_n = \sum_{n_1 + n_2 + 1 = n} B_{n_1} B_{n_2},$$

•

whence

$$B(z) = 1 + z(B(z))^2$$
 (2a)

and

$$B(z) = (1 - \sqrt{1 - 4z})/2z;$$
 $B_n = (n+1)^{-1} {2n \choose n}.$ (2b)

The B_n 's are the Catalan numbers. The Stirling formula implies the classical approximation

$$B_n = (4^n / \sqrt{\pi n^3})(1 + O(1/n)).$$
 (2c)

The same decomposition principle that gives the equation for B applies to the $B^{[h]}$ yielding the recurrence

$$B^{[h+1]}(z) = 1 + z(B^{[h]}(z))^2; \qquad B^{[0]}(z) = 0.$$
(3)

No simple expression is available for the $B_n^{[h]}$ coefficients. The first values of the $B_n^{[h]}(z)$ are

$$B^{[0]}(z) = 0; \quad B^{[1]}(z) = 1; \quad B^{[2]}(z) = 1 + z; \quad B^{[3]}(z) = 1 + z + 2z^2 + z^3;$$
$$B^{[4]}(z) = 1 + z + 2z^2 + 5z^3 + 6z^4 + 6z^5 + 4z^6 + z^7.$$

Obviously, degree $(B^{[h]}(z)) = 2^{h-1} - 1$, and $B_n^{[h]} = B_n$ for n < h. Summarizing the recurrences, we can state

PROPOSITION 1. In the ring of formal power seires,

$$H(z) = \sum_{h \ge 0} (B(z) - B^{[h]}(z)),$$

where B and the $B^{[h]}$ satisfy

$$B(z) = 1 + z(B(z))^2;$$
 $B^{[h+1]}(z) = 1 + z(B^{[h]}(z))^2$ with $B^{[0]}(z) = 0.$

3. Outline of the Method and the First Analytical Continuation of H(z)

Our task is to estimate the coefficients H_n of H(z). The difficulty we face is that we possess neither a closed form expression for H(z) nor even a functional equation satisfied by H(z). This difficulty is due to the nonlinear nature of recurrence (3).

To estimate H_n , we will use Cauchy's theorem which states that

$$H_n = \frac{1}{2i\pi} \int_{\Gamma} H(z) \frac{dz}{z^{n+1}},$$

where Γ is any simple closed curve in the region of analyticity of H(z) that encircles the origin. We shall adopt here for Γ a contour far away from the origin; this has the advantage that even partial information on the growth of H(z) can be used to estimate the Cauchy integral giving H_n .

In the present case, it is easy to show (Proposition 2) that H(z) is analytic in the disk $|z| < \frac{1}{4}$ but in no larger disk. Since H(z) has positive coefficients, this implies that H(z) has a singularity at $\frac{1}{4}$. This singularity, however, turns out to be the only one on the circle $|z| = \frac{1}{4}$. We show in effect that H(z) is analytic in a region of the form

$$D = \{z \colon |z| < \lambda, \operatorname{Arg}(z - \frac{1}{4}) > \omega\}$$

for some constants $\lambda > \frac{1}{4}$ and $\omega \in (0, \pi/2)$. The proof uses both a continuity argument (Proposition 3) and a local study of the recurrence around $\frac{1}{4}$ (Proposition 4).

FLAJOLET AND ODLYZKO

The expansion of H(z) which leads to our estimates of \overline{H}_n is obtained in Section 4. It is shown that in a neighborhood of $z = \frac{1}{4}$ in D, H(z) is the sum of a logarithmic term and a remainder term of smaller order. Most of the complexity of our solution lies in this derivation. This expansion of H(z) is obtained by an extensive analysis of the recurrence of Proposition 1.

The estimates of the coefficients H_n are obtained from the expansion of H(z) in Section 5 with the help of an appropriate contour of integration. This contour, which follows the boundary of a region similar to D (see Fig. 4) has the property that the integral depends almost exclusively on the behavior of H(z) near $z = \frac{1}{4}$. A crucial role is played here by the fact that the contour can essentially include line segements of the form

$$\{re^{\pm i\phi}: 0 < r < \varepsilon\}$$

for some $\varepsilon > 0$ and some fixed $\phi \in (0, \pi/2)$. (If it were not for this fact, we would need a better expansion of H(z).) Proposition 6 gives a general result that applies in many similar situations, and which concludes our proof of Theorem B.

We shall now proceed by proving that the expression for H(z) derived in Section 2 (Proposition 1) is also valid analytically in some domain and is a way of continuing H(z) analytically outside its circle of convergence.

PROPOSITION 2. H(z) has radius of convergence $\frac{1}{4}$ and the equality

$$H(z) = \sum_{h \ge 0} (B(z) - B^{[h]}(z))$$

is valid analytically inside the domain

$$C_0 = \{z : |z| \leq \frac{1}{4}, z \neq \frac{1}{4}\},\$$

the determination of $\sqrt{1-4z}$ in B(z) being positive for real $z < \frac{1}{4}$. Moreover, the sum for H(z) converges absolutely for z in C_0 .

Proof. For each nonempty tree t, we have the obvious inequalities

$$1 \leq \operatorname{height}(t) \leq |t|,$$

which shows that

$$B_n \leqslant H_n \leqslant nB_n$$

From estimate (2c) of B_n it follows that H(z) has radius of convergence equal to $\frac{1}{4}$.

Notice first that the Taylor series of B(z) is absolutely convergent when $|z| \leq \frac{1}{4}$. Indeed, it converges as $\sum n^{-3/2}$ for all z with $|z| = \frac{1}{4}$. Let $R_m(z)$ denote $\sum_{n \geq m} B_n z^n$. Then from simple majorizations we have

$$|B(z) - B^{[h]}(z)| \leq R_h(|z|) \quad \text{when} \quad |z| \leq \frac{1}{4},$$

and $B^{[h]}(z) \to B(z)$ for any z such that $|z| \leq \frac{1}{4}$. The nature of the convergence is obtained by writing

$$B(z) - B^{[h+1]}(z) = z(B(z) - B^{[h]}(z))(B(z) + B^{[h]}(z)).$$

Dividing by 2B(z) and setting

$$e_h(z) = (B(z) - B^{[h]}(z))/(2B(z)),$$

this recurrence is transformed into

$$e_{h+1}(z) = (1 - \sqrt{1 - 4z}) e_h(z)(1 - e_h(z)).$$

We shall also set $\varepsilon = \varepsilon(z) = (1 - 4z)^{4/2}$, the determination of the square root being as above. In this notation,

$$e_{h+1}(z) = (1 - \varepsilon(z)) e_h(z)(1 - e_h(z))$$
 with $e_0(z) = \frac{1}{2}$. (4)

Assuming z to be in C_0 , we have $|1 - \varepsilon| < 1$ and the convergence of the $e_h(z)$ to 0 is geometric with

$$|e_h(z)| < c(z) |1 - \varepsilon(z)|^h$$
 for some $c(z)$;

thus $\sum_{h>0} e_h(z)$ is also convergent and the same holds true for the sum $\sum_{h>0} (B(z) - B^{[h]}(z))$.

As will appear from later considerations, $e_n(\frac{1}{4}) \sim 1/n$ and thus $e_n(\frac{1}{4}) \to 0$ as $n \to \infty$, but at the point $z = \frac{1}{4}$ the series $\sum_{n \ge 0} e_n$ diverges as the harmonic series.

In the sequel we shall mostly work with the functions $e_n(z)$. We shall thus replace Eqs. (3) and (4) by the set

$$e_0(z) = \frac{1}{2}, \qquad e_{n+1}(z) = (1 - \varepsilon(z)) e_n(z)(1 - e_n(z))$$
 (5)

and

$$H(z) = \frac{4}{1 + \varepsilon(z)} \sum_{n \ge 0} e_n(z), \qquad (6)$$

where $\varepsilon(z) = (1 - 4z)^{1/2}$.

We proceed to show that H(z), as given by the previous recurrence equations (5) and (6), is analytic in a domain larger than the circle of convergence. To that purpose, we use an argument which is essentially topological and whose principle is based on some continuity properties of a convergence criterion.

We take the complex plane cut along the ray $z > \frac{1}{4}$, $\varepsilon(z)$ being as before that branch of $(1-4z)^{1/2}$ which is positive for z real, $z < \frac{1}{4}$. For fixed z, consider the function of y

$$f(y) = (1 - \varepsilon(z)) y(1 - y),$$

in which z enters as a parameter.



FIG. 3. A diagram representing the relative positions of the boundaries of C_0 (circle a), of D_0 (curve c) and of a convergence region guaranteed by Propositions 3 and 4 (curve b).

From what we have seen $e_n(z) = f^{(n)}(\frac{1}{2})$, where $f^{(n)}$ is the *n*th iterate of *f*. We are interested in the area in which $e_n(z) \to 0$ in a nondegenerate way. This can only occur if 0 is an *attractive fixed point* of f(y), i.e., if $f'(0) = (1 - \varepsilon)$ has modulus less than 1. In this case any sequence $u_{n+1} = f(u_n)$ converges provided its initial value is close enough to the fixed point.

We thus restrict attention to values of z in the domain

$$D_0 = \{z \colon |1 - \varepsilon(z)| < 1\}.$$

Domain D_0 is the inside of a cardioid-shaped contour that properly contains C_0 (see Fig. 3). The domain of values of z for which $e_n(z) \to 0$ as $n \to \infty$ thus lies somewhere between C_0 and D_0 .

The following lemma is a useful convergence criterion for the sequence $\{e_m(z)\}_{m\geq 0}$.

LEMMA 1 [Convergence criterion for $e_n(z)$]. A necessary and sufficient condition for the sequence $\{e_n(z)\}_{n\geq 0}$ to converge to 0 for $z \in D_0$ is that for some m

$$|e_m(z)| < |1 - \varepsilon(z)|^{-1} - 1.$$

Furthermore, if this condition is satisfied, then the convergence of the $|e_n(z)|$ for $n \ge m$ is monotonic.

Proof. The condition of the lemma is trivially necessary. To obtain its sufficiency, note that applying the triangular inequality to the recurrence of the e_n leads to

$$|e_{n+1}| \leq |1-\varepsilon| |e_n| + |1-\varepsilon| |e_n|^2$$

hence

$$|e_{n+1}| - |e_n| \leq |e_n| |1 - \varepsilon| (|e_n| + 1 - |1 - \varepsilon|^{-1})$$

. '

Thus if $|e_n| < |1 - \varepsilon|^{-1} - 1$, then $|e_{n+1}| < |e_n|$ and a fortiori $|e_{n+1}| < |1 - \varepsilon|^{-1} - 1$ so that the argument can be repeated. We have thus established: if for some m

$$|e_m| < |1-\varepsilon|^{-1}-1,$$

then for all $n \ge m$:

$$|e_{n+1}| < |e_n| < |1-\varepsilon|^{-1} - 1.$$

It remains to prove that $|e_n| \rightarrow 0$ in this case. Assume a contrario

$$|e_n| \to L \neq 0$$
 as $n \to \infty$.

Then, from the basic recurrence

$$e_{n+1} = (1-\varepsilon) e_n (1-e_n),$$

it follows by continuity that

$$|1-e_n| \to 1/|1-\varepsilon|.$$

The conditions

$$|e_n| \rightarrow L < |1-\varepsilon|^{-1} - 1$$
 and $|1-e_n| \rightarrow 1/|1-\varepsilon|$

entail that the only possible accumulation points of the sequence $\{e_n\}$ are points α satisfying

 $|\alpha| = L < |1 - \varepsilon|^{-1} - 1$ and $|1 - \alpha| = 1/|1 - \varepsilon|,$

but these two conditions are clearly contradictory. We must therefore have L = 0, which completes the proof.

Using (5), the first few values of the $e_n(z)$ can be expressed in terms of $\varepsilon(z)$:

$$e_0(z) = \frac{1}{2}, \qquad e_1(z) = \frac{1}{4}(1-\varepsilon), \qquad e_2(z) = \frac{3}{16}(1+\varepsilon/3)(1-\varepsilon)^2.$$

We see, e.g., that e_0 already satisfies the convergence criterion for $z \in \left[-\frac{4}{9}, \frac{2}{9}\right]$.

LEMMA 2 (The open set property for the convergence domain of H(z)). The domain K of values of z in D_0 for which the sequence $\{e_n(z)\}_{n\geq 0}$ converges is an open set. Furthermore the series $\sum_{n\geq 0} e_n(z)$ is analytic in K.

Proof. The proof is based on the continuity of the convergence criterion of Lemma 1. If $z \in K$, then for some m,

$$\phi(z) = |1 - \varepsilon(z)|^{-1} - |e_m(z)| > 1.$$

Now, clearly, $\phi(z)$ is a continuous function of z inside D_0 ; thus there exists a positive real h, such that for all z' satisfying

$$|z'-z| < h,$$

we have $\phi(z') > 1$. Hence, $e_m(z')$ also satisfies the convergence criterion and $e_n(z') \to 0$ as $n \to \infty$.

To prove analyticity we observe that the convergence of $e_n(z)$ to 0 is geometric and uniform. Indeed, since $|1 - \varepsilon(z)| (1 + |e_m(z)|) < d < 1$ for some d, there exists a real δ such that for all z' satisfying $|z' - z| < \delta$,

$$|1 - \varepsilon(z')| (1 + |e_m(z')|) < d < 1.$$

Since for $n \ge m$ the quantities $|e_n(z')|$ decrease with n, we thus have

$$|e_n(z')| \leqslant d^{n-m} |e_m(z')|,$$

hence $|e_n(z')| \leq cd^n$ for some real *c*, uniformly in $|z'-z| < \delta$. This shows $\sum_{n \geq 0} e_n(z')$ to be uniformly convergent in $|z'-z| < \delta$, and so the sum is analytic in $|z'-z| < \delta$.

We can apply Lemma 2 to the points in the disk $|z| \leq \frac{1}{4}$ with $z \neq \frac{1}{4}$. For each such z, there exists a $\delta(z) > 0$ such that H(z) is analytic inside the domain

$$D(z) = \{ z' : |z' - z| < \delta(z) \}.$$

The domain

$$D_1 = \bigcup_{z \in C_0} D(z)$$

is open, properly contains C_0 , and H(z) is analytic inside it.

The point $z = \frac{1}{4}$ is on the boundary of D_1 , but we do not know yet the exact configuration of this boundary at $\frac{1}{4}$. From simple topological considerations (essentially the Borel-Lebesgue lemma), however we have

PROPOSITION 3. For each η , there exists a $\lambda > \frac{1}{4}$ such that H(z) is analytic in the indented crown

$$|\operatorname{Arg}(z)| > \eta$$
 and $|z| < \lambda$.

4. Continuation of H(z) Around the Singularity

We now study the behavior of the sequence $\{e_n(z)\}$ when z lies in a sector around $\frac{1}{4}$ situated inside D_0 . We first show that, in part of the domain, the initial values of

 $e_n(z)$ decrease steadily; we then prove that, at some stage, they satisfy the conditions of the convergence criterion (Lemma 1).

We start with the following lemma:

LEMMA 3. Let g(z) = y(1 - y). If y satisfies $|y| \leq \frac{1}{2}$ and $0 \leq \operatorname{Arg}(y) \leq \operatorname{Arccos} \frac{1}{4}$,

then $|g(y)| \leq |y|$ and $0 \leq \operatorname{Arg} g(y) \leq \operatorname{Arg} (y)$.

Proof. Let $y = re^{it}$. Then

$$g(y) = r(1+r^2-2r\cos t)^{1/2} \exp\left(i\left(t-\operatorname{Arctan}\frac{r\sin t}{1-r\cos t}\right)\right).$$

The hypothesis implies that $2r \cos t \ge r^2$, whence the bound for |g(y)|. On the other hand, as is easy to see,

$$0 \leq \arctan \frac{r \sin t}{1 - r \cos t} \leq \arctan \sin t \leq t,$$

whence the bound for $\operatorname{Arg} g(y)$.

LEMMA 4 (Initial decrease of $|e_m(z)|$). Suppose that $z \in D_0$, Im $z \ge 0$, and let $N(z) = 1 + \lfloor \operatorname{Arccos} \frac{1}{4} / \operatorname{Arg}(1 - \varepsilon(z)) \rfloor$. Then for all n < N(z),

$$|e_{n+1}(z)| \leq |e_n(z)| \leq \frac{1}{4}$$

and $0 \leq \operatorname{Arg}(e_{n+1}) \leq (n+1) \operatorname{Arg}(1-\varepsilon(z))$.

Proof. The proof follows immediately by iterative use of Lemma 4.

The restriction that Im $z \ge 0$ in Lemma 4 and in the sequel is made for notational convenience since

$$e_n(\bar{z}) = \overline{e_n(z)}, \qquad H(\bar{z}) = \overline{H(z)}, \dots.$$

We are now left with the task of proving that for z in a certain sector around $\frac{1}{4}$, $e_N(z)$ satisfies the conditions of Lemma 1.

Our treatment heavily relies on a trick used by De Bruijn [2, p. 157] in the context of nonlinear recurrences of a similar type. We shall express it as follows:

LEMMA 5 (Alternative recurrence for the $e_n(z)$). If all the $e_j(z)$ for j = 0, 1, ..., n-1 are different from 1, then the following relation holds:

$$\frac{(1-\varepsilon)^n}{e_n} = \frac{1-(1-\varepsilon)^n}{\varepsilon} + 2 + \sum_{j < n} \frac{e_j}{(1-e_j)} (1-\varepsilon)^j.$$
(7)

Proof. We start again from the recurrence

$$e_{j+1} = (1-\varepsilon) e_j (1-e_j),$$

and we take out the $(1-\varepsilon)^j$ factor present in e_j ,

$$e_{j+1}/(1-\varepsilon)^{j+1} = (e_j/(1-\varepsilon)^j)(1-e_j).$$

The essential trick now is to take reciprocals,

$$(1-\varepsilon)^{j+1}/e_{j+1} = ((1-\varepsilon)^j/e_j)(1-e_j)^{-1}$$

and use the expansion

$$(1-u)^{-1} = 1 + u + u^2/(1-u),$$

valid provided $u \neq 1$. Here we get

$$(1-\varepsilon)^{j+1}/e_{j+1} = ((1-\varepsilon)^j/e_j)(1+e_j+e_j^2/(1-e_j)),$$

$$(1-\varepsilon)^{j+1}/e_{j+1} = (1-\varepsilon)^j/e_j + (1-\varepsilon)^j + (e_j/(1-e_j))(1-\varepsilon)^j.$$

When we sum these identities for j = 0, ..., n - 1, terms like $(1 - \varepsilon)^j / e_j$ cancel out and using the initial value $1/e_0 = 2$, we get

$$(1-\varepsilon)^n/e_n = \sum_{j < n} (1-\varepsilon)^j + 2 + \sum_{j < n} \frac{e_j}{1-e_j} (1-\varepsilon)^j,$$

from which the lemma follows.

The relation of Lemma 5 suggests $\varepsilon(1-\varepsilon)^n/(1-(1-\varepsilon)^n)$ as a good approximation to e_n and we are going to justify this view in the next few pages. Notice also that this relation between e_{n+1}^{-1} and e_n has the character that an upper bound on the e_j 's for $j \leq n$ is turned into a lower bound on the e_{n+1} 's and vice versa. As an application, we study the sequence $f_n = e_n(\frac{1}{4})$ whose asymptotic behavior will be needed later.

The f_n satisfy the recurrence

$$f_{n+1} = f_n(1 - f_n)$$
 with $f_0 = \frac{1}{2}$.

Hence, from Lemma 5,

$$\frac{1}{f_n} = n + 2 + \sum_{j < n} \frac{f_j}{1 - f_j}.$$
(8)

The f_n 's being positive, it follows that

$$1/f_n > n+2$$
 or $f_n < 1/(n+2)$.

Using this more precise etimate again in (8), we get

$$\frac{1}{f_n} < n+2 + \sum_{j < n} \frac{1}{j+2}.$$

Continuing this procedure, we see that

$$f_n = (n + \log(n) + O(1))^{-1},$$

and more precise estimates can be derived by iteration of this process.

LEMMA 6 (Convergence in sector around $\frac{1}{4}$). There exist positive constants ρ_0 , θ_0 such that the sequence $\{e_n(z)\}$ converges to 0 when z is such that

$$z \in D_0; \qquad |\varepsilon(z)| < \rho_0 \qquad and \qquad -((\pi/4) + \theta_0) < \operatorname{Arg} \varepsilon(z) < -((\pi/4) - \theta_0).$$

Proof. We only have to show that $e_{N(z)}(z)$ is small enough to satisfy the conditions of Lemma 1. For this purpose we use Lemma 5 to provide an upper bound on $|e_{N(z)}(z)|$.

We set $\varepsilon(z) = \rho e^{i\theta}$ and expand $(1 - \varepsilon(z))^{N(z)}$ in terms of ρ for small ρ when θ lies in some interval around $-\pi/4$ not containing 0. The following expansions are valid for ρ small enough and $\operatorname{Arg}(\varepsilon(z)) \neq 0$. They furthermore hold uniformly when θ is in any interval of the form $[-(\pi/4) - \lambda, -(\pi/4) + \lambda]$ with $0 < \lambda < \pi/4$:

$$|1 - \varepsilon(z)| = 1 - \rho \cos \theta + O(\rho^2),$$

$$\operatorname{Arg}(1 - \varepsilon(z)) = -\rho \sin \theta + O(\rho^2),$$

$$N(z) = \frac{-\alpha}{\rho \sin \theta} + O(1) \quad \text{with} \quad \alpha = \arccos \frac{1}{4},$$

$$|1 - \varepsilon(z)|^{N(z)} = e^{\alpha \cot \theta} + O(\rho).$$

In order to get an upper bound on e_N , we shall derive an asymptotic lower bound on the right-hand side of the relation giving $(1-\varepsilon)^n/e_n$ in Lemma 5, which we take as

$$\frac{(1-\varepsilon)^n}{e_n} = \frac{1-(1-\varepsilon)^n}{\varepsilon} + \frac{8}{3} + \frac{1}{3} + \sum_{1 \le j < n} \frac{e_j}{1-e_n} (1-\varepsilon)^j.$$

Since for $1 \leq j \leq N(z)$, $|e_j(z)| \leq \frac{1}{4}$, we have $|e_j/(1-e_j)| \leq \frac{1}{3}$ and

$$\left| \frac{1}{3} + \sum_{1 \leq j < N} e_j (1 - e)^j (1 - \varepsilon_j)^{-1} \right| \leq \frac{1}{3} + \frac{1}{3} \sum_{1 \leq j < N} |1 - \varepsilon|^j$$
$$\leq \frac{1}{3} (1 - |1 - \varepsilon|^N) / (1 - |1 - \varepsilon|)$$
$$< \frac{1}{3} (1 - e^{\alpha \cot \theta}) / (\rho \cos \theta) + O(1).$$

On the other hand

$$\begin{aligned} |(1 - (1 - \varepsilon)^N)/\varepsilon| &\ge (1 - |1 - \varepsilon|^N)/|\varepsilon| \\ &> (1 - e^{\alpha \cot \theta})/\rho + O(1). \end{aligned}$$

Thus for ρ small enough

$$\left|\frac{1-(1-\varepsilon)^{N}}{\varepsilon}\right| > \frac{8}{3} + \left|\frac{1}{3} + \sum_{1 \leq j < N} \frac{e_{j}}{1-e_{j}} (1-\varepsilon)^{j}\right|,$$

an inequality satisfied provided $\cos \theta > \frac{1}{3} + \delta$ for some $\delta > 0$, which we shall now assume.

We have thus shown

$$|1-\varepsilon|^N/|e_N| > (1-e^{\alpha\cot\theta})(\rho\cos\theta)^{-1}(\cos\theta-\frac{1}{3})(1+O(\rho)),$$

or equivalently

$$|e_N| < \rho |1-\varepsilon|^N (\cos \theta)(1-e^{\alpha \cot \theta})^{-1} (\cos \theta-\frac{1}{3})^{-1} (1+O(\rho)).$$

This estimate is to be compared to $|1-\varepsilon|^{-1}-1$ which is

 $|1-\varepsilon|^{-1}-1=\rho\cos\theta+O(\rho^2).$

Thus the convergence criterion is satisfied for ρ small enough provided

$$e^{\alpha \cot \theta} (1 - e^{\alpha \cot \theta})^{-1} (\cos \theta - \frac{1}{3})^{-1} < 1.$$

Equality is achieved for $-\theta = 0.819168... > \pi/4$ and inequality is ensured for all smaller values of $|\theta|$, which completes the proof of the lemma.

Again the convergence under the conditions of Lemma 6 is geometric except at $z = \frac{1}{4}$ and we can restate this lemma as

PROPOSITION 4. The function H(z) is analytic in a sector around $\frac{1}{4}$ defined by

$$\{z \neq \frac{1}{4}; |z - \frac{1}{4}| < \alpha_0 \text{ and } (\pi/2) - \beta_0 < |\operatorname{Arg}(z - \frac{1}{4})|\}$$

for some $\alpha_0, \beta_0 > 0$.

There does not seem to be any more straightforward argument to prove convergence of $e_n(z)$ to 0 in the domain described in Propositions 3 and 4. Actually, numerical computations indicate that the convergence of $e_n(z)$ is not monotonic in the whole of the convergence region, and the e_n 's display fairly erratic behavior away from the point $z = \frac{1}{4}$.

HEIGHTS OF BINARY TREES

5. Estimates of H(z) and the Average Height of Binary Trees

From the results of Sections 3 and 4 as summarized by Propositions 3 and 4, we now know that H(z) is analytic in an indented crown-shaped region depicted in Fig. 3. We proceed to evaluate the Taylor coefficient H_n of H(z) by means of Cauchy's integral formula

$$H_n = \frac{1}{2i\pi} \oint H(z) \frac{dz}{z^{n+1}},$$

selecting a contour inside that region which gives predominance to the behavior of the function around the singularity $\frac{1}{4}$. To do so, further information is required on the growth order of H(z) around $\frac{1}{4}$. After some preparation (Lemmas 7 and 8), we show that H(z) behaves there like a logarithm (Proposition 5). Once this is done, we are able to conclude the proof of Theorem B.

LEMMA 7 (Uniform bounds for $|e_n(z)|$ around $\frac{1}{4}$). There exist constants α_1, β_1 , and c_1 such that

$$|e_n(z)| < c_1/n$$

when $|z - \frac{1}{4}| < \alpha_1$ and $(\pi/2) - \beta_1 < |\operatorname{Arg}(z - \frac{1}{4})| < (\pi/2) + \beta_1$. Moreover, if $n \ge N(z)$, where N(z) is defined as in Lemma 4, then

$$|e_n(z)| < c_1 |\varepsilon(z)| |1 - \varepsilon(z)|^n.$$

Proof. We may suppose without loss of generality that $\text{Im } z \ge 0$. Suppose first that $1 \le n \le N(z)$. Let $\varepsilon(z) = \rho e^{i\theta}$. Proceeding as in the proof of Lemma 6, we find that

$$\frac{8}{3} + \left| \frac{1}{3} + \sum_{j=0}^{n-1} \frac{e_j}{1 - e_j} (1 - \varepsilon)^j \right| \leq \frac{8}{3} + \frac{1}{3} \frac{1 - |1 - \varepsilon|^n}{1 - |1 - \varepsilon|} < \frac{1}{3} \frac{1 - |1 - \varepsilon|^n}{\rho \cos \theta} + O(1).$$

while

$$|(1-(1-\varepsilon)^n)/\varepsilon| \ge (1-|1-\varepsilon|^n)/\rho.$$

Hence if $n \ge c_2$, then

$$\left|3+\sum_{j=0}^{n-1}\frac{e_j}{1-e_j}\left(1-\varepsilon\right)^j\right|<\frac{1}{2}\left|\frac{1-(1-\varepsilon)^n}{\varepsilon}\right|,$$

and so

$$\begin{aligned} |1-\varepsilon|^n/|e_n| \ge \frac{1}{2} |1-(1-\varepsilon)^n|/|\varepsilon|, \\ |e_n| \le 2 |\varepsilon| |1-\varepsilon|^n/(1-|1-\varepsilon|^n) \le 2\rho/(1-|1-\varepsilon|^n). \end{aligned}$$

We are considering $n \leq N(z)$, so

$$|1 - \varepsilon|^n = \exp(-n\rho\cos\theta + O(n\rho^2)) \ge 1 - \delta n\rho$$

for some $\delta > 0$, and so

 $|e_n| \leq 2/(\delta n).$

Since $|e_n| = O(n^{-1})$ for $n < c_2$, we find that

$$|e_n| \leq c_3 n^{-1}$$
 for $n \leq N(z)$.

Let us next suppose that n > N(z). Since we already know that $|e_j|$ is monotone decreasing for $j \ge N(z)$ (Lemmas 1 and 6),

$$|e_j| \leq c_3/N(z)$$
 for $j \geq N(z)$,

and therefore

$$\left| 3 + \sum_{j=0}^{n-1} \frac{e_j}{1 - e_n} (1 - \varepsilon)^j \right| \leq 3 + c_4 \sum_{j=1}^{N(z)} j^{-1} + \frac{c_3}{N(z)} \sum_{j=N(z)+1}^{n-1} |1 - \varepsilon|^j$$
$$\leq c_5 \log N(z) + c_6 N(z)^{-1} \rho^{-1} \leq c_7 \log \rho^{-1}.$$

On the other hand, $|1 - \varepsilon|^n \leq \frac{1}{2}$ for $n \geq N(z)$ and ρ small enough, so

$$|(1-(1-\varepsilon)^n)/\varepsilon| \ge (1-|1-\varepsilon|^n)/\rho \ge (2\rho)^{-1}.$$

Since $(2\rho)^{-1} > 2c_7 \log \rho^{-1}$ for ρ small enough,

$$|e_n| \leq 4\rho |1-\varepsilon|^n = 4 |\varepsilon| |1-\varepsilon|^n$$

for $n \ge N(z)$ if we make α_1 small enough. This proves the last part of Lemma 7. To complete the proof of the first part, we note that for $\varepsilon = \rho e^{i\theta}$, $(\pi/2) - \beta_1 < \operatorname{Arg}(z - \frac{1}{4}) < (\pi/2) + \beta_1$,

$$|\varepsilon| |1 - \varepsilon|^n \leq \rho (1 - \frac{1}{2}\rho)^n$$

and the maximum of $\rho(1-\frac{1}{2}\rho)^n$ as a function of ρ occurs at $\rho = 2(n+1)^{-1}$ and is $\leq 2(n+1)^{-1}$.

LEMMA 8 (Uniform bound for the convergence of $e_n(z)$ to $e_n(\frac{1}{4})$). There exist constants α_2 , β_2 , and c_2 such that

$$|e_n(z) - e_n(\frac{1}{4})| < c_2 |\varepsilon(z)|$$

when $|z - \frac{1}{4}| < \alpha_2$ and $(\pi/2) - \beta_2 < |\operatorname{Arg}(z - \frac{1}{4})| < (\pi/2) + \beta_2$.

Proof. Applying the estimate of Lemma 7 to the expansion given by Lemma 5 yields

$$\frac{(1-\varepsilon)^n}{e_n} = \frac{1-(1-\varepsilon)^n}{\varepsilon} + O\left(\sum_{j=1}^{\infty} j^{-1} |1-\varepsilon|^j\right)$$
$$= (1-(1-\varepsilon)^n)/\varepsilon + O(\log(1-|1-\varepsilon|)^{-1})$$
$$= (1-(1-\varepsilon)^n)/\varepsilon + O(\log|\varepsilon|^{-1}),$$

as well as the already known result

$$\frac{1}{e_n(\frac{1}{4})} = n + O\left(\sum_{j < n} j^{-1}\right) = n + O(\log n).$$

Hence for $n \leq N(z)$

$$(1-\varepsilon)^{n} \{e_{n}(\frac{1}{4})-e_{n}\} e_{n}^{-1}e_{n}(\frac{1}{4})^{-1} = (1-\varepsilon)^{n}/e_{n} - (1-\varepsilon)^{n}/e_{n}(\frac{1}{4})$$
$$= (1-(1-\varepsilon)^{n} - n\varepsilon(1-\varepsilon)^{n})/\varepsilon + O(\log|\varepsilon|^{-1})$$
$$= O(n^{2}|\varepsilon|).$$

Therefore

$$|e_n(\frac{1}{4}) - e_n| = O(n^2 |\varepsilon e_n e_n(\frac{1}{4})(1-\varepsilon)^{-n}|) = O(|\varepsilon|),$$

which proves the lemma for $n \leq N(z)$.

On the other hand, if n > N(z), then

$$|e_n|, |e_n(\frac{1}{4})| = O(n^{-1}) = O(|\varepsilon|),$$

so the lemma is trivial in this case.

With these lemmas, we proceed to determine the behavior of H(z) around $\frac{1}{4}$. Our previous developments suggest approximating $\sum_{n\geq 0} e_n(z)$ by

$$L(z) = \sum_{n \ge 1} \frac{\varepsilon (1-\varepsilon)^n}{1-(1-\varepsilon)^n}$$

in an appropriate region. To do this, we study the difference

$$D(z) = \sum_{n \ge 0} e_n(z) - L(z).$$

571/25/2-5

Using the expression for $e_n(z)$ given in Lemma 5, we see that

$$D(z) = \frac{1}{2} + \sum_{n \ge 1} e_n(z) \frac{S_n(z)}{q(n, z)},$$

where

$$q(n, z) = \frac{1 - (1 - \varepsilon(z))^n}{\varepsilon(z)}$$
 and $S_n(z) = 3 + \sum_{1 \le j < n} \frac{e_j(z)}{1 - e_j(z)} (1 - \varepsilon(z))^j$.

We notice that $D(\frac{1}{4})$ exists since the defining series converges as $\sum (\log n)/n^2$. We will show that $D(z) = D(\frac{1}{4}) + o(1)$ as $z \to \frac{1}{4}$ and will obtain an estimate of this o(1) term.

LEMMA 9 (First approximation lemma). For z in a neighborhood of $\frac{1}{4}$, with $|z - \frac{1}{4}| < \alpha_3$, $(\pi/2) - \beta_3 < \operatorname{Arg}(z - \frac{1}{4}) < (\pi/2) + \beta_3$,

$$D(z) = D(\frac{1}{4}) + O(|1 - 4z|^{1/4 - \eta})$$
 for any $\eta > 0$.

Proof. As in Lemma 8,

$$\left|3 + \sum_{j=1}^{n-1} \frac{e_j}{1 - e_j} (1 - \varepsilon)^j\right| = O\left(\sum_{j=1}^{\infty} j^{-1} |1 - \varepsilon|^j\right) = O(\log |\varepsilon|^{-1}).$$

We also have for $n \ge 3$, however,

$$\left|3 + \sum_{j=1}^{n-1} \frac{e_j}{1 - e_j} (1 - \varepsilon)^j\right| = O\left(3 + \sum_{j=1}^{n-1} j^{-1}\right) = O(\log n).$$

Therefore

$$(1-\varepsilon)^n/e_n = (1-(1-\varepsilon)^n)/\varepsilon + t_n,$$

where

$$t_n = t_n(z) = O(\log(\min(n, |\varepsilon|^{-1}))).$$

Hence, if n exceeds some fixed constant,

$$e_n/(1-\varepsilon)^n = \varepsilon/(1-(1-\varepsilon)^n) + O(|\varepsilon^2 t_n|/|1-(1-\varepsilon)^n|^2),$$

$$d_n = e_n - \varepsilon(1-\varepsilon)^n/(1-(1-\varepsilon)^n) = O(|\varepsilon^2 t_n||1-\varepsilon|^n/|1-(1-\varepsilon)^n|^2).$$

If $|\varepsilon|^{-1/2} \leq n \leq |\varepsilon|^{-1}$, then

$$d_n = O(|\varepsilon^2|\log(n)/|1 - (1-\varepsilon)^n|^2) = O(\log(n)/n^2)$$

If $n > |\varepsilon|^{-1}$,

$$d_n = O(|\varepsilon|^2 | 1 - \varepsilon|^n \log |\varepsilon|^{-1}).$$

Therefore

$$\sum_{n>|\varepsilon|^{-1/2}} d_n = O\left(\sum_{n>|\varepsilon|^{-1/2}} \frac{\log n}{n^2}\right) + O\left(|\varepsilon|^2 \log |\varepsilon|^{-1} \sum_{n\geq 0} |1-\varepsilon|^n\right)$$
$$= O(|\varepsilon|^{1/2} \log |\varepsilon|^{-1}).$$

Since for all $n \ge 2$

$$d_n(\frac{1}{4}) = O(\log(n)/n^2),$$

we find

$$\sum_{n > |\varepsilon|^{-1/2}} d_n - \sum_{n > |\varepsilon|^{-1/2}} d_n(\frac{1}{4}) = O(|\varepsilon|^{1/2} \log |\varepsilon|^{-1})$$

For $1 \leq n < |\varepsilon|^{-1/2}$,

$$\varepsilon(1-\varepsilon)^n/(1-(1-\varepsilon)^n)=n^{-1}+O(|\varepsilon|).$$

Therefore

$$\sum_{n < |\varepsilon|^{-1/2}} \left(d_n - d_n \left(\frac{1}{4} \right) \right) = \sum_{n < |\varepsilon|^{-1/2}} O\left(\left| e_n - e_n \left(\frac{1}{4} \right) \right| + \left| \frac{\varepsilon (1 - \varepsilon)^n}{1 - (1 - \varepsilon)^n} - \frac{1}{n} \right| \right)$$
$$= \sum_{n < |\varepsilon|^{-1/2}} O(|\varepsilon|) = O(|\varepsilon|^{1/2}),$$

which was to be shown.

The constant $D(\frac{1}{4})$ in Lemma 9 can be evaluated numerically as

$$D\left(\frac{1}{4}\right) = \frac{1}{2} + \sum_{n \ge 1}^{n} \left(e_n\left(\frac{1}{4}\right) - \frac{1}{n}\right),$$

and we find $D(\frac{1}{4}) = -1.602...$

To get the final expansion of H(z), we only need to estimate L(z). The observation that

$$\varepsilon/(1-(1-\varepsilon)^n) \to 1/n$$

for fixed n, when $\varepsilon \to 0$, suggests that L(z) behaves like

$$\sum_{n \ge 1} \frac{(1-\varepsilon)^n}{n} = \log \varepsilon,$$

which we are now going to justify rigorously.

Notice also that expanding in powers of $(1 - \varepsilon)$, we obtain

$$L(z) = \varepsilon(z) \sum_{m \ge 1} d(m)(1 - \varepsilon(z))^m,$$

with d(m) the divisor function of $m: d(m) = \sum_{d \mid m} 1$.

PROPOSITION 5 (Main approximation lemma for H(z)). For z in a sector around $\frac{1}{4}$:

$$|z - \frac{1}{4}| < \alpha$$
 and $(\pi/2) - \beta < |\operatorname{Arg}(z - \frac{1}{4})| = (\pi/2) + \beta$,

the following expansion holds for H(z):

$$H(z) = -2\log(1-4z) + K + O(|1-4z|^{\nu}) \quad \text{for any} \quad \nu < \frac{1}{4},$$

with $K \cong -4.1$, a constant.

Proof. It only remains to approximate the function

$$\sum_{n \ge 1} \frac{\varepsilon (1-\varepsilon)^n}{1-(1-\varepsilon)^n},$$

where z is in the specified region. Setting $1 - \varepsilon = e^{-u}$, this amounts to approximating

$$\sum_{n \ge 1} \frac{(1 - e^{-u})}{(1 - e^{-nu})} e^{-nu} = \frac{1 - e^{-u}}{u} \sum_{n \ge 1} u \frac{e^{-nu}}{1 - e^{-nu}}$$

when u is close to 0 and Arg u is close to $\pi/4$. To approximate this sum, we consider it as a Riemann sum relative to the integral

$$\int_u^\infty \frac{e^{-x}}{1-e^{-x}}\,dx.$$

Since the integral from 0 to ∞ is divergent, we split the sum according to whether n|u| < 1 or $n|u| \ge 1$, and compute the error terms separately.

For n such that $n |u| \ge 1$, we use the Taylor expansion

$$\left|\int_{nu}^{(n+1)u} \frac{e^{-x}}{1-e^{-x}} dx - \frac{ue^{-nu}}{1-e^{-nu}}\right| < \frac{|u|^2}{2} \max_{t \in [0,1]} \left|\frac{d}{dx} \frac{e^{-x}}{1-e^{-x}}\right|_{x=(n+t)u}$$

and summing, we see that

$$\sum_{n \ge |u|^{-1}} \frac{ue^{-nu}}{1 - e^{-nu}} = \int_{u[|u|^{-1}]}^{\infty} \frac{e^{-x}}{1 - e^{-x}} dx + O(|u|).$$

For *n* such that n|u| < 1, on the other hand, we expand $e^{-x}/(1-e^{-x})-x^{-1}$ which is differentiable and of bounded derivative over [0, 1], so that

$$\left|\frac{ue^{-nu}}{1-e^{-nu}} - \frac{1}{n} - \int_{nu}^{(n+1)u} \left(\frac{e^{-x}}{1-e^{-x}} - \frac{1}{x}\right) dx\right| < c |u|^2$$

for some constant c. Hence with $n_0 = |u|^{-1}$, we have

$$\sum_{n\geq 1} u \frac{ue^{-nu}}{1-e^{-nu}} = \sum_{n<1/|u|} \frac{1}{n} + \int_{u}^{n_0 u} \left(\frac{e^{-x}}{1-e^{-x}} - \frac{1}{x}\right) dx + \int_{n_0 u}^{\infty} \frac{e^{-x}}{1-e^{-x}} dx + O(|u|).$$

Approximating the harmonic series by the logarithm and changing the bounds of the integrals with only O(|u|) correcting terms, we see that (with γ the Euler constant):

$$\sum_{n \ge 1} u \frac{e^{-nu}}{1 - e^{-nu}} = -\log|u| + \gamma + \int_0^{u/|u|} \left(\frac{e^{-x}}{1 - e^{-x}} - \frac{1}{x}\right) dx + \int_{u/|u|}^{\infty} \frac{e^{-x}}{1 - e^{-x}} dx + O(|u|).$$

Using the Cauchy residue theorem we can change the path of integration to the real axis, and we have

$$\sum_{n \ge 1} u \frac{e^{-nu}}{1 - e^{-nu}} = -\log|u| + \gamma + \int_0^1 \left(\frac{e^{-x}}{1 - e^{-x}} - \frac{1}{x}\right) dx$$
$$+ \int_1^\infty \frac{e^{-x}}{1 - e^{-x}} dx - \int_1^{u/|u|} \frac{dx}{x} + O(|u|)$$
$$= -\log|u| - i\operatorname{Arg}(u) + \delta + O(|u|)$$
$$= -\log u + \delta + O(|u|),$$

with

$$\delta = \int_0^1 \left(\frac{e^{-x}}{1 - e^{-x}} - \frac{1}{x} \right) dx + \int_1^\infty \frac{e^{-x}}{1 - e^{-x}} dx + \gamma.$$

In fact, the two integrals cancel each other and we have $\delta = \gamma$. Since $\varepsilon = u + O(|u|^2)$ and $(1 - e^{-u})/u = 1 + O(|u|)$, we get

$$\sum_{n\geq 1}\frac{\varepsilon(1-\varepsilon)^n}{1-(1-\varepsilon)^n}=-\log\varepsilon+\gamma+O(|\varepsilon|).$$

Combining this with the approximation in Lemma 9 yields the result, with the constant K given by

$$K = 4D_n(\frac{1}{4}) + 4\gamma.$$

FLAJOLET AND ODLYZKO

To estimate the coefficients of H(z), we next translate the approximation of H(z) to an approximation of its coefficients. Since the result is of independent interest, e state it in a slightly more general form than strictly necessary here. The lemma is spired by [15] and may be compared to the classical Darboux method although the inditions of validity differ appreciably.

PROPOSITION 6 (Translation lemma). Let G(z) be analytic in a domain

$$D = \{z: z \neq \rho, |z| < \rho_1, |\operatorname{Arg}(z - \rho)| > \theta \text{ with } \rho_1 > \rho, \theta < \pi/2 \}.$$

ssume G(z) has the asymptotic expansion

$$G(z) = \gamma \log \left(1 - \frac{z}{\rho}\right) + \mu + \sum_{1 \le i \le m} \lambda_i \left(1 - \frac{z}{\rho}\right)^{\alpha_i} + O\left(\left|1 - \frac{z}{\rho}\right|^{\nu}\right)$$

ith $0 < \alpha_1 < \alpha_2 < \cdots < \alpha_m < v$, valid inside D. Then the nth Taylor coefficient G_n of f(z) admits the asymptotic expansion

$$G_n = \rho^{-n} \left[-\frac{\lambda}{n} + \sum_{\alpha_i + j < \nu} \frac{c_{ij}}{n^{\alpha_i + j + 1}} + O\left(\frac{1}{n^{\nu + 1}}\right) \right].$$

ų,

Proof. The *n*th Taylor coefficient can be computed using Cauchy's residue neorem as

$$G_n = \frac{1}{2i\pi} \int_{0^+} G(z) \frac{dz}{z^{n+1}},$$

there the contour simply encircles the origin and is inside the domain of analyticity f the function. For small $\omega > 0$ we take here the specific contour

$$\Gamma(\omega) = \Gamma_{0,\omega} \cup \Gamma_{1,\omega} \cup \Gamma_2,$$

efine for some fixed θ_1 and r_1 satisfying

$$\theta < \theta_1 < \pi/2$$
 and $\rho < r_1 < \rho_1$,

у

$$\begin{split} &\Gamma_{0,\omega} = \{z \colon |z-\rho| = \omega, \, |\operatorname{Arg}(z-\rho)| > \theta_1\}, \\ &\Gamma_{1,\omega} = \{z \colon |z-\rho| \ge \omega, \, |z| < r_1, \, |\operatorname{Arg}(z-\rho)| = \theta_1\}, \\ &\Gamma_2 = \{z \colon |z| = r_1, \, |\operatorname{Arg}(z-\rho)| \ge \theta_1\}. \end{split}$$

'he contour is depicted on Fig. 4.



FIG. 4. A diagram showing the contour $\Gamma(\omega)$.

We first show that we can let ω shrink to zero. As $\omega \to 0$, the integral

$$I(\omega) = \frac{1}{2i\pi} \int_{\Gamma_{0,\omega}} G(z) \frac{dz}{z^{n+1}} \to 0,$$

as can be seen from the inequality

$$|I(\omega)| \leq (\rho - \omega)^{-n-1} \max\{|G(z)| : z \in \Gamma_{0,\omega}\} \omega.$$

From the local expansion it follows that the upper bound vanishes as $\omega \to 0$. Letting $\Gamma = \Gamma(0)$ and $\Gamma_1 = \Gamma_{1,0}$, we thus see that G_n can be computed as

$$G_n = \frac{1}{2i\pi} \int_{\Gamma} G(z) \frac{dz}{z^{n+1}}.$$

The same argument applies to the functions in the local expansion of $G: \log(1-z/\rho)$ and the $(1-z/\rho)^{\alpha}$, showing that

$$-\frac{\rho^{-n}}{n} = \frac{1}{2i\pi} \int_{\Gamma} \log\left(1-\frac{z}{\rho}\right) \frac{dz}{z^{n+1}},$$
$$(-1)^n \rho^{-h} \binom{\alpha}{n} = \frac{1}{2i\pi} \int_{\Gamma} \left(1-\frac{z}{\rho}\right)^{\alpha} \frac{dz}{z^{n+1}}.$$

Hence,

$$G_n = \rho^{-n} \left[-\frac{\lambda}{n} + \sum_{1 \le i \le m} (-1)^n \binom{\alpha_i}{n} \right] + \frac{1}{2i\pi} \int_{\Gamma} R(z) \frac{dz}{z^{n+1}},$$

with

$$R(z) = G(z) - \lambda \log \left(1 - \frac{z}{\rho}\right) - \mu - \sum_{1 \le i \le m} \lambda_i \left(1 - \frac{z}{\rho}\right)^{\alpha_i}.$$

Now R(z) is analytic along Γ_2 and is $O(|1 - z/\rho|^{\nu})$ in D. Consider first the integral of R(z) along Γ_2 ; we have the obvious upper bound

$$\left|\frac{1}{2i\pi}\int_{\Gamma_2} R(z)\frac{dz}{z^{n+1}}\right| < \max\{|R(z)|: z \in \Gamma_2\} r_1^{-n}.$$

R(z) being analytic along Γ_2 is bounded, and this integral is exponentially small compared to ρ^{-n} , since $r_1 > \rho$. We are thus left with estimating integrals of the form

$$I_{\nu}(n) = \int_{\Gamma_2} \left| 1 - \frac{z}{\rho} \right|^{\nu} \frac{dz}{|z|^{n+1}}.$$

We set $z = \rho(1 + te^{i\phi})$ with $\phi = \pm \theta_1$ and t real. Using the symmetry of the contour, we have for some σ

$$I_{\nu}(n) = 2\rho^{-n} \int_{0}^{\sigma} \frac{t^{\nu} dt}{|1 + te^{i\phi}|^{n+1}} < 2\rho^{-n} \int_{0}^{\infty} \frac{t^{\nu} dt}{|1 + te^{i\phi}|^{n+1}}$$

Now $|1 + te^{i\phi}| = (1 + t^2 + 2\cos\phi)^{1/2}$ and since $\cos\phi > 0$, we have $(1 + t^2 + 2t\cos\phi)^{1/2} > 1 + \lambda t$ for some $\lambda > 0$, so that

$$I_{\nu}(n) < 2\rho^{-n} \int_{0}^{\infty} \frac{t^{\nu} dt}{(1+\lambda t)^{n+1}} < O\left(\rho^{-n} \int_{0}^{\infty} \frac{x^{\nu} dx}{(1+x)^{n+1}}\right).$$

To conclude with the bound we only need to show that $\int_0^\infty x^{\nu} dx/(1+x)^{n+1}$ is $O(1/n^{1+\nu})$. Indeed,

$$\int_0^\infty \frac{x^{\nu} \, dx}{(1+x)^{n+1}} = \int_0^1 \frac{x^{\nu} \, dx}{(1+x)^{n+1}} + O(2^{-n})$$

For $x \in [0, 1]$, $(1 + x) > e^{x/2}$, so that

$$\int_0^1 \frac{x^{\nu} \, dx}{(1+x)^{n+1}} < \int_0^1 x^{\nu} e^{-(n+1)x/2} \, dx < \int_0^\infty x^{\nu} e^{-(n+1)x/2} \, dx < \frac{\Gamma(\nu+1) \, 2^{\nu+1}}{(n+1)^{\nu+1}}$$

Hence,

$$I_{\nu}(n) = O(\rho^{-n}n^{-1-\nu}).$$

Putting everything together, we have thus shown that

$$G_n = \rho^{-n} \left[-\frac{\lambda}{n} + \sum_{1 \leq i \leq m} \lambda_i \begin{pmatrix} \alpha_i \\ n \end{pmatrix} (-1)^n \right] + O(\rho^{-n} n^{-\nu-1}).$$

To conclude the proof of the proposition, it only remains to examine the asymptotics of coefficients of the form

$$(-1)^n \binom{\alpha}{n} = (-1)^n \alpha(\alpha - 1) \cdots (\alpha - n + 1)/n! = n^{-1} \Gamma(n - \alpha) \Gamma(-\alpha)^{-1} \Gamma(n)^{-1}.$$

Known properties of the gamma function show the existence of an asymptotic expansion

$$(-1)^n \binom{\alpha}{n} \sim \sum_{j \ge 0} \frac{c_j(\alpha)}{n^{\alpha+j+1}},$$

with, in particular,

.

$$c_0(\alpha) = \Gamma(-\alpha)^{-1}.$$

Inserting these expansions into the estimate for G_n thus completes the proof of **P**roposition 6 with

$$c_{ij} = \lambda_i c_j(\alpha_i).$$

We have thus seen that adequate local information on a function G around its singularity leads to corresponding asymptotic information on its Taylor coefficients. The better the local approximation, the more terms the asymptotic expansion contains.

We can now complete the proof of Theorem B. Proposition 3 shows H(z) to be analytic outside the circle of convergence. Hence,

TABLE II

The Average Height of Binary Trees: Comparison of the Exact Values to the Asymptotic Estimates

		 In the second secon second second sec
n	\bar{H}_n	$\overline{H}_n(2\sqrt{\pi n})^{-1}$
10	7.07	0.631
20	11.29	0.712
50	19.97	0.797
100	29.98	0.846
200	44.29	0.883
500	72.94	0.920
1000	105.42	0.940
2000	151.50	0.956
5000	243.17	0.970
10000	346.64	0.978
16000	440.31	0.982

FLAJOLET AND ODLYZKO

THEOREM B. The average height of binary trees with n internal nodes satisfies

$$\overline{H}_n = 2\sqrt{\pi n} + O(n^{1/4+\eta}) \quad \text{for any} \quad \eta > 0.$$

Proposition 6 also shows that any improvement in the expansion of H(z) will lead to a better error term.

Numerical results corresponding to Theorem B are displayed in Table II. We notice that the convergence of \overline{H}_n to $2\sqrt{\pi n}$ is initially quite slow; however, for sizes of trees about 16,000, the gap appears to be less than 2%.

6. HEIGHTS IN SIMPLE FAMILIES OF TREES

Following Meir and Moon [13], we now consider planar trees with labels attached to nodes. All labels are taken from a fixed label set L

$$L = L_0 \cup L_1 \cup L_2 \cup \cdots,$$

with L_r the set of labels that may be attached to a node of degree r. We assume that each of the L_r is finite and we let c_r denote $|L_r|$; we can also assume without loss of generality that all the L_r 's are disjoint. A family defined in this way is said to be simple (or simply generated [13]). This definition obviously includes all families of unlabeled trees defined by restrictions on the set of allowed node degrees (in which case $c_r = 0$ or 1). It also covers all families of term trees, i.e., tree representations of expressions over an arbitrary set of operators. As examples, we mention

(a) the family of binary trees for which $c_0 = c_2 = 1$ and $c_r = 0$ for $r \neq 0, 2$; these have been considered in the previous sections;

(β) the family of general planar trees for which $c_r = 1$ for all $r \ge 0$: the analysis in [3] deals with these trees;

(γ) the family of unary-binary trees for which $c_0 = c_1 = c_2 = 1$ and $c_r = 0$ for r > 2; they appear as shapes of expression trees when unary as well as binary operations are allowed; the trees are counted by the Motzkin numbers;

(δ) the family of 2-3 trees (unbalanced) for which $c_0 = c_2 = c_3 = 1$ and $c_r = 0$ otherwise; their blanced counterparts are a useful data structure and have been counted by Odlyzko [15];

(c) the family of t-ary trees (which also appear in digital search); for these trees $c_r = 1$ if r = 0 or t and $c_r = 0$, otherwise.

As in the above examples, we shall restrict attention to those simple families for which there exists an absolute constant M such that

$$\forall r, \qquad c_r < M,$$

although our treatment also generalizes to sequences $\{c_r\}$ with a growth rate limited by an exponential.

Up to isomorphism, a simple family of trees is described by the sequence $\{c_r\}_{r\geq 0}$. Given a simple family E, we let y_n denote the number of trees of total size *n*; i.e., the number of trees formed with a total of *n* nodes. The generating function

$$y(z) = \sum_{n \ge 1} y_n z^n$$

satisfies an equation of the form

$$y(z) = z\phi(y(z)),$$
 where $\phi(y) = \sum_{n \ge 0} c_r y^r.$

Also, if we define

 $y_n^{[h]}$ = number of trees of size *n* and height $\leq h$,

with height measured by the number of nodes along the longest branch, then the generating functions

$$y^{[h]}(z) = \sum_{n} y^{[h]}_{n}(z)$$

are given by

$$y^{[0]}(z) = 0, \qquad y^{[h+1]}(z) = z\phi(y^{[h]}(z)).$$

The functions ϕ corresponding to cases $(\alpha)-(\varepsilon)$ are thus respectively,

 $1 + y^2$; $(1 - y)^{-1}$; $1 + y + y^2$; $1 + y^2 + y^3$; $1 + y^t$.

In the case of general planar trees, the $y^{[h]}(z)$ appear as convergents of a continued fraction, and additional algebraic information is available leading to explicit expressions for the $y^{[h]}(z)$; this is the basis of the treatment in [3].

In the binary case, there is a slight difference between the equation we obtain here, namely,

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

and the equation for B(z) which is

$$B(z) = 1 + zB(z)^2.$$

The two functions are related by

$$y(z)=zB(z^2),$$

which reflects the fact that in this section we consider total size measured by the total number of nodes (both nullary and binary).

The case of nonplanar labeled trees (with distinct labels) does not fall into our category of simple trees. It can, however, be subjected to the same analytical treatment since the exponential generating function

$$\hat{y}(z) = \sum y_n \frac{z^n}{n!}$$
 with y_n = number of trees of size n ,

satisfies the equation

$$\hat{y}(z) = z \exp(\hat{y}(z)),$$

with similar expressions relative to trees of bounded height. We shall thus obtain the Renyi and Szekeres result [17] as a consequence of our Theorem S.

We now indicate the lines along which the method employed for binary trees can be extended to these simple families of trees. Let

$$H_n = \sum_{h \ge 0} h(y_n^{[h]} - y_n^{[h-1]})$$

denote the total height of trees of size n, with the generating function

$$H(z)=\sum_{n\geq 0} H_n z^n.$$

We are interested in the average heights defined by

$$H_n = H_n / y_n,$$

provided $y_n \neq 0$. We proceed by proving that y(z) has algebraic singularities on its circle of convergence [13], and that H(z) has corresponding logarithmic singularities.

We have to distinguish two cases based on the value of

$$d = \operatorname{GCD}\{r: c_r \neq 0\}.$$

The situation where d = 1 (planar trees, unary-binary trees,...) is the simplest one since then y has only one singularity on its circle of convergence; in this case, $y_n \neq 0$ for all $n \ge n_0$. The situation where $d \ne 1$ (binary trees, *t*-ary trees,...) requires combining results relative to each of the d singularities of y on its circle of convergence; in that case, $y_n = 0$ if $n \ne 1 \pmod{d}$.

Case 1 (Unicity of singularity). We start again with the equation

$$y(z) = z\phi(y(z))$$

and look for the point where the implicit function theorem ceases to apply. This occurs when

$$\frac{d}{dy}\left(\frac{y}{\phi(y)}\right) = 0,$$
 i.e., $\phi(y) = y\phi'(y).$

Let τ be the value of smallest modulus such that $\phi(\tau) = \tau \phi'(\tau)$. The GCD condition implies that τ is unique and real: let $\rho = \tau/\phi(\tau)$ be the corresponding value of z. For (z, y) in a neighborhood of (ρ, τ) satisfying $y = z\phi(y)$, a local expansion shows that

$$z - \rho = -(y - \tau)^2 \phi''(\tau) \tau / (2\phi^2(\tau)) + O(|y - \tau|^3).$$

Hence, around $z = \rho$, y behaves as

$$\tau - (2\phi(\tau)/\phi''(\tau))^{1/2} (1 - z/\rho)^{1/2}$$

and its nth Taylor coefficient is asymptotic to

$$c_1 \rho^{-n} n^{-3/2}$$
 with $c_1 = (\phi(\tau)/2\pi\phi''(\tau))^{1/2}$.

This is essentially the Darboux-Polya theorem applied to tree enumerations (see [13]).

Starting from the two equations

$$y(z) = z\phi(y(z));$$
 $y^{[h+1]}(z) = z\phi(y^{[h]}(z)),$

and subtracting, we get

$$y(z) - y^{[h+1]}(z) = z(\phi(y(z)) - \phi(y^{[h]}(z))).$$

Using the Taylor expansion of the right-hand side of this equation around y(z), we see that

$$y - y^{[h+1]} = z(y - y^{[h]}) \phi'(y) [1 - (y - y^{[h]}) \phi''(y)/(2\phi'(y)) + O(|y - y^{[h]}|^2)].$$

When $z = \rho$, $z\phi'(y) = 1$; expanding $z\phi'(y)$ around ρ , we get

$$z\phi'(y) = 1 + (y - \tau)\tau\phi''(\tau)/\phi(\tau) + O(|y - \tau|^2)$$

= 1 - (1 - z/\rho)^{1/2}\tau(2\phi''(\tau)/\phi(\tau))^{1/2} + O(|y - \tau|^2)

Thus setting $e_h(z) = y(z) - y^{[h]}(z)$, and $1 - z\phi'(y) = \varepsilon(z)$, we see that

$$e_{h+1}(z) = (1 - \varepsilon(z)) e_h(z)(1 - \phi''(\tau) e_h(z)/(2\phi'(\tau)) + O(|e_h^2(z)| + |e_h(z)(y - \tau)|)),$$

where

$$\varepsilon(z) = (1 - z/\rho)^{1/2} \tau (2\phi''(\tau)/\phi(\tau))^{1/2} + O(|y - \tau|^2).$$

The situation is now quite similar to what we had before. Taking reciprocals and applying the old trick leads to the approximate expression

$$e_n(z) \sim c_2 \varepsilon(z)(1-\varepsilon(z))^n/(1-(1-\varepsilon(z))^n)$$

with $c_2 = 2\phi'(\tau)/\phi''(\tau)$. Hence $H(z) = \sum_{n \ge 0} e_n(z)$ behaves around its singularity $z = \rho$ like $c_2 \log \varepsilon(z)$ and

$$H_n \sim \frac{1}{2} c_2 \rho^{-n} n^{-1},$$

or equivalently

$$\overline{H}_n \sim \frac{1}{2} (c_2/c_1) n^{1/2}.$$

Case 2 (Multiple singularities). We now assume that $d = \text{GCD}\{n: c_n \neq 0\}$ is nontrivial $(d \neq 1)$. The equation

$$y = z\phi(y)$$

can then be put in the form

$$y = z\psi(y^d)$$

with $\psi(u) = \phi(u^{1/d})$ a power series in u. The previous computations apply here: if τ is the smallest positive root of the equation

$$\phi(\tau) = \tau \phi'(\tau),$$

then y(z) has an algebraic singularity at τ . Now, since $\phi(y)$ depends only on y^d , we see that y(z) also has singularities at the points

$$\tau_j = \omega^j \tau$$
 for $j = 0, 1, ..., d-1$,

where ω is a primitive *d*th root of unity. Setting as before

$$\rho = \tau/\phi(\tau),$$

these singularities correspond to values of z

$$\rho_i = \omega^j \rho$$
.

Local expansions for y can also be carried out around the ρ_i showing that

$$z - \rho_j = -\omega^j (y - \tau_j)^2 \, \phi''(\tau) \, \tau/(2\phi^2(\tau)) + O(|y - \tau_j|^3).$$

Hence, around $z = \rho_j$, the approximation of y is

$$\tau_j - \omega^j (2\phi(\tau)/\phi''(\tau))(1-z/\rho_j)^{1/2}.$$

The nth Taylor coefficient of this expansion is approximated by

$$c_1 \rho^{-n} \omega^{-j(n-1)} n^{-3/2}$$
 with $c_1 = (\phi(\tau)/(2\pi\phi''(\tau)))^{1/2}$,

and provided $n \equiv 1 \pmod{d}$ —which is to be assumed since $y_n = 0$ if $n \not\equiv 1 \pmod{d}$ —these terms add up to

$$dc_1\rho^{-n}n^{-3/2}.$$

The same phenomenon occurs for H(z) which also has d singularities on its circle of convergence. Around $z = \rho_j$, H(z) behaves as

$$\frac{1}{2}c_2\omega^j\log(1-z/\rho_j),$$

so that for $n \equiv 1 \pmod{d}$

$$H_n \sim (d/2) c_2 \rho^{-n} n^{-1}.$$

Hence again

$$\bar{H}_n \sim \frac{1}{2} (c_2/c_1) n^{1/2}.$$

We can thus state:

THEOREM S. For simple families of trees corresponding to the equation $y = z\phi(y)$, and for $n \equiv 1 \pmod{d}$ with $d = \text{GCD}\{r: c_r \neq 0\}$, the average heights satisfy

$$\overline{H}_n \sim \lambda n^{1/2},$$

where

$$\lambda = (2\pi/(\phi(\tau) \phi''(\tau)))^{1/2} \phi'(\tau),$$

and τ is the smallest positive root of the equation

$$\phi(\tau)-\tau\phi'(\tau)=0.$$

COROLLARY,

(i) The average height of a unary-binary tree with n nodes is asymptotic to

 $\sqrt{3\pi n}$.

(ii) The average height of an unbalanced 2-3 tree with n nodes is asymptotic to

$$\sqrt{\pi n(2+3\tau)/(1+3\tau)},$$

where τ is the positive root of the equation $2\tau^3 + \tau^2 - 1 = 0$.

(iii) The average height of a t-ary tree with n internal (t-ary) nodes is asymptotic to

$$\sqrt{2\pi nt/(t-1)}$$
.

(iv) The average height of a (planar rooted) tree with n nodes [3] is asymptotic to

 $\sqrt{\pi n}$.

(v) The average height of a labeled nonplanar tree with n nodes [17] is asymptotic to

$$\sqrt{2\pi n}$$
.

7. DISTRIBUTION RESULTS

In this section, we shall show that our methods can be extended to derive information about the distribution of heights in simple families of trees. We shall deal with the binary case giving asymptotic equivalents for moments of higher order (variance, etc.). The distribution of heights in trees appears to obey a limiting theta distribution. A similar result has been proved by Renyi and Szekeres [17] in the case of labeled nonplanar trees using a rather different method, and in the case of general planar trees by Kemp [10] using the explicit enumeration results available in that particular case. We prove

THEOREM MB (Moments of the distribution of height in binary trees). The nth moment of the distribution of heights in binary trees of size n satisfies, for $r \ge 2$,

$$\overline{M}_{r,n} \sim 2^r r(r-1) \, \Gamma(r/2) \, \zeta(r) \, n^{r/2} \qquad as \quad n \to \infty.$$

Proof. The *r*th moment of the distribution of heights intrees of size *n* is giving by

$$\overline{M}_{r,n} = \frac{M_{r,n}}{B_n}$$
 with $M_{r,n} = \sum_{h \ge 1} h^r (B_n^{[h]} - B_n^{[h-1]}).$

The quantities $M_{r,n}$ are estimated from their generating functions:

$$M_r(z) = \sum_{n \ge 0} M_{r,n} z^n,$$

with

$$M_{r}(z) = \sum_{h \ge 1} h^{r} (B^{[h]}(z) - B^{[h-1]}(z)).$$

We only need to consider here the case where r > 1. Expressing M_r in terms of the e_n 's and ε , we get

$$M_{r}(z) = \frac{4}{1+\varepsilon(z)} \sum_{h \ge 1} h^{r}(e_{h-1}(z) - e_{h}(z)) = \frac{4}{1+\varepsilon(z)} \sum_{h \ge 0} ((h+1)^{r} - h^{r}) e_{h}(z),$$

using summation by parts. Hence setting

$$S_r(z) = \sum_{h \ge 1} h^r e_h(z),$$

we see that

$$M_r(z) = (4/(1+\varepsilon))\left[rS_{r-1} + \binom{r}{2}S_{r-2} + \binom{r}{3}S_{r-3} + \cdots\right].$$

The problem thus reduces (for each r) to estimating the order of $S_r(z)$ around the singularity $\frac{1}{4}$. From this information, the asymptotic behavior of the $M_{r,n}$ is recovered by methods similar to Proposition 6.

We first compare $S_r(z)$ with the simpler function

$$T_r(z) = \sum_{n \ge 1} n^r \frac{\varepsilon (1-\varepsilon)^n}{1-(1-\varepsilon)^n}.$$

To do so, we study the difference $S_r - T_r$ using the tools of Lemma 9. The summation giving $S_r - T_r$ is split into

$$S_r - T_r = \sum_{n < |\varepsilon|^{-1/2}} n^r d_n + \sum_{|\varepsilon|^{-1/2} \le n < |\varepsilon|^{-1}} n^r d_n + \sum_{n \ge |\varepsilon|^{-1}} n^r d_n = U_1 + U_2 + U_3$$

with

$$d_n = e_n - \varepsilon (1-\varepsilon)^n / (1-(1-\varepsilon)^n).$$

With the estimates for d_n previously derived, we find:

(i) $U_1 = O(\sum_{n < |\varepsilon|^{-1/2}} n^r \log((n)/n^2))$, using $|\varepsilon(1-\varepsilon)^n/(1-(1-\varepsilon)^n)| = n^{-1} + O(\varepsilon)$ and $t_n = O(\log\min(n, |\varepsilon|^{-1}))$.

(ii) $U_2 = O(\sum_{|\varepsilon|^{-1/2} \le |\varepsilon|^{-1}} n^r \log(n)/n^2)$, using $d_n = O(\log(n)/n^2)$ in this range. Hence, $U_1 + U_2 = O(|\varepsilon|^{-r+1} \log |\varepsilon|^{-1})$.

(iii) $U_3 = O(|\varepsilon|^2 \log |\varepsilon|^{-1} \sum_{n > |\varepsilon|^{-1}} n^r |1 - \varepsilon|^n) = O(|\varepsilon|^{-r+1} \log |\varepsilon|^{-1}), \text{ using } d_n = O(|\varepsilon|^2 |1 - \varepsilon|^n \log |\varepsilon|^{-1}).$

We have thus shown

$$|S_r - T_r| = O(|\varepsilon|^{-r+1} \log |\varepsilon|^{-1}),$$

a difference of a smaller order than T_r , as we now prove.

\$71/25/2-6

Notice first in expanding T_r that

$$T_r = \varepsilon \sum_{n \ge 1} n^r \frac{(1-\varepsilon)^n}{1-(1-\varepsilon)^n} = \varepsilon \sum_{n \ge 1} \sigma_r(n)(1-\varepsilon)^n$$

where $\sigma_r(n)$ is the sum of the kth powers of the divisors of n

$$\sigma_r(n) = \sum_{d \mid n} d^r,$$

with corresponding Dirichlet generating function $\zeta(s) \zeta(s-r)$.

A function like

$$F_r(u) = \sum \sigma_r(n) e^{-nu}$$

can be evaluated asymptotically, for real $u \rightarrow 0$ by appealing to properties of the Mellin transform as in [3]. The Mellin transform is readily found to be

$$F_r^*(s) = \zeta(s) \, \zeta(s-r) \, \Gamma(s),$$

whose rightmost pole is at s = r + 1. Residue computation now shows that

$$F_r(u) = \zeta(r+1) \, \Gamma(r+1) \, u^{-r-1} + O(|u|^{-1}),$$

from which $T_r(z)$ can be estimated when ε is real.

To extend this evaluation to complex z and ε , we use the method of Lemma 10 We set again $e^{-u} = 1 - \varepsilon$, and

$$T_{r} = \varepsilon \sum_{n \ge 1} n^{r} \frac{e^{-nu}}{1 - e^{-nu}} = \varepsilon u^{-r-1} \sum_{n \ge 1} (nu)^{r} \frac{e^{-nu}}{1 - e^{-nu}} u.$$

The sum is a Riemann sum relative to the integral

$$c_r = \int_0^\infty x^r \, \frac{e^{-x}}{1 - e^{-x}} \, dx,$$

the integrand being of bounded derivative over the interval. We thus have

$$T_r = c_r \varepsilon u^{-r-1} (1 + O(|u|)),$$

and translating back in terms of ε , we get

$$T_r(z) = c_r \varepsilon^{-r} + O(|\varepsilon|^{-r+1}).$$

To compute c_r , it suffices to expand $(1 - e^{-x})^{-1}$, and determine separately each integral. One finds

$$c_r = \Gamma(r+1)\,\zeta(r+1).$$

Returning to M_r , we have thus obtained the local expansion

$$M_r(z) = 4r\zeta(r) \Gamma(r) \varepsilon^{-r+1} + O(|\varepsilon|^{-r+2} \log |\varepsilon|).$$

To conclude with the asymptotic growth of the $M_{r,n}$, we again need a translation lemma analogous to Proposition 6. In Proposition 6, the remainder term in the expansion of the function is small near the singularity. This is no longer the case now, and so we use a different contour to obtain the following result:

PROPOSITION 7. Suppose that g(z) is analytic in

$$E = \{z \colon |z| \leq \rho, \ z \neq \rho\}$$

for some $\rho > 0$, and that for $z \in E$,

$$g(z) = O(|1 - z/\rho|^{\alpha})$$

for some $\alpha < 0$. Then, the nth Taylor series coefficient g_n of g(z) satisfies

$$g_n = O(\rho^{-n}n^{-\alpha-1})$$

Proof. We use Cauchy's theorem with the contour $\Gamma = \Gamma_0 \cup \Gamma_1$, where

$$\Gamma_0 = \{z \colon |z - \rho| = 1/n, \, |z| \leq \rho\}, \qquad \Gamma_1 = \{z \colon |z| = \rho, \, |z - \rho| \ge 1/n\}.$$

On Γ_0 ,

$$|g(z)/z^{n+1}| = O((\rho - n^{-1})^{-n-1}(n\rho)^{-\alpha}) = O(\rho^{-n}n^{-\alpha}),$$

and so

$$\frac{1}{2\pi i} \int_{\Gamma_0} \frac{g(z) \, dz}{z^{n+1}} = O(\rho^{-n} n^{-\alpha - 1}).$$

Let θ_0 be determined by $0 < \theta_0 < \pi/2$,

$$\rho |1-e^{i\theta_0}|=1/n.$$

Then

$$\frac{1}{2\pi i}\int_{\Gamma_1}\frac{g(z)\,dz}{z^{n+1}}=O\left(\rho^{-n}\int_{\theta_0}^{\pi}|1-e^{i\theta}|^{\alpha}\,d\theta\right).$$

Now $|1 - e^{i\theta}| > c |\theta|$ for some fixed c > 0 if $|\theta| \le \pi$, so the term on the right side above is

$$O\left(\rho^{-n}\int_{\theta_0}^{\pi}\theta^{\alpha} d\theta\right) = O(\rho^{-n}\theta_0^{\alpha+1}).$$

Since $\theta_0 \sim c' n^{-1}$ as $n \to \infty$ for some c' > 0, we obtain the claim of the proposition.

Applying this proposition to the error term in the expansion of $M_r(z)$, and using the explicit expressions for the coefficients of ε^{-r} , we obtain

$$M_{r,n} = 4r\zeta(r) \Gamma(r) 4^n \binom{(1-r)/2}{n} + O(4^n n^{r/2 - 5/2 + n}),$$

for any $\eta > 0$. Since for fixed nonintegral α

$$\binom{\alpha}{n} = \Gamma(-\alpha)^{-1} n^{-\alpha-1} (1 + O(1/n)),$$

we find

$$M_{r,n} \sim 4r\zeta(r) \Gamma(r) \Gamma((r-1)/2)^{-1} 4^n n^{r/2-3/2}$$

Dividing by B_n , we finally get

$$\overline{M}_{r,n} \sim 4\pi^{1/2} r \frac{\Gamma(r) \zeta(r)}{\Gamma((r-1)/2)} n^{r/2},$$

which using the duplication formula for the gamma function yields

$$\overline{M}_{r,n} \sim 2^r r(r-1) \Gamma(r/2) \zeta(r) n^{r/2}.$$

For $n = 10^4$, the asymptotic estimates of the 2nd, 3rd, and 4th moment are within 10% of the actual values.

Now we consider the *normalized height* defined for a binary tree of size n by

$$\bar{h}(t) = \text{height}(t)/(2\sqrt{n}).$$

The *r*th moment $\mu_{r,n}$ of \overline{h} on trees of size *n* satisfies

$$\mu_{r,n} \to r(r-1) \Gamma(r/2) \zeta(r)$$
 as $n \to \infty$,

with error terms essentially O(1/n). (The formula is seen to be still valid for r = 1, if we take limits.) We thus see that normalized height converges to a distribution whose rth moment is given by

$$r(t-1) \Gamma(r/2) \zeta(r).$$

The limit distribution is identified by comparing these quantities with the moments of the theta distribution [17], whose cumulative distribution function is

$$H(x) = 4x^{-3}\pi^{5/2} \sum_{k \ge 0} k^2 e^{-k^2\pi^2/x^2} = \sum_{-\infty < k < +\infty} e^{-k^2x^2} (1 - 2k^2x^2)$$

with corresponding density

$$h(x) = 4x \sum_{k \ge 1} k^2 (2k^2 x^2 - 3) e^{-k^2 x^2}.$$

The rth moment of this distribution is precisely

$$\mu_r = r(r-1) \, \Gamma(r/2) \, \zeta(r).$$

We have thus proved

COROLLARY. The normalized height

$$\bar{h}(t) = \text{height}(t)/(2\sqrt{n})$$

on trees of size n admits a limiting theta distribution with density function

$$h(x) = 4x \sum_{k \ge 1} k^2 (2k^2 x^2 - 3) e^{-k^2 x^2}$$

as $n \to \infty$.

The same principle applies to simple families of trees, and one finds for the rth moment relative to trees of size n an asymptotic expression of the form

$$\xi^{r/2} r(r-1) \Gamma(r/2) \zeta(r) n^{r/2}$$

which again shows that, suitably normalized, the distributions of heights tend to a theta distribution.

THEOREM MS (Moments of the distribution of height in simple trees). For simple families of trees corresponding to the equation $y = z\phi(y)$, the rth moment of height in trees of size n is asymptotic to

$$r(r-1) \Gamma(r/2) \zeta(r) \xi^{r/2} n^{r/2}$$
 with $\xi = 2\phi'(\tau)^2/(\phi(\tau) \phi''(\tau)).$

The distribution of the normalized height in trees of size n

$$\bar{h}(t) = \text{height}(t) / \sqrt{\xi n}$$

tends to the limiting theta distribution of density h(x).

8. Conclusions

To conclude, we observe that many combinatorial problems—especially tree enumerations—have generating functions associated to functional equations of the form

$$f(z) = \Phi(z, f(z)),$$

FLAJOLET AND ODLYZKO

where Φ is a functional reflecting the structural definition of the objects. The approximations provided by the iterative scheme

$$f^{[0]}(z) = 0;$$
 $f^{[h+1]}(z) = \Phi(z, f^{[h]}(z))$

are often of combinatorial significance, representing a partition of the objects according to some form of *height*. In this paper we dealt with equations of the form

$$f(z) = z\phi(f(z))$$

corresponding to simple families of trees.

The enumeration of nonplanar unlabeled rooted trees corresponds to functional equations of the form

$$f(z) = z e^{f(z) + (1/2)f(z^2) + (1/3)f(z^3) + \cdots}$$

as appears from developments in Polya theory. The present approach is applicable since the occurrence of $f(z^2)$; $f(z^3)$,... is known not to affect singularities too much and f(z) still has an algebraic singularity on its circle of convergence (see Polya [16]).

On the other hand, the statistics about binary search trees and tournament trees represent equations of a different nature with probable singularities of the type of $(1/(1-z))\log(1/(1-z))$. We mention here the two equations

$$T(z) = 1 + \int_0^z T^2(z) dz$$
 and $T(z) = \exp \int_0^z T(z) dz$,

whose approximations provided by the iterative scheme are associated with, respectively, height and one-sided height. The methods developed here do not seem to apply to these problems.

Another line of extension of our methods is to look at different limit distributions. In another work, the authors have shown that the limit distribution of binary trees of given height by size is Gaussian. The proof there is achieved by applying the saddle point method and investigating the analytical properties of the $B^{[h]}(z)$ outside the circle of convergence where they display a doubly exponential growth.

Finally we mention that other methods applicable to large classes of trees have already received some attention: Meir and Moon [13] have shown that path length in simple families of trees is essentially $\sim \alpha n \sqrt{n}$; Odlyzko [15] has dealt with functional equations of a general nature relative to balanced trees; Flajolet and Steyaert [7] have shown that the simple backtracking algorithm for tree matching has linear average time when inputs are taken from any simple family of trees.

HEIGHTS OF BINARY TREES

REFERENCES

- 1. E. A. BENDER, Asymptotic methods in enumeration, SIAM Rev. 16 (1974), 485-515.
- 2. N. DE BRUIJN, "Asymptotic Methods in Analysis," North-Holland, Amsterdam, 1961.
- 3. N. DE BRUIJN, D. KNUTH, AND S. RICE, The average height of planted plane trees, in "Graph Theory and Computing" (R-C. Read, Ed.), pp. 15-22, Academic Press, New York, 1972.
- 4. P. FLAJOLET, "On the Performance Evaluation of Extendible Hashing and Tree Searching," to be published.
- 5. P. FLAJOLET AND A. M. ODLYZKO, Exploring binary trees and other simple trees, in "Proceedings of 21st IEEE Found. Computer Sci. Symposium," New York, 1980, pp. 207–216.
- 6. P. FLAJOLET, J. C. RAOULT, AND J. VUILLEMIN, The number of registers required to evaluate arithmetic expressions, *Theoret. Comput. Sci.* 9 (1979), 99-125.
- 7. P. FLAJOLET AND J. M. STEYAERT, On the analysis of tree matching algorithms, in "Proceedings, 7th ICALP Conf.," Amsterdam, 1980.
- 8. R. KEMP, The average number of registers needed to evaluate a binary tree optimally, Acta Inform. 11 (1979), 363-372.
- 9. R. KEMP, "The Average Height of a Derivation Tree Generated by a Linear Grammar in a Special Chomsky Normal Form," Saarbrucken University Report A 78/01, 1978.
- 10. R. KEMP, On the stack size of regularly distributed binary trees, in "Proceedings, 6th ICALP Conf.," Udine 1979.
- 11. D. E. KNUTH, "The Art of Computer Programming: Fundamental Algorithms," Addison-Wesley, Reading, Mass., 1968.
- 12. D. E. KNUTH, "The Art of Computer Programming: Sorting and Searching," Addison-Wesley, Reading, Mass., 1973.
- 13. A. MEIR AND J. W. MOON, On the altitude of nodes in random trees, Canad. J. Math 30 (1978), 997-1015.
- 14. A. MEIR, J. W. MOON, AND J. R. POUNDER, On the order of random channel networks, SIAM J. Algebraic Discrete Math. 1 (1980), 25-33.
- 15. A. ODLYZKO, Periodic oscillations of coefficients of power series that satisfy functional equations, Adv. in Math. 44 (1982), 180-205.
- 16. G. POLYA, Kombinatorische Anzahlbestimmungen für Graphen, Gruppen, und Chemische Verbindungen, Acta Math. 68 (1937), 145-254.
- 17. A. RENYI AND G. SZEKERES, On the height of trees, Austral. J. Math. 7 (1967), 497-507.
- 18. J. RIORDAN, The enumeration of trees by height and diameter, IBM J. Res. Dev. 4 (1960), 473–478.
- 19. J. M. ROBSON, The height of binary search trees, Austral. Comput. J. 11 (1979), 151-153.
- 20. J. M. RONSON, "The Asymptotic Behaviour of the Height of Binary Search Trees," to be published.