

AD 731038

THE AVERAGE HEIGHT OF PLANTED PLANE TREES

Details of illustrations in
this document may be be
studied on microfiche

BY

N. G. DE BRUIJN

D. E. KNUTH

S. O. RICE

GCT 30-100

STAN-CS-218-71

MAY, 1971



COMPUTER SCIENCE DEPARTMENT

School of Humanities and Sciences

STANFORD UNIVERSITY

Incorporated by
NATIONAL TECHNICAL
INFORMATION SERVICE
Springfield, Va. 22161



Unclassified

Security Classification

DOCUMENT CONTROL DATA - R & D

(Security classification of title, body of abstract and indexing annotation must be entered when the overall report is classified)

| | |
|--|------------------------------------|
| 1. ORIGINATING ACTIVITY (Corporate Author) | 2a. REPORT SECURITY CLASSIFICATION |
| Stanford University | Unclassified |
| | 2b. GROUP |

3. REPORT TITLE

THE AVERAGE HEIGHT OF PLANTED PLANE TREES

4. DESCRIPTIVE NOTES (Type of report and inclusive dates)

Technical Report, May 1971

5. AUTHORED BY (First name, middle initial, last name)

N. G. de Bruijn, D. E. Knuth, and S. O. Rice

| | | |
|--|---|-----------------|
| 6. REPORT DATE | 7a. TOTAL NO. OF PAGES | 7b. NO. OF REPS |
| May 1971 | 7 | 2 |
| 8a. CONTRACT OR GRANT NO. | 8b. ORIGINATOR'S REPORT NUMBER (If any) | |
| ONR N-00014-67-A-0112-0057 NR 044-402. | STAN-CS-71-218 | |
| 8b. PROJECT NO. | 8c. OTHER REPORT NO(S) (Any other numbers that may be assigned to the report) | |
| NSF GJ-992 | None | |

10. DISTRIBUTION STATEMENT

Releasable without limitations on dissemination.

| | |
|---|--|
| 11. SUPPLEMENTARY NOTES | 12. SPONSORING MILITARY ACTIVITY |
| Details of Illustrations in this document may be best plotted on microfiche | Mathematics Program Office of Naval Research Arlington, Virginia 22217 |

13. ABSTRACT

An asymptotic expression for the average height of a planted plane tree is derived, based on an asymptotic series for sums such as $\sum_{k \geq 1} \binom{2n}{n+k} d(k)$ and $\sum_{k \geq 1} e^{-k^2/n} d(k)$, where $d(n)$ is the number of divisors of n .

$$1. \quad \dots + \frac{1}{k} + \dots + \left(\dots + \frac{1}{k} + \dots + \frac{1}{d(k)} \right) + \dots$$

Unclassified

Security Classification

| 14 KEY WORDS | LINK A | | LINK B | | LINK C | |
|--|--------|----|--------|----|--------|----|
| | ROLE | WT | ROLE | WT | ROLE | WT |
| Planted plane tree, height of tree, ballot problem, random walk, asymptotic series, divisor function, gamma function, zeta function. | | | | | | |

DD FORM 1 NOV 1973 (BACK)
(PAGE 2)

Unclassified

Security Classification

The Average Height of Planted Plane Trees

by N. G. de Bruijn, D. E. Knuth,^{*} and S. O. Rice

Abstract

An asymptotic expression for the average height of a planted plane tree is derived, based on an asymptotic series for sums such as $\sum_{k \geq 1} \binom{2n}{n+k} d(k)$ and $\sum_{k \geq 1} e^{-k^2/n} d(k)$, where $d(n)$ is the number of divisors of n .

AMS 1970 subject classifications: Primary 05C05, 10H25; 95A15, 10A40

Keywords: Planted plane tree, height of tree, ballot problem, random walk, asymptotic series, divisor function, gamma function, zeta function.

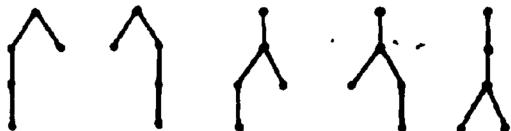
* This research was supported in part by the National Science Foundation, under grant number GJ-992, and the Office of Naval Research under grant number N-00014-67-A-0112-0057 NR 044-402. Reproduction in whole or in part is permitted for any purpose of the United States Government.

The Average Height of Planted Plane Trees

by N. G. de Brujin, D. E. Knuth, and S. O. Rice

A planted plane tree (sometimes called an ordered tree) is a rooted tree which has been embedded in the plane so that the relative order of subtrees at each branch is part of its structure. In this paper we shall say simply "tree" instead of "planted plane tree", following the custom of computer scientists.

The height of a tree is the number of nodes on a maximal simple path starting at the root. For example, there are exactly 5 trees with five nodes and height 4, namely



The height of a tree is of interest in computing because it represents the maximum size of a stack used in algorithms that traverse the tree [5, p. 317-318]. Our goal in this paper is to study the average height of a tree with n nodes, assuming that all n -node trees are equally likely. The corresponding problem for oriented (i.e., rooted, unordered) trees has been solved by Rényi and Szekeres [6]. Our principal results are stated in equations (32) and (34) below.

Trees appear in many disguises, and in particular there is a natural correspondence between trees of height $\leq h$ and discrete random walks in a straight line, with absorbing barriers at 0 and $h+1$. If we "wander around" a tree with n nodes, as shown by the dotted lines in Figure 1, the vertical component of successive positions describes a path of length $2n-1$ from 1 to 0; for example, the path in Figure 1 is 1,2,3,2,1,2,3,2,3,4,3,4,3,4,3,2,1,2,3,2,3,2,1,0. (This is one way a gambler can lose \$1 before winning \$5.) This construction, suggested by T. E. Harris in 1952 [2], is clearly reversible.

The height of trees plays a similar role in the classical ballot problem: How many ways are there to arrange n ballots for candidate A and m for candidate B in such a way that the number of votes for A never lags behind the number for B, as the ballots are counted, but A is never more than h votes ahead? The answer is the number of trees with $n+1$ nodes and height $\leq h+1$, again by the construction indicated in Figure 1. The ballot sequence corresponding to that tree is AABBAABAAABABABBAABAB .



Figure 1. A tree as a random walk.

We shall begin our study of the asymptotic properties of height by reviewing some known results. Let A_{nh} be the number of trees with n nodes and height $\leq h$, and let

$$A_h(z) = \sum A_{nh} z^n \quad (1)$$

be the corresponding generating function. We obtain all trees with height $\leq h+1$ by taking a root node and attaching zero or more subtrees each of which has height $\leq h$; therefore

$$\begin{aligned} A_{h+1}(z) &= z(1 + A_h(z) + A_h(z)^2 + A_h(z)^3 + \dots) \\ &= z / (1 - A_h(z)) , \quad \text{for } h \geq 0 . \end{aligned} \quad (2)$$

Clearly $A_0(z) = 0$. This relation yields a simple recurrence for the numbers A_{nh} ,

$$A_{n,h+1} = A_{n-1,h+1} A_{1,h} + A_{n-2,h+1} A_{2,h} + \dots + A_{1,h+1} A_{n-1,h} , \quad \text{for } n \geq 2, h \geq 0 . \quad (3)$$

From which it is easy to prepare a table of the first few values:

| $n = 1$ | 2 | 3 | 4 | 5 | 6 | 7 | 8 |
|---------|-----|-----|-----|-----|-----|-----|-----|
| $h = 1$ | 1 | 0 | 0 | 0 | 0 | 0 | 0 |
| $h = 2$ | 1 | 1 | 1 | 1 | 1 | 1 | 1 |
| $h = 3$ | 1 | 1 | 2 | 4 | 8 | 16 | 32 |
| $h = 4$ | 1 | 1 | 2 | 5 | 15 | 34 | 89 |
| $h = 5$ | 1 | 1 | 2 | 5 | 14 | 41 | 122 |
| $h = 6$ | 1 | 1 | 2 | 5 | 14 | 42 | 131 |

Since no tree with n nodes can have a height greater than n , we have

$$A_{nh} = A_{nh} - \binom{2n-2}{n-1} \frac{1}{n} , \quad h \geq n , \quad (4)$$

the well-known formula for the total number of trees with n nodes [cf. 3, p. 389].

Iteration of (2) yields a continued fraction representation of $A_h(z)$, e.g.

$$A_h(z) = \cfrac{z}{1 - \cfrac{z}{1 + \cfrac{1}{1 - z}}} . \quad (5)$$

This suggests expressing the generating function as a quotient of polynomials,

$$A_h(z) = \frac{p_h(z)}{p_{h+1}(z)} , \quad (6)$$

where

$$p_0(z) = 0 , \quad p_1(z) = 1 , \quad p_{h+1}(z) = p_h(z) - sp_{h-1}(z) . \quad (7)$$

The solution to this recurrence is

$$p_h(z) = \frac{1}{\sqrt{1-4z}} \left(\left(\frac{1+\sqrt{1-4z}}{2} \right)^h - \left(\frac{1-\sqrt{1-4z}}{2} \right)^h \right) , \quad (8)$$

and the form of this solution suggests setting $z = 1/(b \cos^2 \theta)$. We obtain

$$\begin{aligned} P_h\left(\frac{1}{b \cos^2 \theta}\right) &= \frac{\sin \theta}{\sin \theta (2 \cos \theta)^{b-1}} \\ A_h\left(\frac{1}{b \cos^2 \theta}\right) &= \frac{\sin \theta}{2 \cos \theta \sin(b+1)\theta} \end{aligned} \quad (9)$$

Incidentally it is easy to verify that $P_h(-1)$ is the Fibonacci number F_h , and that

$$P_h(z) = \sum_{0 \leq k \leq h} \binom{h-1-k}{k} (-z)^k, \quad \text{for } h \geq 1; \quad (10)$$

this leads to another recurrence for the A_{nh} .

Since $P_h(z)^2 = P_{h+1}(z)P_{h-1}(z) = z^{h-1}$, there is a simple generating function for the number of trees with n nodes and height exactly h ,

$$A_h(z) = A_{h-1}(z) = \frac{z^h}{P_{h+1}(z)P_h(z)}; \quad (11)$$

This formula was recently derived by Kremeras [5, p. 37].

Since P_h is a polynomial of degree $\lfloor (h-1)/2 \rfloor$, the roots of $P_h(z) = 0$ are $1/(b \cos^2(j\pi/b))$, for $1 \leq j < h/2$; and we obtain a partial fraction expansion of the generating function,

$$A_h(z) = \sum_{1 \leq j \leq h/2} \frac{\tan^2 \theta_{jh}}{(b+1)(1 - (b \cos^2 \theta_{jh})z)} + a_h + b_h z, \quad (12)$$

where

$$\begin{aligned} \theta_{jh} &= j\pi/(b+1); \\ a_{2m} &= 0; \quad b_{2m} = 0; \quad a_{2m+1} = \frac{-\pi(2m+1)}{b(b+1)}; \quad b_{2m+1} = \frac{1}{b^2+1}, \quad \text{for } m \geq 1. \end{aligned} \quad (13)$$

This leads immediately to the "explicit" formula

$$A_{nh} = \frac{1}{b+1} \sum_{1 \leq j \leq h/2} b^n \sin^2 \frac{j\pi}{b+1} \cos^{2n-2} \frac{j\pi}{b+1}, \quad \text{for } n \geq 2. \quad (14)$$

(It is rather remarkable that this formula gives a constant value for fixed n and all $b \geq n$. It is perhaps even more remarkable that Lagrange derived a formula in 1775 which essentially includes this as a special case! See [5, p. 247]; Feller [1, p. 322] observes that the formula has been rediscovered many times, although it appears in many texts on probability in connection with the equivalent "gambler's ruin" problem.) As a special case of (14) we have the asymptotic formula

$$A_{nh} \sim \frac{b^n}{b+1} \tan^2\left(\frac{\pi}{b+1}\right) \cos^{2n}\left(\frac{\pi}{b+1}\right), \quad \text{fixed } b, \quad n \rightarrow \infty. \quad (15)$$

Another interesting expression for A_{nh} can be derived by applying complex variable theory. We have

$$\begin{aligned} A_{nh} &= \frac{1}{2\pi i} \int_{\Gamma}^{(0+)} \frac{dz}{z^{n+1}} A_n(z) \\ &= \frac{1}{2\pi i} \int_{\Gamma}^{(0+)} \frac{dz}{z^n} (1-u) \frac{1-u^{-h}}{1-u^{h+1}} \end{aligned} \quad (16)$$

where

$$u = \frac{1 - \sqrt{1 - bz}}{1 + \sqrt{1 - bz}} \quad , \quad (17)$$

by (6) and (8). Since

$$z = \frac{u}{(1+u)^2} \quad , \quad (18)$$

we have $u \approx z$ when $|z| \ll 1$; hence we may change variables in (16) to obtain

$$A_{nh} = \frac{1}{2\pi i} \int_{\Gamma}^{(0+)} \frac{du}{u^n} (1-u)(1+u)^{2n-2} \frac{1-u^{-h}}{1-u^{h+1}} \quad . \quad (19)$$

In other words A_{nh} is the coefficient of u^{n-1} in $(1-u)(1+u)^{2n-2}(1-u^{-h}) / (1-u^{h+1})$. Some simplification now occurs when we consider the number of trees with height greater than h :

$$\begin{aligned} B_{nh} &= A_{nh} - A_{nh} \\ &= \frac{1}{2\pi i} \int_{\Gamma}^{(0+)} \frac{du}{u^{n+1}} (1-u)^2 (1+u)^{2n-2} \frac{u^{-h+1}}{1-u^{h+1}} \quad . \end{aligned} \quad (20)$$

It follows that

$$B_{n+1, h+1} = \sum_{k \geq 1} \left(\binom{2n}{n+1-k} + \binom{2n}{n-k} + \binom{2n}{n-1-k} \right) \quad . \quad (21)$$

The average height of a tree with n nodes is S_n/A_{nn} , where S_n is the (finite) sum

$$\begin{aligned} S_n &= \sum_{h \geq 1} h(A_{nh} - A_{n,h-1}) = \sum_{h \geq 1} h(B_{n,h-1} - B_{nh}) \\ &= \sum_{h \geq 0} B_{nh} \\ &= \frac{1}{2\pi i} \int_{\Gamma}^{(0+)} \frac{du}{u^{n+1}} (1-u)^2 (1+u)^{2n-2} \sum_{h \geq 1} \frac{u^{-h}}{1-u^h} \\ &= \frac{1}{2\pi i} \int_{\Gamma}^{(0+)} \frac{du}{u^{n+1}} (1-u)^2 (1+u)^{2n-2} \sum_{k \geq 1} d(k) u^k \quad . \end{aligned} \quad (22)$$

(As usual, $d(k)$ denotes the number of positive divisors of k .) Therefore

$$s_{n+1} = \sum_{k \geq 1} d(k) \left(\binom{2n}{n+k} - 2 \binom{2n}{n-k} + \binom{2n}{n-1-k} \right) \quad (25)$$

We shall now proceed to obtain an asymptotic series for the sum

$$f_n(n) = \sum_{k \geq 1} \frac{\binom{2n}{n+k}}{\binom{2n}{n}} d(k) , \quad \text{fixed } a, \quad n \rightarrow \infty \quad (26)$$

and this will lead to an asymptotic series for s_n .

Let $x = (k-a)/n$. By Stirling's approximation we have

$$\frac{\binom{2n}{n+k}}{\binom{2n}{n}} = \exp \left(-2n \left(\frac{x^2}{1/2} + \frac{x^4}{3/4} + \dots \right) + \left(\frac{x^2}{2} + \frac{x^4}{4} + \dots \right) - \frac{1}{x} (x^2 + x^4 + \dots) + O(x^2 n^{-3}) \right), \quad (25)$$

when $-1/2 < x < 1/2$, and

$$\frac{\binom{2n}{n+k}}{\binom{2n}{n}} = O(n^{-2\epsilon}) \quad \text{when } k \geq n^{1/2+\epsilon} + a ,$$

for all fixed $\epsilon > 0$. Therefore the sum of all terms for $k \geq n^{1/2+\epsilon} + a$ in (26) is negligible, being $O(n^{-\infty})$ for all $n > 0$, and we may take $x = O(n^{-1/2+\epsilon})$ in (25).

We now turn to the asymptotic behavior of the function

$$g_b(n) = \sum_{k \geq 1} k^b d(k) e^{-k^2/n} , \quad \text{fixed } b, \quad n \rightarrow \infty . \quad (26)$$

Again the terms for $k \geq n^{1/2+\epsilon}$ are negligible, so we can use (25) to express f_n in terms of g_b :

$$\begin{aligned} f_n(n) &= g_0(n) + \frac{2n}{n} g_1(n) - \frac{a^2}{n} g_0(n) + \frac{4a^2 + 1}{2n^2} g_2(n) - \frac{1}{6n^3} g_4(n) \\ &\quad - \frac{2a^3 + a}{n^2} g_1(n) + \frac{4a^3 + 5a}{3n^3} g_3(n) - \frac{a}{3n} g_5(n) + O(n^{-2\epsilon} g_0(n)) . \end{aligned} \quad (27)$$

In principle such an expansion could be carried out as far as we like, hence the problem of obtaining an asymptotic expansion for $f_n(n)$ reduces to the analogous problem for $g_b(n)$.

The behavior of $g_b(n)$ can be derived by starting with the well-known formula

$$e^{-x} = \frac{1}{2\pi i} \int_{c-i\infty}^{c+i\infty} \Gamma(z) x^{-z} dz , \quad c > 0, \quad x > 1 , \quad (28)$$

obtained, for example, by Fourier inversion of $\Gamma(c+iz)$. Then since $\Gamma(z) = \sum_{k \geq 1} d(k) / k^z$ we find

$$\begin{aligned}
a_b(n) &= \sum_{k \geq 1} \frac{1}{2\pi i} \int_{c-i\infty}^{c+i\infty} n^z \Gamma(z) z^{b-2k} d(z) dz \\
&= \frac{1}{2\pi i} \int_{c-i\infty}^{c+i\infty} n^z \Gamma(z) \zeta((2z-b)^2) dz
\end{aligned} \tag{29}$$

where now $c > \frac{1}{2}(b+1)$. Let q be a fixed positive number. When $\operatorname{Re}(z) \geq -q$, $\zeta(z) = O(|z|^{q+\frac{1}{2}})$ as $z \rightarrow -$; and since $n^z \Gamma(z)$ gets small on vertical lines we can shift the line of integration to the left as far as we please if we only take the residues into account. There is a double pole at $z = \frac{1}{2}(b+1)$, and possibly some simple poles at $z = 0, -1, -2, \dots$. Let $w = z - \frac{1}{2}(b+1)$; we have

$$n^z \Gamma(z) \zeta((2z-b)^2) = n^{\frac{1}{2}(b+1)} \Gamma(\frac{1}{2}(b+1)) (1 + w(\ln n) + O(w^2)) \cdot (1 + w(\frac{1}{2}(b+1)) + O(w^2)) \cdot (\frac{1}{4w^2} + \frac{1}{w} + O(1))$$

where $\psi(z) = \Gamma'(z)/\Gamma(z)$, hence the residue at the double pole is

$$n^{\frac{1}{2}(b+1)} \Gamma(\frac{1}{2}(b+1)) (\frac{1}{4} \ln n + \frac{1}{4} \psi(\frac{1}{2}(b+1)) + \gamma) \quad . \tag{30}$$

The residue at $z = -k$ is

$$\frac{n^{-k} (-1)^k \zeta(-2k-b)^2}{k!} = \frac{n^{-k} (-1)^k \pi^2}{(2k+b+1)^2 k!} \tag{31}$$

which is almost always zero when b is even. The sum of (30) and (31) for all $k \geq 0$ gives an asymptotic series for $a_b(n)$. Hence we have, for all $n > 0$,

$$\begin{aligned}
a_0(n) &= \frac{1}{4} \sqrt{\pi n} \ln n + (\frac{1}{4} \gamma - \frac{1}{2} \ln 2) \sqrt{\pi n} + \frac{1}{4} + O(n^{-1}) \quad ; \\
a_1(n) &= \frac{1}{8} n \ln n + \frac{1}{8} \gamma n + \frac{1}{32n} - \frac{1}{1600} n^{-1} + O(n^{-2}) \quad ; \\
a_2(n) &= \frac{5}{8} \sqrt{\pi n} \ln n + (\frac{1}{8} \gamma + \frac{1}{8} \ln 2) \sqrt{\pi n} + O(n^{-1}) \quad ;
\end{aligned} \tag{32}$$

etc. These formulas have been verified by computer calculation; for example, when $n = 10$, $a_0(n) \approx 3.96042$ and $\frac{1}{4} \sqrt{\pi n} \ln n + (\frac{1}{4} \gamma - \frac{1}{2} \ln 2) \sqrt{\pi n} + \frac{1}{4} \approx 3.96041$.

Returning to our original problem about trees, we have

$$\begin{aligned}
\frac{s_{n+1}}{(n+1) \prod_{m=1, m \neq 1}^n m} &= f_1(n) - 2f_0(n) + f_{-1}(n) \\
&= \frac{1}{n} a_0(n) + \frac{1}{n^2} a_2(n) + O(n^{-3/2} \log n)
\end{aligned} \tag{33}$$

by (4), (25), (24), and (27), and this equals $\sqrt{\pi} n^{-\frac{1}{2}} - \frac{1}{8} n^{-1} + O(n^{-3/2} \log n)$. We have proved, in particular, the following result:

Theorem. The average height of a planted plane tree with n nodes, considering all such trees to be equally likely, is

$$\sqrt{2n} - \frac{1}{2} + O(n^{-\frac{1}{2}} \log n) . \quad (34)$$

The same method can be used to obtain as many further terms of the expansion as desired. The factor $\log n$ in the error term turns out to be unnecessary.

We wish to thank Prof. John Riordan for pointing out references [2] and [4].

- [1] Feller, William, An Introduction to Probability Theory and its Applications, 1, 2nd edition (Wiley, 1957).
- [2] Harris, T. E., "First passage and recurrence distributions," Trans. Amer. Math. Soc. 73 (1952), 471-486.
- [3] Knuth, Donald E., The Art of Computer Programming, 1 (Addison-Wesley, 1968).
- [4] Kreweras, G., "Sur les éventails de segments," Cahiers du Bureau Universitaire de Recherche Opérationnelle 15 (1970), 1-41.
- [5] Lagrange, J. L., "Recherches sur les suites récurrentes," Oeuvres de Lagrange 1 (Paris, 1867), 149-251.
- [6] Rényi, A., and Szekeres, G., "On the height of trees," Australian J. Math. 7 (1967), 477-507.
- [7] Riordan, John, "The enumeration of trees by height and diameter," IBM J. Res. Develop. 4 (1960), 475-478.
- [8] Riordan, John, "Ballots and Trees," J. Comb. Theory 6 (1969), 408-411.