

ON THE VARIANCE OF THE HEIGHT OF RANDOM BINARY SEARCH TREES*

LUC DEVROYE[†] AND BRUCE REED[†]

Abstract. Let H_n be the height of a random binary search tree on n nodes. We show that there exists a constant $\alpha = 4.31107\dots$ such that $\mathbf{P}\{|H_n - \alpha \log n| > \beta \log \log n\} \rightarrow 0$, where $\beta > 15\alpha/\ln 2 = 93.2933\dots$. The proof uses the second moment method and does not rely on properties of branching processes. We also show that $\text{Var}\{H_n\} = O((\log \log n)^2)$.

Key words. binary search tree, probabilistic analysis, random tree, asymptotics, height, second moment method

AMS subject classifications. 68Q25, 60C05

1. Introduction. The height H_n of a random binary search tree on n nodes, constructed in the usual manner, starting from a random equiprobable permutation of $1, \dots, n$, is known to be close to $\alpha \log n$, where $\alpha = 4.31107\dots$ is the unique solution on $[2, \infty)$ of the equation $\alpha \log((2e)/\alpha) = 1$. First, Pittel [12] showed that $H_n/\log n \rightarrow \gamma$ almost surely as $n \rightarrow \infty$ for some positive constant γ . This constant was known not to exceed α (Robson [15]), and it was shown in Devroye [4] that $\gamma = \alpha$ as a consequence of the fact that $\mathbf{E}H_n \sim \alpha \log n$. Robson [16] has found that H_n does not vary much from experiment to experiment and seems to have a fixed range of width not depending on n . Devroye [5] proved that $H_n - \alpha \log n = O(\sqrt{\log n \log \log n})$ in probability, but this does not quite confirm Robson's findings. It is the purpose of this paper to prove the following theorem.

THEOREM.

$$\mathbf{E}H_n = \alpha \log n + O(\log \log n)$$

and

$$\text{Var}\{H_n\} = O((\log \log n)^2).$$

While this is a major step forward, we still do not know whether $\text{Var}\{H_n\} = O(1)$. For more information on random binary search trees, one may consult Knuth [7], [8], Aho, Hopcroft, and Ullman [1], [2], Mahmoud and Pittel [10], Devroye [6], Mahmoud [9], and Pittel [13].

Finally, we note that this paper contains the first proof of the asymptotic properties of H_n that is not based upon the theory of branching processes or branching random walks. We merely employ a well-known representation of random binary search trees from Devroye [4], and combine it with the second moment method, which has found so many other applications in the theory of random graphs (see, e.g., Palmer [11]).

2. Notation and definitions. Let T_∞ be the complete infinite binary tree. Each node x has a right son $r(x)$ and a left son $l(x)$. We consider a random labelled tree R_∞ obtained from T_∞ by choosing a uniform $[0, 1]$ random variable $U(x)$ for each node x of T_∞ and labelling the edge $(x, r(x))$ by $U(x)$ and the edge $(x, l(x))$ by $1 - U(x)$. The label of edge a is denoted $L(a)$. We let R_k be the random tree consisting of the first k edge levels of R_∞ .

For each node y of R_∞ , we let $f(y)$ be the product of the labels of the edges on the unique path from the root to y . We remark that for each $x \in R_\infty$, $-\log U(x)$ is an exponential

*Received by the editors September 24, 1992; accepted for publication (in revised form) April 21, 1994. This research was supported by Natural Sciences and Engineering Research Council of Canada grant A3456.

[†]School of Computer Science, McGill University, Montreal, Quebec H3A 2K6, Canada (luc@crodo.cs.mcgill.ca).

random variable with mean 1. If the labels on the path from the root to a node y of R_∞ are U_1, \dots, U_i , then we define

$$h_n(y) = \lfloor \dots \lfloor \lfloor nU_1 \rfloor U_2 \rfloor \dots U_i \rfloor.$$

Also, $-\log f(y)$ is distributed as the sum of i independently and identically distributed (i.i.d.) exponential random variables with mean 1, i.e., it is gamma distributed with parameter i .

Fact 1. It is well known that we can construct a random binary search tree T_n on n nodes by taking a copy R of R_∞ and letting T_n consist of those nodes y of R with $h_n(y) \geq 1$. (See, e.g., Devroye [4].)

Fact 2. Let y be a node of R_∞ at depth i (i.e., at edge-distance i from the root). Then

$$nf(y) - i \leq h_n(y) \leq nf(y).$$

Facts 1 and 2 basically allow us to obtain refined information regarding H_n merely by studying R_∞ . The inequality in Fact 2 introduces a certain looseness; in fact, it will limit the accuracy of the results on H_n to be $O(\log \log n)$.

3. Lemmas regarding the gamma distribution. The sum S_n of n i.i.d. exponential random variables with mean 1 is gamma (n) distributed. Its density is given by

$$g(t) = \frac{t^{n-1} e^{-t}}{(n-1)!}, \quad t > 0.$$

LEMMA 1. Let $\{t_n\}$ be a sequence of numbers such that $t_n \sim cn$ as $n \rightarrow \infty$ for some $c \in (0, 1)$. Then

$$\mathbf{P}\{S_n < t_n\} \sim \frac{1}{1-c} \frac{e^{-t_n} (t_n)^n}{n!}.$$

Proof. By integration by parts,

$$\begin{aligned} \mathbf{P}\{S_n < t_n\} &= \int_0^{t_n} \frac{t^{n-1} e^{-t}}{(n-1)!} dt \\ &= e^{-t_n} \left(\frac{t_n^n}{n!} + \frac{t_n^{n+1}}{(n+1)!} + \frac{t_n^{n+2}}{(n+2)!} + \dots \right) \\ &\sim \frac{1}{1-c} \frac{e^{-t_n} (t_n)^n}{n!}. \quad \square \end{aligned}$$

LEMMA 2. Let $t \in (0, 1)$ be a fixed constant. Then

$$\frac{e^{-tn} (tn)^n}{n!} \leq \mathbf{P}\{S_n < tn\} \leq \frac{1}{1-t} \frac{e^{-tn} (tn)^n}{n!}.$$

Proof. The lower bound follows directly by integration by parts as in the proof of Lemma 1. For the upper bound, note that

$$\begin{aligned} \mathbf{P}\{S_n < tn\} &\leq e^{-tn} \left(\frac{(tn)^n}{n!} + \frac{(tn)^{n+1}}{(n+1)!} + \frac{(tn)^{n+2}}{(n+2)!} + \dots \right) \\ &\leq \frac{e^{-tn} (tn)^n}{n!} \left(1 + \frac{tn}{n+1} + \left(\frac{tn}{n+1} \right)^2 + \dots \right) \\ &\leq \frac{e^{-tn} (tn)^n}{n!} \left(\frac{1}{1-t} \right). \quad \square \end{aligned}$$

LEMMA 3.

$$A \leq \sqrt{n}2^n \mathbf{P} \{S_n < n/\alpha\} \leq B,$$

where $A = e^{-1/12}/\sqrt{2\pi}$ and $B = \alpha/((\alpha - 1)\sqrt{2\pi})$.

Proof. From Lemma 2,

$$\frac{e^{-n/\alpha}(n/\alpha)^n}{n!} \leq \mathbf{P} \{S_n < n/\alpha\} \leq \frac{1}{1 - 1/\alpha} \frac{e^{-n/\alpha}(n/\alpha)^n}{n!}.$$

Use the fact that $n! = (n/e)^n \sqrt{2\pi n} e^{\theta/(12n)}$ for some $\theta \in (0, 1)$ and the definition of α . □

LEMMA 4. *There exists a universal constant C such that*

$$\mathbf{P} \{S_n \geq Cn\} \leq 2^{-2n}.$$

$C = 5$ will do.

Proof. Take $C > 1$. By Chernoff's exponential bounding method (Chernoff [3]), for $t > 0$,

$$\mathbf{P} \{S_n \geq Cn\} \leq \mathbf{E} e^{tS_n} e^{-tCn} = (1 - t)^{-n} e^{-tCn} = (Ce^{1-C})^n,$$

where we take $1 - t = 1/C$. For C large enough (e.g., $C \geq 5$), this is less than 4^{-n} . □

LEMMA 5. *Let E_1, E_2, \dots, E_n be i.i.d. random variables with a density, and let a be a fixed constant. Then*

$$\mathbf{P} \{E_1 < a, E_1 + E_2 < 2a, \dots, E_1 + \dots + E_n < na \mid E_1 + \dots + E_n < na\} \geq \frac{1}{n}.$$

Proof. Define $F_i = E_i - a$ for all i . Define $E_r = E_{r-n}$, when $n < r \leq 2n$. Then, by symmetry,

$$\begin{aligned} & \mathbf{P} \{E_1 < a, E_1 + E_2 < 2a, \dots, E_1 + \dots + E_n < na \mid E_1 + \dots + E_n < na\} \\ &= \mathbf{P} \{F_1 < 0, F_1 + F_2 < 0, \dots, F_1 + \dots + F_n < 0 \mid F_1 + \dots + F_n < 0\} \\ &= \frac{1}{n} \sum_{i=1}^n \mathbf{P} \{F_i < 0, F_i + F_{i+1} < 0, \dots, F_i + \dots + F_{i+n-1} < 0 \mid F_1 + \dots + F_n < 0\} \\ &= \mathbf{P} \{F_S < 0, F_S + F_{S+1} < 0, \dots, F_S + \dots + F_{S+n-1} < 0 \mid F_1 + \dots + F_n < 0\}, \end{aligned}$$

where S is independent of the E_i 's and uniformly distributed on $\{1, \dots, n\}$. Now, fix E_1, \dots, E_n , and let $s \in \{1, \dots, n\}$ be the (unique) value at which $\sum_{i>0, i<s} F_i$ is maximal. If $s = 1$, then $\sum_{i=1}^j F_i < 0$ for all $j > 0$. If $s > 1$, then, as $\sum_{i=1}^n F_i < 0$, we see that $\sum_{i=s}^{s+j} F_i = \sum_{i=1}^{s+j} F_i - \sum_{i=1}^{s-1} F_i < 0$ for all $j \geq 0$. Thus,

$$\begin{aligned} & \mathbf{P} \{F_S < 0, F_S + F_{S+1} < 0, \dots, F_S + \dots + F_{S+n-1} < 0 \mid F_1 + \dots + F_n < 0\} \\ & \geq \mathbf{P} \{S = s\} = \frac{1}{n}. \quad \square \end{aligned}$$

4. Proof of the theorem.

LEMMA 6. *Consider positive integers $n > k$. Then*

$$\mathbf{P} \{H_n \geq k\} \geq \mathbf{P} \{\exists \text{ leaf } y \in R_k : f(y) \geq (k + 1)/n\}.$$

Proof. This follows immediately from Facts 1 and 2. □

LEMMA 7. *There exists a constant $d > 0$ such that for sufficiently large j ,*

$$\mathbf{P} \{\exists \text{ leaf } y \in R_j : f(y) \geq (j + 1)/\exp(j/\alpha + d \log(j/\alpha))\} \geq 1 - \frac{1}{j^3}.$$

We may pick $d = \epsilon + 15/\log 2$ for any small $\epsilon > 0$.

Proof. The proof is contained in §5. \square

LEMMA 8. Let d be the constant of Lemma 7. Then, for sufficiently large n ,

$$\mathbf{P} \{H_n \geq \alpha \log n - d\alpha \log \log n - 1\} \geq 1 - \frac{1}{(\alpha \log n)^3}.$$

We may choose $d = \epsilon + 15/\log 2$ for any small $\epsilon > 0$.

Proof. The proof follows from Lemmas 6 and 7 by setting $j = k = \lfloor \alpha \log n - d\alpha \log \log n \rfloor$. \square

LEMMA 9.

$$\mathbf{P} \{H_n \geq \lceil \alpha \log n + i \rceil\} \leq \left(\frac{2}{\alpha}\right)^i, \quad i \geq 0.$$

Proof. See Devroye [4, p. 492]. \square

Note that the theorem follows from Lemmas 8 and 9 without work.

5. Proof of Lemma 7.

LEMMA 10. For every i with probability at least $1 - 2^{-i}$, every leaf of R_i has $f(y) \geq e^{-5i}$.

Proof. The probability that, for some leaf y of R_i , we have $f(y) < e^{-5i}$ is at most 2^i times $\mathbf{P}\{S_i \geq 5i\}$, where S_i is gamma i distributed. By Lemma 4, this does not exceed $2^i/4^i = 2^{-i}$. \square

LEMMA 11. For sufficiently large k with probability at least $1/k^3$, there is a leaf y of R_k with $f(y) \geq e^{-k/\alpha}$.

Lemma 11 will be proved in §6. If Lemma 11 is true, then we can proceed with the proof of Lemma 7 as follows: First note that we can obtain a copy of R_{i+k} by making each leaf of R_i a root of a copy of R_k , where all these trees are independently labelled. Define $k = \lfloor j - A \log j \rfloor$ and $i = \lceil A \log j \rceil$ so that $j = k + i$ with some constant A to be picked further on. Note first that for j large enough, if $A > \alpha$,

$$\frac{k}{\alpha} + 5i \leq \frac{j}{\alpha} + 5A \log \left(\frac{j}{\alpha}\right) - \log(j + 1).$$

Then,

$$\begin{aligned} & \mathbf{P} \{ \nexists \text{ leaf } y \in R_j \text{ with } f(y) \geq 1/\exp(j/\alpha + 5A \log(j/\alpha) - \log(j + 1)) \} \\ & \leq \mathbf{P} \{ \exists \text{ leaf } y \in R_i \text{ with } f(y) < e^{-5i} \} \\ & \quad + \mathbf{P} \{ \exists \text{ leaf } y \in R_j \text{ with } f(y) \geq 1/\exp(j/\alpha + d \log(j/\alpha) - \log(j + 1)) \\ & \quad \quad | \forall \text{ leaf } y \in R_i : f(y) \geq e^{-5i} \} \\ & \leq 2^{-i} + \mathbf{P} \{ \text{every copy of } R_k \text{ contains no leaf } y \text{ with } f(y) \geq 1/\exp(k/\alpha) \} \\ & \quad \text{(by Lemma 10)} \\ & \leq 2^{-i} + (1 - k^{-3})^{2^i} \quad \text{(by Lemma 11)} \\ & \leq 2^{-i} + \exp(-2^i k^{-3}) \\ & \leq j^{-A \log 2} + \exp(-j^{A \log 2 - 3}) \\ & \leq j^{-3} \end{aligned}$$

for j large enough, provided that $A \log 2 > 3$. This proves Lemma 7. We note that we can pick $d = 5A$, where $A = \epsilon + \max(\alpha, 3/\log 2)$ for any small $\epsilon > 0$.

6. Proof of Lemma 11. Let P be a path from the root to a leaf y of R_k . The condition $f(y) \geq 1/e^{k/\alpha}$ is equivalent to

$$\sum_{e \in P} (-\log L(e)) \leq \frac{|P|}{\alpha}.$$

We call a leaf y special if, in addition to the above condition, it satisfies

$$\sum_{e \in P'} (-\log L(e)) \leq \frac{|P'|}{\alpha}$$

for every subpath P' of P that originates at a terminal vertex y . Such subpaths are called terminal. Let S be the collection of special leaves of R_k . By Lemma 5, the expected number of special leaves is at least $1/k$ times $\mathbf{P}\{S_k < k/\alpha\}$ times 2^k . By Lemma 3,

$$\mathbf{E}|S| \geq \frac{e^{-1/12}}{\sqrt{2\pi}k^{3/2}}.$$

Next, we consider the expected number of pairs of special leaves to be able to apply the second moment method. We fix a leaf z of R_k and count $|S|$, given that $z \in S$. To this end, let w be another leaf of R_k . Let P_z and P_w denote the paths from the root of R_k to z and w , respectively. Then P_z and P_w have an initial common subsequence, i.e., the join $P_z \cap P_w$. Let e_1, e_2, \dots, e_k be the edges on the path from the root to z and define $Q_i = \{e_1, \dots, e_i\}$. For any j , the number of leaves of R_k whose join with P_z is Q_j is 2^{k-j} . Furthermore, the probability that a leaf $w \in R_k$ is a special leaf, given that $z \in S$ and $P_z \cap P_w = Q_j$, is bounded above by the probability that for the terminal path $P' \subseteq P_w - Q_j$ with $|P'| = \max(0, k - j - 1)$, we have

$$\sum_{e \in P'} (-\log L(e)) \leq \frac{|P'|}{\alpha}.$$

Note that P' contains one edge less than $P_w - Q_j$. Later, this allows us to work out a conditional probability, given $z \in S$, without much trouble. By Lemma 3, the probability of the event mentioned above is at most

$$\frac{\alpha}{(\alpha - 1)\sqrt{2\pi}(k - j - 1)2^{k-j-1}}.$$

Thus,

$$\begin{aligned} \mathbf{E} \{ |\{w \in S : P_w \cap P_z = Q_j\}| \mid z \in S \} & \\ & \leq \frac{\alpha 2^{k-j}}{(\alpha - 1)\sqrt{2\pi}(k - j - 1)2^{k-j-1}} \\ & = \frac{2\alpha}{(\alpha - 1)\sqrt{2\pi}(k - j - 1)} \\ & \leq 2, \end{aligned}$$

when $k - j \geq 2$. The previous expected value is bounded by 2 when $k - j \in \{0, 1\}$. Therefore,

$$\begin{aligned} \mathbf{E} \{|\mathcal{S}| \mid z \in \mathcal{S}\} &= \sum_{j=0}^k \mathbf{E} \{|\{w \in \mathcal{S} : P_w \cap P_z = Q_j\}| \mid z \in \mathcal{S}\} \\ &\leq \sum_{j=0}^k 2 = 2k + 2. \end{aligned}$$

Hence, by the second moment method,

$$\begin{aligned} \mathbf{P} \{|\mathcal{S}| \geq 1\} &\geq \frac{\mathbf{E}|\mathcal{S}|}{1 + \sup_{z \text{ leaf of } R_k} \mathbf{E}\{|\mathcal{S}| \mid z \in \mathcal{S}\}} \\ &\geq \frac{\mathbf{E}|\mathcal{S}|}{2k + 3} \\ &\geq \frac{e^{-1/12}}{\sqrt{2\pi}(2k + 3)k^{3/2}} \\ &\geq \frac{1}{k^3} \end{aligned}$$

for all k large enough. This concludes the proof of Lemma 11.

Acknowledgments. The authors thank Colin McDiarmid and an anonymous referee for helpful comments.

REFERENCES

- [1] A. V. AHO, J. E. HOPCROFT, AND J. D. ULLMAN, *The Design and Analysis of Computer Algorithms*, Addison-Wesley, Reading, MA, 1975.
- [2] ———, *Data Structures and Algorithms*, Addison-Wesley, Reading, MA, 1983.
- [3] H. CHERNOFF, *A measure of asymptotic efficiency for tests of a hypothesis based on the sum of observations*, Ann. Math. Statist., 23 (1952), pp. 493–507.
- [4] L. DEVROYE, *A note on the height of binary search trees*, J. Assoc. Comput. Mach., 33 (1986), pp. 489–498.
- [5] ———, *Branching processes in the analysis of the heights of trees*, Acta Inform., 24 (1987), pp. 277–298.
- [6] ———, *On the height of random m -ary search trees*, Random Structures Algorithms, 1 (1990), pp. 191–203.
- [7] D. E. KNUTH, *The Art of Computer Programming, Vol. 1: Fundamental Algorithms*, Addison-Wesley, Reading, MA, 1973.
- [8] ———, *The Art of Computer Programming, Vol. 3: Sorting and Searching*, Addison-Wesley, Reading, MA, 1973.
- [9] H. M. MAHMOUD, *Evolution of Random Search Trees*, John Wiley, New York, 1992.
- [10] H. MAHMOUD AND B. PITTEL, *On the most probable shape of a search tree grown from a random permutation*, SIAM J. Algebraic Discrete Meth., 5 (1984), pp. 69–81.
- [11] E. M. PALMER, *Graphical Evolution*, John Wiley, New York, 1985.
- [12] B. PITTEL, *On growing random binary trees*, J. Math. Anal. Appl., 103 (1984), pp. 461–480.
- [13] ———, *Note on the heights of random recursive trees and random m -ary search trees*, Tech. report, Department of Mathematics, Ohio State University, 1992.
- [14] R. PYKE, *Spacings*, Roy. Statist. Soc. Ser. B, 7 (1965), pp. 395–445.
- [15] J. M. ROBSON, *The height of binary search trees*, Austral. Comput. J., 11 (1979), pp. 151–153.
- [16] ———, *The asymptotic behaviour of the height of binary search trees*, Austral. Comput. Sci. Comm., Queensland Univ. Tech., Brisbane, 1982, p. 88.