

On the expected height of fringe-balanced trees

Luc Devroye*

School of Computer Science, McGill University, 3480 University Street, Montreal, Canada H3A 2A7

Received October 26, 1992 / February 16, 1993

Abstract. We study the effect of a well-known balancing heuristic on the expected height of a random binary search tree. After insertion of an element, if any node on the insertion path has a subtree of size precisely $2t+1$ for a fixed integer t , then the subtree rooted at that node is destroyed and replaced by a new subtree in which the median of the $2t+1$ elements is the new root. If H_n denotes the height of the resulting random tree, we show that $H_n/\log n \rightarrow c(t)$ in probability for some function $c(t)$. In particular, $c(0)=4.31107\dots$ (the ordinary binary search tree), $c(1)=3.192570\dots$, $c(3)=2.555539\dots$, $c(10)=2.049289\dots$ and $c(100)=1.623695\dots$.

1 Introduction

Consider an ordinary binary search constructed by standard consecutive insertions of values X_1, \dots, X_n . It is well-known that when the input forms a random permutation of $\{1, \dots, n\}$ (or equivalently, when the input sequence is independent and identically distributed) that the height H_n of the tree satisfies the following convergence property:

$$\frac{H_n}{\log n} \rightarrow 4.31107 \dots \text{ a.s.}$$

(Robson 1979; Devroye 1986). The purpose of this note is to investigate what happens to H_n when we apply a very simple heuristic during the insertion process.

Bell (1965) and Walker and Wood (1976) introduced the following heuristic applied to the fringe of the tree: fix an integer t . After insertion of an element in the tree, verify whether one of the nodes on the insertion path is the root of a subtree of size precisely $2t+1$. If so, we perform the operation MEDIAN-ROOT on that subtree. This operation consists of finding the median of the

* Research of the author was sponsored by NSERC Grant A3456 and by FCAR Grant 90-ER-0291

$2t+1$ elements in the subtree, and making it the root. This can be done in a number of ways, via a splay operation, via a tree splitting method, or by complete reorganization of the subtree. As t is typically small – $t=1$ is the most frequently studied case –, the reorganization of the subtree is less important.

Another way of looking at this, following Poblete and Munro (1985), is to consider internal nodes and external nodes, where internal nodes hold one data point and external nodes are bags of capacity $2t$. Insertion proceeds as usual. As soon as an external node overflows (i.e., when it would grow to size $2t+1$), its bag is split about the median, leaving two new external nodes (bags) of size t each, and an internal node holding the median. After the insertion process is completed, we may wish to expand the bags into balanced trees. In fact, what we do with the bags does not matter: if H'_n is the maximal distance between an internal node and the root, and H_n is the height of the tree, i.e., the maximal distance between any expanded external node and the root, it is easy to see that

$$H'_n < H_n \leq H'_n + 2t.$$

Asymptotically, this is unimportant.

Using the branching process method of proof (Devroye 1986, 1987, 1990; see also Mahmoud 1992) we obtain the almost sure limit value of $H_n/\log n$ for all t . For another possible proof method, see Pittel (1992). The improvement in H_n is important for small values of t . Not surprisingly, for every $\varepsilon > 0$, we can find a t such that

$$\lim_{t \rightarrow \infty} \mathbf{P}\{H_n > (1 + \varepsilon) \log_2 n\} = 0.$$

Unfortunately, the value of t increases very rapidly as $\varepsilon \downarrow 0$.

2 The expected depth of a node

The depth D_n of the last node when the fringe heuristic is used has been studied by the theory of Markov processes or urn models in a series of papers, notably by Poblete and Munro (1985), Aldous et al. (1988). See also Gonnet and Baeza-Yates (1991, p. 109). Poblete and Munro (1985) showed that

$$\lim_{n \rightarrow \infty} \frac{\mathbf{E}D_n}{\log n} = d(t)$$

and

$$\frac{D_n}{\log n} \rightarrow d(t) \text{ i.p.,}$$

where

$$\frac{1}{d(t)} = \sum_{i=t+1}^{2t+1} \frac{1}{1+i}.$$

It is a simple exercise to show that $d(t) \rightarrow 1/\log 2$ as $t \uparrow \infty$. The values of H_n we will obtain below are larger than these. A comparison of the limits will be provided.

3 The main result

Theorem 1 *A random binary search tree constructed with the aid of the fringe heuristic with parameter t has the following property:*

$$\frac{H_n}{\log n} \rightarrow c(t) \text{ i.p.,}$$

where $c(t)$ is the unique solution greater than $d(t)$ of the equation

$$(1) \quad \lambda(c) - c \sum_{i=t+1}^{2t+1} \log \left(1 + \frac{\lambda(c)}{i} \right) + c \log 2 = 0,$$

and $\lambda(c)$ is defined by the equation

$$\frac{1}{c} = \sum_{i=t+1}^{2t+1} \frac{1}{\lambda + i}.$$

A table of values for $d(t)$ and $c(t)$ is given below.

t	$c(t)$	$d(t)$	t	$c(t)$	$d(t)$
0	4.311070	2.000000	10	2.049289	1.490455
1	3.192570	1.714286	20	1.863726	1.467601
2	2.779633	1.621622	30	1.782617	1.459539
3	2.555539	1.575985	40	1.734851	1.455420
4	2.411554	1.548863	50	1.702554	1.452920
5	2.309726	1.530900	60	1.678898	1.451241
6	2.233133	1.518130	70	1.660617	1.450035
7	2.172976	1.508587	80	1.645976	1.449128
8	2.124195	1.501186	90	1.633883	1.448420
9	2.083648	1.495279	100	1.623695	1.447853

Remark 1 *The existence of a solution.* The constant $\lambda(c)$ is well-defined whenever

$$\frac{1}{c} \leq \sum_{i=t+1}^{2t+1} \frac{1}{1+i} \stackrel{\text{def}}{=} \frac{1}{d}.$$

This follows from the fact that the left-hand-side of (1) is a concave function of c , and that at $c=d$ (note that $\lambda(c)=1$), the left-hand-side of (1) is positive (and takes the value 1). The derivative of (1) is

$$\begin{aligned} \lambda'(c) - \sum_{i=t+1}^{2t+1} \log \left(1 + \frac{\lambda(c)}{i} \right) - c \lambda'(c) \sum_{i=t+1}^{2t+1} \frac{1}{\lambda(c)+i} + \log 2 \\ = - \sum_{i=t+1}^{2t+1} \log \left(1 + \frac{\lambda(c)}{i} \right) + \log 2, \end{aligned}$$

and the second derivative of (1) is

$$-\frac{1}{\partial c/\partial \lambda} \sum_{i=t+1}^{2t+1} \frac{1}{\lambda(c)+i} = -\frac{1}{c \partial c/\partial \lambda} < 0.$$

Thus, the left-hand side of (1) is a concave function in c .

4 A property of spacings

As in Devroye (1986, 1987), we establish a crucial link between random binary search trees and trees of random variables. The root of a random binary search tree splits the $n-1$ remaining elements in two sets, one for each subtree, where the size of the left subtree is distributed as $\lfloor nU \rfloor$, and U is a uniform $[0, 1]$ random variable. In the model that we are considering, the split is not according to a uniform $[0, 1]$ random variable. Rather, the value U should be replaced by V , the median of $2t+1$ i.i.d. uniform $[0, 1]$ random variables. Such a random variable will be called a median-of-uniform (MOU) random variable, and we will denote all such random variables by V or V_i . The sizes of the subtrees are jointly distributed as $(N, n-1-N)$, where $nV-t-1 \leq N \leq nV+t$, and \leq denotes stochastic ordering.

Proof. Let U denote the minimum of $2t+1$ i.i.d. uniform $[0, 1]$ random variables, and let M be the minimum of $2t+1$ integers drawn without replacement from $\{1, \dots, n\}$. A simple argument shows that $M \leq nU+1$, where \leq denotes stochastic ordering. Thus, $n-M \geq nW-1$, where W is distributed as the maximum of $2t+1$ i.i.d. uniform $[0, 1]$ random variables. Iterating this $t+1$ times, an induction argument shows that $N \geq nV-(t+1)$, where V is distributed as the median of $2t+1$ i.i.d. uniform $[0, 1]$ random variables. Since $n-1-N$ is distributed as N , we see that $N \leq nV+t$ as well, where all the inequalities still refer to stochastic dominance.

Each of the subtrees can be split in a similar fashion, requiring this time two new MOU random variables. This process can be repeated at all levels and it leads to a tree in which node values are subtree sizes that are approximately obtained as (truncated) products of MOU random variables. More formally, let T_k be a complete binary tree with k full levels of edges. The total number of edges is $2^1+2^2+\dots+2^k=2^{k+1}-2$. We will use the symbol p for a path from root to leaf (there are 2^k such paths in T_k). Consider all edges pairwise in level order and from left to right, and associate with each pair an independent random vector distributed as $(V, 1-V)$ where V is distributed as a MOU random variable. Take one of the 2^k leaves of T_k and let p be the path from that leaf to the root of T_k . Let $\{V_1, V_2, \dots, V_k\}$ denote the collection of MOU random variables encountered on the path from leaf to root. The quantity

$$nV_1 2 V_2 \dots V_k$$

is approximately equal to the size of the subtree rooted at a leaf in the random binary search tree. A node corresponds to a real internal node in the search

tree when its values is $\geq 2t + 1$. Indexing the MOU random variables in an obvious manner, we define

$$Z_k = \begin{cases} 1 & (k=0) \\ \max_p \prod_{i \in p} V_i & (k>0). \end{cases}$$

Lemma 1 *Let $n \geq 1, k \geq 1$ be given integers. Then the internal node height H'_n of a random binary search tree on n nodes with fringe reorganization is related to the random variables Z_k via the following inequalities:*

$$\begin{aligned} \mathbf{P}\{H'_n \geq k\} &\geq \mathbf{P}\left\{Z_k \geq \frac{1+2t+k(t+1)}{n}\right\}, \\ \mathbf{P}\{H'_n \geq k+2l\} &\leq \mathbf{P}\left\{Z_k \geq \frac{1}{n}\right\} + 2^l \mathbf{P}\left\{\prod_{i=1}^l V_i \geq \frac{l}{k-l}\right\}, \quad 0 < l < k. \end{aligned}$$

Here V_1, \dots, V_k are i.i.d. MOU random variables.

Proof. A simple induction argument shows that the size of the subtree rooted at a node with path p to the root is at least equal to

$$n \prod_{i \in p} V_i - k(t+1).$$

If this is at least $2t + 1$ for some node, then $H'_n \geq k$. On the other hand, if $H'_n \geq k + 2l$, then some node at distance k from the root has a subtree of size at least equal to $2l(t+1) + 1$. On the other hand, if $N(p)$ denotes the size of the subtree rooted at the node at distance k from the root, where p is the path to the root, and V_1, \dots, V_k are the MOU variables associated with the edges on this path, with V_k nearest the root, then $N(p)$ is in distribution not greater than

$$\begin{aligned} &n \prod_{i=1}^k V_i + t + t \sum_{j=2}^k \prod_{i=j}^k V_i \\ &\leq n \prod_{i=1}^k V_i + lt + (k-l)t \prod_{i=k-l+1}^k V_i. \end{aligned}$$

The last product is shared by precisely 2^{k-l} paths. Therefore, by Bonferroni's inequality,

$$\begin{aligned} \mathbf{P}\{H'_n \geq k+2l\} &\leq \mathbf{P}\left\{\max_p N(p) \geq 2l(t+1) + 1\right\} \\ &\leq \mathbf{P}\left\{Z_k \geq \frac{1}{n}\right\} + 2^l \mathbf{P}\left\{\prod_{i=1}^l V_i \geq \frac{l}{k-l}\right\}. \quad \square \end{aligned}$$

Next, we need a simple result about the MOU distribution related to spacings of uniform random variables (see Pyke 1965).

Lemma 2 A MOU random variable V is distributed as

$$\prod_{i=t+1}^{2t+1} U_i^{1/i},$$

where the U_i 's are i.i.d. uniform $[0, 1]$.

Proof. Note that the maximum of $2t+1$ i.i.d. uniform $[0, 1]$ random variables is distributed as $U_{2t+1}^{1/(2t+1)}$. Apply this rule recursively to obtain the given property. \square

Finally, we recall a theorem regarding trees of random variables from the theory of branching random walks (Hammersley 1974; Kingman 1973; Biggins 1976, 1977): note that

$$\log Z_k = \begin{cases} 0 & (k=0) \\ \max_p \sum_{i \in p} \log V_i & (k>0). \end{cases}$$

The maximum is over all 2^k paths in the tree of length 2^k . Then, we have

$$\frac{\log Z_k}{k} \rightarrow \gamma \text{ a.s.}$$

as $k \rightarrow \infty$, where

$$\gamma = \inf \left\{ x > - \sum_{j=t+1}^{2t+1} \frac{1}{j} : 2M(x) < 1 \right\}$$

and

$$M(x) = \inf_{\lambda > 0} \mathbf{E} \{ V^\lambda e^{-\lambda x} \}.$$

In our case, by Lemma 2, we see that

$$M(x) = \inf_{\lambda > 0} \prod_{j=t+1}^{2t+1} \frac{1}{1 + \lambda/j} e^{-\lambda x}.$$

The minimum of the latter expression is obtained when λ is the solution of the equation

$$x = - \sum_{j=t+1}^{2t+1} \frac{1}{j + \lambda}.$$

This only has a positive solution when

$$0 > x > - \sum_{j=t+1}^{2t+1} \frac{1}{j}.$$

The function $M(x)$ is decreasing in x . Over the range of interest here, it decreases from 1 to 0. Therefore, γ is well-defined.

Another way of viewing this is from the exponential inequality angle. Indeed, for $\lambda > 0$,

$$\begin{aligned} \mathbf{P}\{\log Z_k > kx\} &\leq 2^k \mathbf{P}\left\{\log \prod_{i=1}^k V_i > kx\right\} \\ &= 2^k \mathbf{P}\left\{\sum_{i=1}^k \log V_i > kx\right\} \\ &\leq 2^k \mathbf{E}\left\{\prod_{i=1}^k e^{\lambda \log V_i} e^{-\lambda kx}\right\} \\ &= (2 \mathbf{E}\{V_1^\lambda e^{-\lambda x}\})^k \\ &= (2M(x))^k \end{aligned}$$

if we choose λ as outlined above. The upper bound tends to zero as $k \rightarrow \infty$ when $2M(x) < 1$.

5 Proof of Theorem 1

From Lemma 1 and the last result of the previous section, we see that

$$\mathbf{P}\left\{Z_k \geq \frac{1}{n}\right\} = \mathbf{P}\left\{\frac{\log Z_k}{k} \geq \frac{-\log n}{k}\right\} \rightarrow 0$$

when we choose $k \sim c \log n$ and $-1/c > \gamma$. Also, as each V_i (notation of Lemma 1) is stochastically smaller than $(1 + U)/2$, where U is uniform $[0, 1]$, we see that

$$\mathbf{E} V_i^\lambda \leq \mathbf{E}(1 + U)^\lambda / 2^\lambda \leq 2/(1 + \lambda).$$

Thus, for any $\lambda > 0$, by Jensen's inequality, and $0 < l < k$,

$$\begin{aligned} 2^l \mathbf{P}\left\{\prod_{i=1}^l V_i \geq \frac{l}{k-l}\right\} &\leq 2^l \frac{\mathbf{E}^l V_i^\lambda (k-l)^\lambda}{l^\lambda} \\ &\leq \left(\frac{4}{(l+\lambda)}\right)^l \left(\frac{k-l}{l}\right)^\lambda \rightarrow 0 \end{aligned}$$

if we choose $l = \lfloor \sqrt{k} \rfloor$, $k \sim c \log n$ (as above), and $\lambda = 4$. By Lemma 1, we thus conclude that $\mathbf{P}\{H'_n \geq k + 2\sqrt{k}\} \rightarrow 0$. Furthermore,

$$\mathbf{P}\{H'_n \geq k\} \geq \mathbf{P}\left\{\frac{\log Z_k}{k} \geq \frac{\log(2t+1+k(t+1)) - \log n}{k}\right\} \rightarrow 1$$

when we choose $k \sim c \log n$ and $-1/c < \gamma$. Putting both statements together, we obtain

$$\frac{H'_n}{\log n} \rightarrow -\frac{1}{\gamma} \text{ i.p.}$$

Set $-1/\gamma \equiv c(t)$. This leads precisely to the definition of λ and to equation (1) (which is equivalent to $2M(-1/c) = 1$) in Theorem 1. \square

As a by-product of the inequality

$$\mathbf{P}\{Z_k \geq 1/n\} \leq \left(2M\left(\frac{-\log n}{k}\right)\right)^k,$$

we see that $\mathbf{E}H'_n{}^q/\log^q n \rightarrow (c(t))^q$ for all $q > 0$. Using arguments as in Pittel (1984), we may also deduce strong convergence of $H'_n/\log n$ to $c(t)$.

References

- Aldous, D., Flannery, B., Palacios, J.L.: Two applications of urn processes: the fringe analysis of search trees and the simulation of quasi-stationary distributions of Markov chains. *Probab. Eng. Inf. Sci.* **2**, 293–307 (1988)
- Bell, C.J.: An investigation into the principles of the classification and analysis of data of an automatic digital computer. Doctoral Dissertation, Leeds University, 1965
- Biggins, J.D.: The first and last-birth problems for a multitype age-dependent branching process. *Adv. Appl. Probab.* **8**, 446–459 (1976)
- Biggins, J.D.: Chernoff's theorem in the branching random walk. *J. Appl. Probab.* **14**, 630–636 (1977)
- Devroye, L.: A note on the height of binary search trees. *J. ACM* **33**, 489–498 (1986)
- Devroye, L.: Branching processes in the analysis of the heights of trees. *Acta Inf.* **24**, 277–298 (1987)
- Devroye, L.: On the height of random m -ary search trees. *Random Struct. Algorithms* **1**, 191–203 (1990)
- Gonnet, G.H., Baeza-Yates, R.: Handbook of algorithms and data structures. Reading: Addison-Wesley 1991
- Hammersley, J.M.: Postulates for subadditive processes. *Ann. Probab.* **2**, 652–680 (1974)
- Kingman, J.F.C.: Subadditive ergodic theory. *Ann. Probab.* **1**, 883–909 (1973)
- Mahmoud, H.M.: Evolution of Random Search Trees. New York: John Wiley 1992
- Pittel, B.: On growing random binary trees. *J. Math. Anal. Appl.* **103**, 461–480 (1984)
- Pittel, B.: Note on the heights of random recursive trees and random m -ary search trees. Technical Report, Department of Mathematics, The Ohio State University 1992
- Poblete, P.V., Munro, J.I.: The analysis of a fringe heuristic for binary search trees. *J. Algorithms* **6**, 336–350 (1985)
- Pyke, R.: Spacings. *J. R. Stat. Soc. Ser. B* **7**, 395–445 (1965)
- Robson, J.M.: The height of binary search trees. *Aust. Comput. J.* **11**, 151–153 (1979)
- Walker, A., Wood, D.: Locally balanced binary trees. *Comput. J.* **19**, 322–325 (1976)