

# The height of q-Binary Search Trees

Michael Drmota, Helmut Prodinger

► **To cite this version:**

Michael Drmota, Helmut Prodinger. The height of q-Binary Search Trees. Discrete Mathematics and Theoretical Computer Science, DMTCS, 2002, 5, pp.97-108. <hal-00958975>

**HAL Id: hal-00958975**

**<https://hal.inria.fr/hal-00958975>**

Submitted on 13 Mar 2014

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

# The Height of $q$ -Binary Search Trees

Michael Drmota<sup>1</sup> and Helmut Prodinger<sup>2</sup>

<sup>1</sup> Department of Geometry, Technical University of Vienna, Wiedner Hauptstrasse 8–10, A-1040 Vienna, Austria  
e-mail: michael.drmota@tuwien.ac.at

<sup>2</sup> The John Knopfmacher Centre for Applicable Analysis and Number Theory, School of Mathematics, University of the Witwatersrand, P. O. Wits, 2050 Johannesburg, South Africa  
e-mail: helmut@staff.ms.wits.ac.za, WWW: <http://www.wits.ac.za/helmut/index.htm>

received Aug 10, 2001, accepted Apr 29, 2002.

---

$q$ -binary search trees are obtained from words, equipped with the geometric distribution instead of permutations. The average and variance of the height are computed, based on random words of length  $n$ , as well as a Gaussian limit law.

**Keywords:** Binary search tree,  $q$ -analog, height

---

## 1 Introduction

The paper [8] introduces for the first time a meaningful  $q$ -model of binary search trees: instead of binary search trees, one considers tournament trees, which differ only marginally from binary search trees; if one starts from a permutation  $(\pi_1 \pi_2 \dots \pi_n)$ , then one inserts the number  $i$  instead of the number  $\pi_i$ . Thus, traversing the tree in inorder, we might think of the associated permutation as  $\rho 1 \sigma$ , where 1 goes into the root, and  $\rho$  resp.  $\sigma$  form (recursively) the left resp. right subtree. We could have called this paper “The Height of  $q$ -Tournament Trees;” however we decided not to do so since binary search trees are by far better known, both, in the community of theoretical computer scientists, and combinatorialists. A nice reference for tournament trees and increasing trees in general is [2].

Now instead of considering permutations  $\pi_1 \pi_2 \dots \pi_n$ , we consider words over the alphabet  $\{1, 2, 3, \dots\}$ , and (geometric) probabilities attached to the letters, i. e., the probability of letter  $i$  is  $pq^{i-1}$ , with  $p + q = 1$ . The binary search tree is then constructed by writing a nonempty word  $w$  as  $w = xay$ , where  $a$  is the smallest letter occurring, and  $x \in \{a + 1, a + 2, \dots\}^*$  and  $y \in \{a, a + 1, \dots\}^*$ . The letter  $a$  goes into the root and  $x$  resp.  $y$  form the left resp. right subtree.

The paper [8] dealt with the path length; here we consider the height. The height of a binary search tree (and thus of the trees in our model) is defined to be the largest number of nodes in a path from the root to a leaf; the empty tree (related to the empty word) has height 0.

We will prove that the expected height, when considering random words of length  $n$ , is asymptotic to  $pn$ ; the variance will also be computed as well as a Gaussian limit law. (The letter  $p$  will always denote  $1 - q$  in this paper.) Recall that the result for traditional binary search trees is  $\sim c \log n$  with  $c = 4.31107$ , see [3]; unfortunately we do not get that result as the limit  $q \rightarrow 1$  as it so happened for the path length. However, nothing

is wrong here, since the limit  $q \rightarrow 1$  in  $pn$  simply tells us that the average height should be less than linear. Although this original hope of getting the classical result as a corollary did not work out, we nevertheless think that the results presented here are of independent interest. Note also that  $pn$  is the expected number of letters 1 in a (random) string of length  $n$ . In our asymmetric model, they must all lie on one path.

In more detail, we will obtain the following results:

Let  $H_n$  denote the height of  $q$ -binary search trees with  $n$  (internal) nodes (this is a random variable, defined on words of length  $n$ ).

**Theorem 1** *For every positive  $q < 1$  the height  $H_n$  of  $q$ -binary search trees with  $n$  (internal) nodes satisfies a central limit theorem of the form*

$$\sup_{x \in \mathbb{R}} \left| \Pr\{H_n \leq x\} - \Phi\left(\frac{x - pn}{\sqrt{npq}}\right) \right| = O(n^{-1/2}), \quad (1)$$

where  $\Phi(x)$  denotes the normal distribution function

$$\Phi(x) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^x e^{-t^2/2} dt.$$

The expected value is given by

$$\mathbb{E}H_n = pn + O\left(\frac{1}{p}\right), \quad (2)$$

in which the  $O$ -constant is uniform for  $0 < q < 1$ . The variance can be estimated by

$$\mathbb{V}H_n = pqn + O_q\left(n^{1/2} \log^2 n\right). \quad (3)$$

The  $O$ -constant in the error term of the variance is not uniform for  $0 < q < 1$ . The dependency on  $q$  has not been worked out since the order of magnitude of the error term is surely not optimal. However, although the error term for the expected value is uniform, this theorem cannot be used to cover the case  $q \rightarrow 1$  where it is known that the expected value is given by

$$\mathbb{E}H_n = c \log n + O(\log \log n)$$

(with  $c = 4.31107 \dots$ , see [3]) and the variance is bounded:

$$\mathbb{V}H_n = O(1)$$

(see [5, 6]).

Intuitively, this result (and its proof) says that the height is dominated by the number of 1's in the sequence. This seems to be due to the fact that most 1's appear consecutively in the sequence, producing rather skew subtrees. The contribution from 2's is asymptotically negligible since there are much fewer of them, etc.

Studying the combinatorics (of words) of geometrically distributed random variables has been a long term project of one of us (H. P.), and further papers can be found on this author's webpage. We do not want to give more details here, but in all previous cases the correspondence between the model of words and its limit (permutations) led to very satisfactory results.

We would like to cite the paper [1] which is of general interest in this context.

## 2 Weak Convergence

**Lemma 1** Let  $y_k(x)$  denote the generating function

$$y_k(x) = \sum_{n \geq 0} \Pr\{H_n \leq k\} x^n.$$

Then we have  $y_0(x) \equiv 0$ ,  $y_k(0) = 1$  and

$$y_{k+1}(x) = px y_k(x) y_k(qx) + y_{k+1}(qx) \quad (4)$$

for  $k \geq 0$ .

*Proof.* The proof is a straightforward translation of the basic decomposition  $w = xay$ : If  $a = 1$  (described by  $px$ ), then the left subtree consists only of letters in  $\{2, 3, \dots\}$ , described by  $y_k(qx)$ , but the right subtree can have any letters, which is described by  $y_k(x)$ . However, if  $a > 1$ , we get the term  $y_{k+1}(qx)$  (we might then think of all letters being reduced by 1).  $\square$

If we write this in the form

$$\frac{y_{k+1}(x) - y_{k+1}(qx)}{(1-q)x} = y_k(qx)y_k(x),$$

then the limit  $q \rightarrow 1$  gives

$$y'_{k+1}(x) = y_k^2(x),$$

the usual recursion in the instance of binary search trees. (Recall that  $\frac{f(x) - f(qx)}{(1-q)x} =: (D_q f)(x)$ , and  $D_q$  is called the  $q$ -difference operator.) Also, if we write

$$y_k(x) = \sum_{0 \leq j < 2^k} a_{k,j} x^j,$$

then we get, by comparing coefficients,

$$a_{k+1,j} = \frac{p}{1-q^j} \sum_{0 \leq i < j} a_{k,i} a_{k,j-1-i} q^i \quad \text{for } j \geq 1,$$

and  $a_{k,0} = 1$  for all  $k \geq 0$ . Thus, Theorem 1 may be reformulated in terms of  $a_{k,j}$ . However, in this paper we will not make use of this notion.

**Lemma 2** The generating functions  $y_k(x)$  ( $k \geq 0$ ,  $0 \leq x < 1$ ) are bounded above by

$$y_k(x) \leq \frac{1}{1-x} - \frac{x}{1-x} \left( \frac{px}{1-qx} \right)^k. \quad (5)$$

Furthermore, this inequality is also true on the level of coefficients, i. e. for  $n \geq 1$  and  $k \geq 1$

$$\begin{aligned} \Pr\{H_n \leq k\} &= [x^n] y_k(x) \\ &\leq 1 - [x^n] \frac{x}{1-x} \left( \frac{px}{1-qx} \right)^k \\ &= \sum_{l=0}^{k-1} \binom{n-1}{l} p^l q^{n-1-l}. \end{aligned} \quad (6)$$

Note that this upper bound is an exact binomial distribution function which is asymptotically normal with mean  $p(n-1) + 1 = pn + q$  and variance  $pq(n-1)$ .

*Proof.* Since  $y_k(x) \leq 1/(1-x)$  we get from (4) that

$$y_{k+1}(x) \leq \frac{px}{1-qx} y_k(x) + \frac{1}{1-qx}.$$

Thus, (5) follows by induction. Note that any step in these calculations is also true on the level of coefficients. Thus,

$$\Pr\{H_n \leq k\} \leq 1 - [x^n] \frac{x}{1-x} \left( \frac{px}{1-qx} \right)^k.$$

Set

$$Y_k(x) = \frac{1}{1-x} - \frac{x}{1-x} \left( \frac{px}{1-qx} \right)^k$$

and

$$\Delta_k(x) = Y_k(x) - Y_{k-1}(x) = \frac{1}{p} \left( \frac{px}{1-qx} \right)^k.$$

Since

$$[x^n] \Delta_k(x) = \binom{n-1}{k-1} p^{k-1} q^{n-k},$$

(6) follows immediately. □

**Lemma 3** *The generating functions  $y_k(x)$  ( $k \geq 0$ ,  $0 \leq x < 1$ ) are bounded below by*

$$\begin{aligned} y_k(x) &\geq \frac{1}{1-x} - \frac{x}{1-x} \left( \frac{px}{1-qx} \right)^k \\ &\quad - \frac{(qx)^2 \left( \frac{px}{1-qx} \right)^k - (qx)^k}{1-qx \frac{px}{1-qx} - qx} \\ &\quad - \frac{qx}{1-x} \left( \frac{px}{1-qx} \right) \frac{\left( \frac{px}{1-qx} \right)^k - (qx)^k}{\frac{px}{1-qx} - qx} \\ &\quad + \frac{qx^2}{1-x} \left( \frac{px}{1-qx} \right)^k \frac{1 - (qx)^k}{1-qx}. \end{aligned} \tag{7}$$

Furthermore, this inequality is also true on the level of coefficients.

*Proof.* First we use the trivial lower bound

$$y_k(x) \geq \frac{1-x^{k+1}}{1-x}$$

and the upper bound (5) to obtain the inequality

$$\begin{aligned}
y_{k+1}(x) &= pxy_k(x)y_k(qx) + y_{k+1}(qx) \\
&\geq x(1-q)y_k(x) \left( \frac{1}{1-qx} - \frac{(qx)^{k+1}}{1-qx} \right) + \frac{1}{1-qx} - \frac{(qx)^{k+2}}{1-qx} \\
&\geq \frac{px}{1-qx} y_k(x) + \frac{1}{1-qx} \\
&\quad - px \frac{(qx)^{k+1}}{1-qx} \left( \frac{1}{1-x} - \frac{x}{1-x} \left( \frac{px}{1-qx} \right)^k \right) - \frac{(qx)^{k+2}}{1-qx}.
\end{aligned}$$

Now (7) follows by induction. As in the proof of Lemma 2 it can be observed that (7) is also true on the level of coefficients.  $\square$

In order to obtain proper error terms for the lower bound for  $\Pr\{H_n \leq k\}$  which may derived from Lemma 3 we make use of the following lemma.

**Lemma 4** *Let  $0 < q < 1$  be given and let  $F(x)$  be a function which is analytic for  $|x| < 1 + \varepsilon$  for some  $\varepsilon > 0$ . Then there exist  $c > 0$  and  $\eta > 0$  (depending on  $q$ ) such that*

$$[x^n]F(x) \left( \frac{px}{1-qx} \right)^k = O\left(n^{-1/2} \exp\left(-\frac{c}{n}(k-pn)^2\right)\right) \quad (8)$$

and

$$[x^n]F(x) \frac{\left(\frac{px}{1-qx}\right)^k - (qx)^k}{\frac{px}{1-qx} - qx} = O\left(n^{-1/2} \exp\left(-\frac{c}{n}(k-pn)^2\right)\right) + O(q^k) \quad (9)$$

uniformly for  $|k-pn| \leq \eta n$  as  $n \rightarrow \infty$ .

*Proof.* By using standard saddle point asymptotics (compare with [4]) it follows that

$$\begin{aligned}
[x^n]F(x) \left( \frac{px}{1-qx} \right)^k &= \frac{1}{2\pi i} \int_{|z|=x_0} F(z) \left( \frac{pz}{1-qz} \right)^k \frac{dz}{z^{n+1}} \\
&= O\left(n^{-1/2} x_0^{-n} \left( \frac{px_0}{1-qx_0} \right)^k\right),
\end{aligned}$$

where the saddle point  $x = x_0 = (1 - \frac{k}{n})/q$  is determined by the equation

$$\frac{x \frac{d}{dx} \left( \frac{px}{1-qx} \right)}{\frac{px}{1-qx}} = \frac{n}{k}.$$

In particular we have  $x_0 = 1$  if  $k/n = p$ . Finally a local expansion of

$$x_0^{-n} \left( \frac{px_0}{1-qx_0} \right)^k = \frac{q^{n-k} p^k}{\left(1 - \frac{k}{n}\right)^{n-k} \left(\frac{k}{n}\right)^k}$$

completes the proof of (8).

The proof of (9) runs along similar lines. However, we have to be a little bit more careful. The reason is that the denominator

$$\frac{pz}{1-qz} - qz$$

is singular for  $z_1 = 0$  and for  $z_2 = (-1 + 2q)/q^2$ . Nevertheless the whole part

$$\frac{\left(\frac{pz}{1-qz}\right)^k - (qz)^k}{\frac{pz}{1-qz} - qz}$$

is regular for  $|z| < 1/q$ . Note that  $|z_2| > 1$  for  $0 < q < \sqrt{2} - 1$  and that  $|z_2| < 1$  for  $\sqrt{2} - 1 < q < 1$ . Thus, if  $q \neq \sqrt{2} - 1$  we get similarly to the above

$$\begin{aligned} [x^n]F(x) \frac{\left(\frac{px}{1-qx}\right)^k - (qx)^k}{\frac{px}{1-qx} - qx} &= \frac{1}{2\pi i} \int_{|z|=x_0} F(z) \frac{\left(\frac{pz}{1-qz}\right)^k - (qz)^k}{\frac{pz}{1-qz} - qz} \frac{dz}{z^{n+1}} \\ &= \frac{1}{2\pi i} \int_{|z|=x_0} F(z) \frac{\left(\frac{pz}{1-qz}\right)^k}{\frac{pz}{1-qz} - qz} \frac{dz}{z^{n+1}} - \frac{1}{2\pi i} \int_{|z|=x_0} F(z) \frac{(qz)^k}{\frac{pz}{1-qz} - qz} \frac{dz}{z^{n+1}} \\ &= \frac{1}{2\pi i} \int_{|z|=x_0} F(z) \frac{\left(\frac{pz}{1-qz}\right)^k}{\frac{pz}{1-qz} - qz} \frac{dz}{z^{n+1}} - \frac{1}{2\pi i} \int_{|z|=1} F(z) \frac{(qz)^k}{\frac{pz}{1-qz} - qz} \frac{dz}{z^{n+1}} \\ &= O\left(n^{-1/2} x_0^{-n} \left(\frac{px_0}{1-qx_0}\right)^k\right) + O(q^k). \end{aligned}$$

Note that we just shifted the paths of integration in regions of analyticity if  $\eta$  is chosen sufficiently small such that  $|x_0 - 1| < |1 - z_2|$ .

If  $q = \sqrt{2} - 1$  then there are polar singularities at  $z_2 = -1$  after splitting the integral. However, the residues of the functions involved (at  $z = -1$ ) are both of order

$$O((\sqrt{2} - 1)^k) = O(q^k),$$

which implies that we get the same error term as in the case  $q \neq \sqrt{2} - 1$ .  $\square$

Note further that (with a little bit more effort) we could have been much more precise. However, the bound given by Lemma 4 is sufficient for our purposes.

**Lemma 5** For every  $0 < q < 1$  there exist  $c > 0$  and  $\eta > 0$  (depending on  $q$ ) such that

$$\left| \Pr\{H_n \leq k\} - \sum_{l=0}^k \binom{n-1}{l} p^l q^{n-1-l} \right| = O\left(n^{-1/2} \exp\left(-\frac{c}{n}(k-pn)^2\right)\right) \quad (10)$$

uniformly for  $|k - pn| \leq \eta n$  and

$$|\Pr\{|H_n - pn| \geq m\}| = O\left(\exp\left(-c\frac{m^2}{n}\right)\right) \quad (11)$$

uniformly for  $|m| \leq \eta n$  as  $n \rightarrow \infty$ .

*Proof.* First we note that (9) of Lemma 4 applies to the function

$$\frac{(qx)^2 \left(\frac{px}{1-qx}\right)^k - (qx)^k}{1 - qx \frac{px}{1-qx} - qx}.$$

Thus we obtain

$$[x^n] \frac{(qx)^2 \left(\frac{px}{1-qx}\right)^k - (qx)^k}{1 - qx \frac{px}{1-qx} - qx} = O\left(n^{-1/2} \exp\left(-\frac{c}{n}(k - pn)^2\right)\right) + O(q^k)$$

uniformly for  $|k - pn| \leq \eta n$  as  $n \rightarrow \infty$ .

The remaining terms will now be treated in the same way as in the proof of Lemma 4:

$$\begin{aligned} & [x^n] \left( \frac{qx}{1-x} \left( \frac{px}{1-qx} \right) \frac{\left(\frac{px}{1-qx}\right)^k - (qx)^k}{\frac{px}{1-qx} - qx} + \frac{qx^2}{1-x} \left( \frac{px}{1-qx} \right)^k \frac{1 - (qx)^k}{1 - qx} \right) \\ &= \frac{1}{2\pi i} \int_{|z|=1-\varepsilon} \left( \frac{qz}{1-z} \left( \frac{pz}{1-qz} \right) \frac{\left(\frac{pz}{1-qz}\right)^k - (qz)^k}{\frac{pz}{1-qz} - qz} - \frac{qz^2}{1-z} \left( \frac{pz}{1-qz} \right)^k \frac{1 - (qz)^k}{1 - qz} \right) \frac{dz}{z^{n+1}} \\ &= \frac{1}{2\pi i} \int_{|z|=x_0} \left( \frac{qz}{1-z} \left( \frac{pz}{1-qz} \right) \frac{1}{\frac{pz}{1-qz} - qz} - \frac{qz^2}{1-z} \frac{1}{1 - qz} \right) \left( \frac{pz}{1-qz} \right)^k \frac{dz}{z^{n+1}} \\ &\quad - \frac{1}{2\pi i} \int_{|z|=1-k^{-1}} \frac{qz}{1-z} \left( \frac{pz}{1-qz} \right) \frac{(qz)^k}{\frac{pz}{1-qz} - qz} \frac{dz}{z^{n+1}} \\ &\quad + \frac{1}{2\pi i} \int_{|z|=1-k^{-1}} \frac{qz^2}{1-z} \left( \frac{pz}{1-qz} \right)^k \frac{(qz)^k}{1 - qz} \frac{dz}{z^{n+1}} \\ &= O\left(n^{-1/2} x_0^{-n} \left( \frac{px_0}{1-qx_0} \right)^k\right) + O(kq^k). \end{aligned}$$

For the first integral we use the fact that the function

$$F(z) = \frac{qz}{1-z} \left( \frac{pz}{1-qz} \right) \frac{1}{\frac{pz}{1-qz} - qz} - \frac{qz^2}{1-z} \frac{1}{1 - qz}$$



is regular at  $z = 1$  und thus bounded in a vicinity of  $z = 1$  (compare with the proof of Lemma 4). The second and third integrals are easy to estimate. We only have to apply the trivial bounds

$$\begin{aligned} \max_{|z|=1-k^{-1}} \left| \frac{1}{1-z} \right| &= k, \\ \max_{|z|=1-k^{-1}} |z^k| &= O(1) \end{aligned}$$

and

$$\max_{|z|=1-k^{-1}} \left| \frac{pz}{1-qz} \right|^k = (1 + O(k^{-1}))^k = O(1).$$

Finally observe that in the range  $|k - pn| \leq \eta n$  we surely have

$$kq^k = O\left(n^{-1/2} \exp\left(-\frac{c}{n}(k - pn)^2\right)\right)$$

which completes the proof of (10).

The proof of (11) is now easy. We just have to combine a proper tail estimate for the binomial distribution, i. e.

$$\sum_{|l-pn| \geq m} \binom{n-1}{l} p^l q^{n-1-l} = O\left(\exp\left(-c \frac{m^2}{n}\right)\right),$$

with (10). □

We are now able to complete the first part of Theorem 1. It is well known that the distribution function of the binomial distribution can be estimated by the distribution function of the normal distribution up to a uniform error of order  $O(n^{-1/2})$ , i. e., as  $n \rightarrow \infty$ ,

$$\sup_{x \in \mathbb{R}} \left| \sum_{l \leq x} \binom{n-1}{l} p^l q^{n-1-l} - \Phi\left(\frac{x - pn}{\sqrt{pqn}}\right) \right| = O(n^{-1/2}), \quad (12)$$

compare with [7, p. 542].

Furthermore, by (10) we get a similar result for the distribution of  $H_n$ :

$$\sup_{x \in \mathbb{R}} \left| \Pr\{H_n \leq x\} - \sum_{l \leq x} \binom{n-1}{l} p^l q^{n-1-l} \right| = O(n^{-1/2}). \quad (13)$$

Note that (10) provides (13) just for  $x$  with  $|x - pn| \leq \eta n$ . However, by (11) we know that

$$|\Pr\{|H_n - pn| \geq \eta n\}| = O(\exp(-c\eta^2 n)).$$

By the monotonicity of the distribution function this implies that for  $x$  with  $|x - pn| \geq \eta n$  we also have

$$|\Pr\{|H_n - pn| \geq x\}| = O(\exp(-c\eta^2 n)).$$

Thus (13) follows.

The first part of Theorem 1, i. e. (1), is now a trivial consequence of (12) and (13).

### 3 Convergence of Moments

Lemmata 2 and 3 can be easily used to get quite tight bounds for the expected value. Note that

$$\sum_{n \geq 0} \mathbb{E} H_n x^n = \sum_{k \geq 0} \left( \frac{1}{1-x} - y_k(x) \right). \quad (14)$$

**Lemma 6** *The expected value of  $H_n$  can be bounded by*

$$pn + q \leq \mathbb{E} H_n \leq pn + O\left(\frac{1}{p}\right). \quad (15)$$

*Proof.* Set

$$E(x) = \sum_{k \geq 0} \left( \frac{1}{1-x} - y_k(x) \right).$$

Then by (5) we obtain

$$E(x) \geq \frac{x}{1-x} \sum_{k \geq 0} \left( \frac{px}{1-qx} \right)^k = \frac{x(1-qx)}{(1-x)^2},$$

which is also true on the level of coefficients. Thus, we have

$$\mathbb{E} H_n \geq pn + q.$$

Similarly we get the upper bound. From

$$\begin{aligned} E(x) &\leq \frac{x(1-qx)}{(1-x)^2} + \frac{(qx)^2}{1-qx} \frac{1}{\left(\frac{px}{1-qx}\right) - qx} \left( \frac{1}{1 - \frac{px}{1-qx}} - \frac{1}{1-qx} \right) \\ &\quad + \frac{qx}{1-x} \left( \frac{px}{1-qx} \right) \frac{1}{\left(\frac{px}{1-qx}\right) - qx} \left( \frac{1}{1 - \frac{px}{1-qx}} - \frac{1}{1-qx} \right) \\ &\quad - \frac{qx^2}{1-x} \frac{1}{1-qx} \left( \frac{1}{1 - \frac{px}{1-qx}} - \frac{1}{1 - \frac{qx^2(1-q)}{1-qx}} \right) \\ &= \frac{x(1-qx)}{(1-x)^2} + \frac{(qx)^2}{(1-x)(1-qx)} \\ &\quad + \frac{pqx^2}{(1-x)^2(1-qx)} - \frac{qx^2}{(1-x)^2} + \frac{qx^2}{(1-x) \left(1 - \frac{pqx^2}{1-qx}\right)} \\ &= \frac{x(1-qx)}{(1-x)^2} + \frac{qx^2(1-qx)}{(1-x)(1-qx - pqx^2)} \end{aligned}$$

we directly obtain

$$\mathbb{E} H_n \leq pn + O\left(\frac{1}{p}\right).$$

This completes the proof of Lemma 6.  $\square$

Unfortunately it seems that we cannot prove similarly tight estimates for the variance. However, we can obtain a non-trivial result:

**Lemma 7** *The variance of  $H_n$  can be bounded by*

$$\mathbb{V}H_n = pqn + O_q\left(n^{1/2} \log^2 n\right). \quad (16)$$

*Proof.* If  $F(x) = \Pr\{X \leq x\}$  denotes the distribution function of a random variable  $X$  (of compact support) then the variance of  $X$  can be represented by

$$\mathbb{V}X = 2 \int_{-\infty}^{\mathbb{E}X} (\mathbb{E}X - y)F(y) dy + 2 \int_{\mathbb{E}X}^{\infty} (y - \mathbb{E}X)(1 - F(y)) dy.$$

We apply this formula for the height  $H_n$  where we have very precise estimates for  $\mathbb{E}H_n$  and its distribution function  $F_n(x) = \Pr\{H_n \leq x\}$ , compare with Lemma 5 and Lemma 6.

Let  $u \leq \eta\sqrt{n}$  be a parameter to be defined later. By applying Lemma 5 and Lemma 6 we get

$$\begin{aligned} 2 \int_0^{\mathbb{E}H_n - u\sqrt{n}} (\mathbb{E}H_n - y)F_n(y) dy &= O\left(ne^{-cu^2}\right), \\ 2 \int_{\mathbb{E}H_n - u\sqrt{n}}^{\mathbb{E}H_n} (\mathbb{E}H_n - y)F_n(y) dy &= \frac{pqn}{2} + O(u^2\sqrt{n}), \\ 2 \int_{\mathbb{E}H_n}^{\mathbb{E}H_n - u\sqrt{n}} (y - \mathbb{E}H_n)(1 - F_n(y)) dy &= \frac{pqn}{2} + O(u^2\sqrt{n}), \end{aligned}$$

and

$$2 \int_{\mathbb{E}H_n - u\sqrt{n}}^n (y - \mathbb{E}H_n)(1 - F_n(y)) dy = O\left(ne^{-cu^2}\right).$$

Choosing  $u = \log^2 n$  gives the result.  $\square$

It should be further mentioned that it is quite easy to get asymptotic relations for all central moments of the form

$$\mathbb{E}(H_n - \mathbb{E}H_n)^k \sim c_k (pqn)^{k/2} \quad (n \rightarrow \infty)$$

by the method as described in Lemma 7, where  $c_{2k} = (2k)!/(2^k k!)$  and  $c_{2k+1} = 0$ . Thus we have not only a weak convergence result for  $H_n$  (properly normalized) but convergence of moments, too.

## References

- [1] M. BARLOW, R. PEMANTLE AND E. PERKINS, *Diffusion limited aggregation on a homogeneous tree*, Prob. Theory and Related Fields **107** (1997), 1–60.
- [2] F. BERGERON, P. FLAJOLET AND B. SALVY, *Varieties of increasing trees*, Lecture Notes in Computer Science **581** (1992), 24–48.
- [3] L. DEVROYE AND B. REED, *On the variance of the height of random binary search trees*, SIAM J. Comput. **24** (1995), 1157–1162.
- [4] M. DRMOTA, *A bivariate asymptotic expansion of coefficients of powers of generating functions*, Europ. J. Combinatorics **15** (1994), 139–152.
- [5] M. DRMOTA, *The variance of the height of binary search trees*, Theoret. Comput. Sci. **270** (2002), 913–919.
- [6] B. REED, *The height of a random binary search tree*, J. Assoc. Comput. Mach., to appear.
- [7] W. FELLER, *An Introduction to Probability Theory and Its Applications, Vol. II*, 2<sup>nd</sup> ed., J. Wiley, New York, 1971.
- [8] H. PRODINGER, *A  $q$ -analogue of the path length of binary search trees*, Algorithmica **31** (2001), 433–441.

