# A $q$–ANALOGUE OF THE PATH LENGTH OF BINARY SEARCH TREES

## HELMUT PRODINGER

ABSTRACT. A reformulation of the path length of binary search trees is given in terms of permutations, allowing to extend the definition to the instance of words, where the letters are obtained by independent geometric random variables (with parameter $q$). In this way, expressions for expectation and variance are obtained which in the limit for $q \to 1$ are the classical expressions.

The path length $\rho(t)$ of a binary search tree $t$ satisfies the recursion $\rho(t) = \rho(t_L) + \rho(t_R) + |t_L| + |t_R|$ where $t_L$ and $t_R$ are the left resp. right subtree of the root. ($|t|$ denotes the size of the tree $t$, i. e. the number of nodes.)

Binary search trees are obtained from permutations. For some background see [4, 1, 2]. Our aim is to rewrite the definition of the path length in terms of permutations, since then we are able to obtain $q$–analogues: This is done by considering words over the alphabet $\{1, 2, \dots\}$ instead, with probabilities $p, pq, pq^2, \dots$, where $p + q = 1$ (geometric probabilities). In the limit $q \to 1$, this model turns into the model of random permutations, as equal letters appear with probability 0 and each relative ordering is equally likely.

For a permutation $\pi = \pi_1 \dots \pi_n$ we define $\rho(\pi)$ by

$$\rho(\pi) = \Big|\{(j,k) \mid 1 \le j < k \le n, \ \pi_j = \min\{\pi_j, \dots, \pi_k\} \quad \text{or} \quad \pi_k = \min\{\pi_j, \dots, \pi_k\}\}\Big|.$$

Then $\rho(\square) = 0$ and, if $\pi = \sigma 1 \tau$, then $\rho(\pi) = \rho(\sigma) + \rho(\tau) + |\sigma| + |\tau|$, as pairs with the left coordinate in $\sigma$ and the right coordinate in $\tau$ are definitely not counted.

But this definition of $\pi$ can be taken as it is where $\pi_1 \dots \pi_n$ now denotes a *word* over the alphabet $\{1, 2, \dots\}$. This will be our starting point.

We want to point out that our previous paper [3] contains easier but related parameters.

In the sequel we want to compute the expectation and the variance of the parameter $\rho$, for random words of length $n$. We define random variables

$$L_{jk} = \begin{cases} 1 & \text{if } \pi_j = \min\{\pi_j, \ldots, \pi_k\} \\ 0 & \text{otherwise,} \end{cases}$$

$$R_{jk} = \begin{cases} 1 & \text{if } \pi_k = \min\{\pi_j, \ldots, \pi_k\} \\ 0 & \text{otherwise,} \end{cases}$$

$$B_{jk} = L_{jk} \cdot R_{jk},$$

$$N_{jk} = (1 - L_{jk}) \cdot (1 - R_{jk}).$$

(The letters $L, R, B, N$ are chosen to indicate *left, right, both, not.*)

Then the parameter $\rho$ may be described as

$$\rho = \sum_{1 \le j < k \le n} \Big[ L_{jk} + R_{jk} - B_{jk} \Big].$$

Now we can introduce the generating function

$$f(v) = \left(\frac{p}{q}\right)^n \sum_{i_1,\ldots,i_n \ge 1} q^{i_1 + \cdots + i_n} \prod_{1 \le j < k \le n} \Big[ L_{jk}v + R_{jk}v - B_{jk}v + N_{jk} \Big];$$

the coefficient of $v^k$ in $f(v)$ is the probability that parameter $\rho$ has value $k$, assuming random words of length $n$.

As always, the expected value is obtained via $\mathbb{E} = f'(1)$;

$$\mathbb{E} = \left(\frac{p}{q}\right)^n \sum_{i_1,\ldots,i_n \ge 1} q^{i_1 + \cdots + i_n} \sum_{1 \le j < k \le n} \left( L_{jk} + R_{jk} - B_{jk} \right)$$

$$= \sum_{1 \le j < k \le n} \left(\frac{p}{q}\right)^{k+1-j} \Big[ \sum_{L_{jk}=1} q^{i_j + \cdots + i_k} + \sum_{R_{jk}=1} q^{i_j + \cdots + i_k} - \sum_{B_{jk}=1} q^{i_j + \cdots + i_k} \Big]$$

$$= \sum_{1 \le j < k \le n} \left(\frac{p}{q}\right)^{k+1-j} \Big[ 2 \sum_{i \ge 1} q^{i(k+1-j)} \frac{1}{p^{k-j}} - \sum_{i \ge 1} q^{i(k+1-j)} \frac{1}{p^{k-j-1}} \Big]$$

$$= (2p - p^2) \sum_{1 \le j < k \le n} \left(\frac{1}{q}\right)^{k+1-j} \sum_{i \ge 1} q^{i(k+1-j)}$$

$$= p(2-p) \sum_{1 \le j < k \le n} \frac{1}{1 - q^{k+1-j}}$$

$$= p(2-p) \sum_{2 \le i \le n} \frac{n+1-i}{1 - q^i}$$

$$= p(2-p) \sum_{1 \le i \le n} \frac{n+1-i}{1 - q^i} - n(2-p).$$

The terms that would survive the limit $q \to 1$ are

$$2p \sum_{1 \leq i \leq n} (n + 1 - i) \frac{1}{1 - q^i} - 2n,$$

and the limit is

$$\lim_{q \to 1} \mathbb{E} = 2 \sum_{1 \leq i \leq n} \frac{n + 1 - i}{i} - 2n = 2(n + 1)H_n - 4n,$$

as is of course well known.

Now we turn to the variance, and this is much harder, since we must first compute the second factorial moment, which is obtained by a second derivative;

$$\mathbb{E}^{\underline{2}} = \left(\frac{p}{q}\right)^n \sum_{i_1, \ldots, i_n \geq 1} q^{i_1 + \cdots + i_n} \times$$

$$\times \sum_{1 \leq j < k \leq n, 1 \leq l < m \leq n, (j,k) \neq (l,m)} \left(L_{jk} + R_{jk} - B_{jk}\right)\left(L_{lm} + R_{lm} - B_{lm}\right)$$

$$= \left(\frac{p}{q}\right)^n \sum_{i_1, \ldots, i_n \geq 1} q^{i_1 + \cdots + i_n} \times$$

$$\times \sum_{1 \leq j < k \leq n, 1 \leq l < m \leq n, (j,k) \neq (l,m)} \left(2L_{jk}L_{lm} + 2L_{jk}R_{lm} - 4L_{jk}B_{lm} + B_{jk}B_{lm}\right)$$

$$= 2\Xi^{\mathrm{LL}} + 2\Xi^{\mathrm{LR}} - 4\Xi^{\mathrm{LB}} + \Xi^{\mathrm{BB}}$$

(using several symmetries).

The range $\Lambda = \{1 \leq j < k \leq n, 1 \leq l < m \leq n, (j,k) \neq (l,m)\}$ must be split into the following 12 disjoint subranges:

$$\Lambda_1 = \{1 \leq j < k < l < m \leq n\},$$
$$\Lambda_2 = \{1 \leq j < l < m < k \leq n\},$$
$$\Lambda_3 = \{1 \leq j < l < k < m \leq n\},$$
$$\Lambda_4 = \{1 \leq j < k = l < m \leq n\},$$
$$\Lambda_5 = \{1 \leq j < l < m = k \leq n\},$$
$$\Lambda_6 = \{1 \leq j = l < k < m \leq n\},$$
$$\Lambda_7 = \{1 \leq l < m < j < k \leq n\},$$
$$\Lambda_8 = \{1 \leq l < j < k < m \leq n\},$$
$$\Lambda_9 = \{1 \leq l < j < m < k \leq n\},$$
$$\Lambda_{10} = \{1 \leq l < m = j < k \leq n\},$$
$$\Lambda_{11} = \{1 \leq l < j < k = m \leq n\},$$
$$\Lambda_{12} = \{1 \leq l = j < m < k \leq n\}.$$

And we will have contributions $\Theta_i^{\mathrm{LL}}$, $\Theta_i^{\mathrm{LR}}$, $\Theta_i^{\mathrm{LB}}$, $\Theta_i^{\mathrm{BB}}$, to $\Xi^{\mathrm{LL}}$, $\Xi^{\mathrm{LR}}$, $\Xi^{\mathrm{LB}}$, $\Xi^{\mathrm{BB}}$, for $i = 1, \ldots, 12$, according to the 12 ranges $\Lambda_i$.

Therefore we must compute 48 (not necessarily) different contributions.

For convenience, we state them as a lemma.

**Lemma 1.** *The contributions* $\Theta_i^{LL}$, $\Theta_i^{LR}$, $\Theta_i^{LB}$, $\Theta_i^{BB}$, *for* $i = 1, \ldots, 12$, *are given by*

$$\Theta_1^{LL} = \Theta_7^{LL} = p^2 \sum_{1 \le j < k < l < m \le n} \frac{1}{1 - q^{k+1-j}} \frac{1}{1 - q^{m+1-l}},$$

$$\Theta_2^{LL} = \Theta_8^{LL} = p^2 \sum_{1 \le j < l < m < k \le n} \frac{1}{1 - q^{k+1-j}} \frac{1}{1 - q^{m+1-l}},$$

$$\Theta_3^{LL} = \Theta_9^{LL} = p^2 \sum_{1 \le j < l < k < m \le n} \frac{1}{1 - q^{m+1-j}} \frac{1}{1 - q^{m+1-l}},$$

$$\Theta_4^{LL} = \Theta_{10}^{LL} = p^2 \sum_{1 \le j < k = l < m \le n} \frac{1}{1 - q^{m+1-j}} \frac{1}{1 - q^{m+1-k}},$$

$$\Theta_5^{LL} = \Theta_{11}^{LL} = p^2 \sum_{1 \le j < l < m = k \le n} \frac{1}{1 - q^{m+1-j}} \frac{1}{1 - q^{m+1-l}},$$

$$\Theta_6^{LL} = \Theta_{12}^{LL} = p \sum_{1 \le j = l < k < m \le n} \frac{1}{1 - q^{m+1-j}};$$

$$\Theta_1^{LR} = \Theta_7^{LR} = p^2 \sum_{1 \le j < k < l < m \le n} \frac{1}{1 - q^{k+1-j}} \frac{1}{1 - q^{m+1-l}},$$

$$\Theta_2^{LR} = \Theta_8^{LR} = p^2 \sum_{1 \le j < l < m < k \le n} \frac{1}{1 - q^{k+1-j}} \frac{1}{1 - q^{m+1-l}},$$

$$\Theta_3^{LR} = \Theta_9^{LR} = p^2 \sum_{1 \le j < l < k < m \le n} \left[ \frac{1}{1 - q^{m+1-j}} \frac{1}{1 - q^{m+1-l}} + \frac{1}{1 - q^{m+1-j}} \frac{1}{1 - q^{k+1-j}} - \frac{1}{1 - q^{m+1-j}} \right],$$

$$\Theta_4^{LR} = p^2 \sum_{1 \le j < k = l < m \le n} \left[ \frac{1}{1 - q^{m+1-j}} \frac{1}{1 - q^{m+1-k}} + \frac{1}{1 - q^{m+1-j}} \frac{1}{1 - q^{k+1-j}} - \frac{1}{1 - q^{m+1-j}} \right],$$

$$\Theta_5^{LR} = p^2 \sum_{1 \le j < l < m \le n} \frac{1}{1 - q^{m+1-j}} \frac{1}{1 - q^{m+1-l}},$$

$$\Theta_6^{LR} = p^2 \sum_{1 \le j < k < m \le n} \frac{1}{1 - q^{m+1-j}} \frac{1}{1 - q^{k+1-j}},$$

$$\Theta_{10}^{LR} = p \sum_{1 \le l < m = j < k \le n} \frac{1}{1 - q^{k+1-l}},$$

$$\Theta_{11}^{LR} = \Theta_{12}^{LR} = p^2 \sum_{1 \le l = j < m < k \le n} \frac{1}{1 - q^{k+1-l}};$$

$$\Theta_1^{LB} = \Theta_7^{LB} = p^3 \sum_{1 \leq j < k < l < m \leq n} \frac{1}{1 - q^{k+1-j}} \frac{1}{1 - q^{m+1-l}},$$

$$\Theta_2^{LB} = \Theta_8^{LB} = p^3 \sum_{1 \leq j < l < m < k \leq n} \frac{1}{1 - q^{k+1-j}} \frac{1}{1 - q^{m+1-l}},$$

$$\Theta_3^{LB} = \Theta_9^{LB} = p^3 \sum_{1 \leq j < l < k < m \leq n} \frac{1}{1 - q^{m+1-j}} \frac{1}{1 - q^{m+1-l}},$$

$$\Theta_4^{LB} = p^3 \sum_{1 \leq j < k = l < m \leq n} \frac{1}{1 - q^{m+1-j}} \frac{1}{1 - q^{m+1-k}},$$

$$\Theta_5^{LB} = p^3 \sum_{1 \leq j < l < m \leq n} \frac{1}{1 - q^{m+1-j}} \frac{1}{1 - q^{m+1-l}},$$

$$\Theta_6^{LB} = p^2 \sum_{1 \leq j < k < m \leq n} \frac{1}{1 - q^{m+1-j}},$$

$$\Theta_{10}^{LB} = p^2 \sum_{1 \leq l < m = j < k \leq n} \frac{1}{1 - q^{k+1-l}},$$

$$\Theta_{11}^{LB} = p^3 \sum_{1 \leq l < j < k = m \leq n} \frac{1}{1 - q^{k+1-l}},$$

$$\Theta_{12}^{LB} = p^2 \sum_{1 \leq l = j < m < k \leq n} \frac{1}{1 - q^{k+1-l}};$$

$$\Theta_1^{BB} = \Theta_7^{BB} = p^4 \sum_{1 \leq j < k < l < m \leq n} \frac{1}{1 - q^{k+1-j}} \frac{1}{1 - q^{m+1-l}},$$

$$\Theta_2^{BB} = \Theta_8^{BB} = p^4 \sum_{1 \leq j < l < m < k \leq n} \frac{1}{1 - q^{k+1-j}} \frac{1}{1 - q^{m+1-l}},$$

$$\Theta_3^{BB} = \Theta_9^{BB} = p^4 \sum_{1 \leq j < l < k < m \leq n} \frac{1}{1 - q^{m+1-j}},$$

$$\Theta_4^{BB} = \Theta_{10}^{BB} = p^3 \sum_{1 \leq j < k = l < m \leq n} \frac{1}{1 - q^{m+1-j}},$$

$$\Theta_5^{BB} = \Theta_{11}^{BB} = p^3 \sum_{1 \leq j < l < m \leq n} \frac{1}{1 - q^{m+1-j}} \frac{1}{1 - q^{m+1-l}},$$

$$\Theta_6^{BB} = \Theta_{12}^{BB} = p^3 \sum_{1 \leq j = l < k < m \leq n} \frac{1}{1 - q^{m+1-j}}.$$

*Proof.* The computations are as (or slightly more complicated than) the one for the expected value. We don't give more details.  □

We simplify those sums and write $a_i = \frac{1}{1-q^i}$ for convenience.

**Lemma 2.**

$$\Xi^{LL} = 2p^2 \sum_{2 \le i,j \le n-2, i+j \le n} a_i a_j \binom{n+2-i-j}{2}$$

$$+ 2p^2 \sum_{2 \le i < j \le n} a_i a_j (n+1-j)(j-i-1)$$

$$+ 2p^2 \sum_{3 \le i < j \le n} a_i a_j (n+1-j)(i-2)$$

$$+ 4p^2 \sum_{2 \le i < j \le n} a_i a_j (n+1-j)$$

$$+ 2p \sum_{3 \le i \le n} a_i (i-2)(n+1-i),$$

$$\Xi^{LR} = 2p^2 \sum_{2 \le i,j \le n-2, i+j \le n} a_i a_j \binom{n+2-i-j}{2}$$

$$+ 2p^2 \sum_{2 \le i < j \le n} a_i a_j (n+1-j)(j-i-1)$$

$$+ 4p^2 \sum_{3 \le i < j \le n} a_i a_j (n+1-j)(i-2)$$

$$- 2p^2 \sum_{4 \le i \le n} a_i \binom{i-2}{2}(n+1-i)$$

$$+ 4p^2 \sum_{2 \le i < j \le n} a_i a_j (n+1-j)$$

$$+ p^2 \sum_{3 \le i \le n} a_i (i-2)(n+1-i)$$

$$+ p \sum_{3 \le i \le n} a_i (i-2)(n+1-i),$$

$$\Xi^{LB} = 2p^3 \sum_{2 \le i,j \le n-2, i+j \le n} a_i a_j \binom{n+2-i-j}{2}$$

$$+ 2p^3 \sum_{2 \le i < j \le n} a_i a_j (n+1-j)(j-i-1)$$

$$+ 2p^3 \sum_{3 \le i < j \le n} a_i a_j (n+1-j)(i-2)$$

$$+ 2p^3 \sum_{2 \le i < j \le n} a_i a_j (n+1-j)$$

$$+ (3p^2 + p^3) \sum_{3 \le i \le n} a_i (i-2)(n+1-i),$$

$$\Xi^{BB} = 2p^4 \sum_{2 \le i,j \le n-2, i+j \le n} a_i a_j \binom{n+2-i-j}{2}$$
$$+ 2p^4 \sum_{2 \le i < j \le n} a_i a_j (n+1-j)(j-i-1)$$
$$+ 2p^3 \sum_{2 \le i < j \le n} a_i a_j (n+1-j)$$
$$+ 2p^4 \sum_{4 \le i \le n} a_i \binom{i-2}{2}(n+1-i)$$
$$+ 4p^3 \sum_{3 \le i \le n} a_i (i-2)(n+1-i).$$

$\square$

The variance is given by

$$\mathbb{V} = 2\Xi^{LL} + 2\Xi^{LR} - 4\Xi^{LB} + \Xi^{BB} + \mathbb{E} - \left(\mathbb{E}\right)^2.$$

In order to simplify this expression, we note the following formulæ:

**Lemma 3.**

$$p^2 \sum_{2 \le i,j \le n-2, i+j \le n} a_i a_j \binom{n+2-i-j}{2}$$
$$= 2p^2 \sum_{1 \le i < j \le n} a_i a_j \binom{n+2-j}{2} - p^2 \sum_{1 \le i \le n} a_i \binom{n+2-i}{2}(i-1)$$
$$- 2p \sum_{1 \le i \le n} a_i \binom{n+1-i}{2} + \binom{n}{2}.$$

*Proof.* Note that

$$\frac{1}{1-q^i} \frac{1}{1-q^j} = \frac{1}{1-q^{i+j}} \left( \frac{1}{1-q^i} + \frac{1}{1-q^j} - 1 \right)$$

and do some trivial rearrangements. $\square$

**Lemma 4.**

$$\left( \sum_{1 \le i \le n} a_i (n+1-i) \right)^2 = 2 \sum_{1 \le i < j \le n} a_i a_j (n+1-i)(n+1-j) + \sum_{1 \le i \le n} a_i^2 (n+1-i)^2.$$

*Proof.* Obvious. $\square$

Using these lemmata and numerous simplifications that were partially supported by Maple, we can state our main result:

**Theorem 1.** *The expectation and the variance of the q–ified path length in words of length n, generated by n independent geometric random variables are given by*

$$\mathbb{E} = p(2 - p) \sum_{1 \leq i \leq n} \frac{n + 1 - i}{1 - q^i} - n(2 - p)$$

*and*

$$\mathbb{V} = 2p^2 \sum_{1 \leq i < j \leq n} \frac{(n + 1 - j)(4i + p(5 - 4i))}{(1 - q^i)(1 - q^j)}$$
$$- p^2(2 - p)^2 \sum_{1 \leq i \leq n} \frac{(n + 1 - i)^2}{(1 - q^i)^2}$$
$$+ p \sum_{1 \leq i \leq n} \frac{n + 1 - i}{1 - q^i} \Big( 6i - 2 + p(-4ni + 4n - 19i + 3i^2 + 7)$$
$$+ 4p^2(ni - n + 3i - 1 - i^2) + p^3(-ni + n + 2i^2 - 8i + 8) \Big)$$
$$+ 5pn - 3p^2n - 2p^3n.$$

$\square$

The terms in the variance that would survive the limit $q \to 1$ are these:

$$8p^2 \sum_{1 \leq i < j \leq n} \frac{(n + 1 - j)i}{(1 - q^i)(1 - q^j)} - 4p^2 \sum_{1 \leq i \leq n} \frac{(n + 1 - i)^2}{(1 - q^i)^2} + 2p \sum_{1 \leq i \leq n} \frac{(n + 1 - i)(3i - 1)}{1 - q^i}.$$

The limit is

$$\lim_{q \to 1} \mathbb{V} = 8 \sum_{1 \leq j \leq n} \frac{(n + 1 - j)(j - 1)}{j} - 4 \sum_{1 \leq i \leq n} \frac{(n + 1 - i)^2}{i^2} + 2 \sum_{1 \leq i \leq n} \frac{n + 1 - i}{i}(3i - 1)$$
$$= 4(n + 1)n - 8(n + 1)H_n + 8n - 4(n + 1)^2 H_n^{(2)} + 8(n + 1)H_n - 4n$$
$$+ 3n(n + 1) - 2(n + 1)H_n + 2n$$
$$= 7n^2 - 4(n + 1)^2 H_n^{(2)} - 2(n + 1)H_n + 13n.$$

This is (of course!) the variance in the classical case.

## References

[1] D. E. Knuth. *The Art of Computer Programming*, volume 3: Sorting and Searching. Addison-Wesley, 1973. Second edition, 1998.

[2] H. M. Mahmoud. *Evolution of Random Search Trees*. John Wiley, New York, 1992.

[3] H. Prodinger. Combinatorics of geometrically distributed random variables: Inversions and a parameter of Knuth. *submitted*, 2000.

[4] R. Sedgewick and P. Flajolet. *An Introduction to the Analysis of Algorithms*. Addison-Wesley, 1996.

HELMUT PRODINGER, CENTRE FOR APPLICABLE ANALYSIS AND NUMBER THEORY, DEPARTMENT OF MATHEMATICS, UNIVERSITY OF THE WITWATERSRAND, P. O. WITS, 2050 JOHANNESBURG, SOUTH AFRICA, EMAIL: `helmut@gauss.cam.wits.ac.za`.
HOMEPAGE: `http://www.wits.ac.za/helmut/index.htm`