

RANDOM RECORDS AND CUTTINGS IN COMPLETE BINARY TREES

SVANTE JANSON

ABSTRACT. We study the number of records in a complete binary tree with randomly labeled vertices or edges. Equivalently, we may study the number of random cuttings required to eliminate a complete binary tree.

The distribution is, after normalization, asymptotically a periodic function of $\lg n - \lg \lg n$; thus there is no true asymptotic distribution but a family of limits of different subsequences; these limits are similar to a 1-stable distribution but have some periodic fluctuations.

1. INTRODUCTION

Let each vertex v in a rooted tree T have a random value λ_v attached to it, and assume that these values are i.i.d. with a continuous distribution (so that a.s. there are no ties). Say that a value λ_v is a *record* if it is the smallest value in the path from the root to v . Let $X_v(T)$ denote the (random) number of records. Note that this generalizes the classical record problem (which is the case when T is a path), see for example [9].

Alternatively, we may attach random values to the edges, and let $X_e(T)$ denote the number of edges with record values (along the path from the root).

It is obvious that the choice of common distribution of the labels does not affect the result, and that we as well can count the values that are largest. We can also let the labels be a random permutation of $\{1, \dots, n\}$.

The same random variables appears when we consider random cuttings of the tree T defined as follows, see [6]. Make a random cut by choosing one vertex [edge] at random. Delete this vertex [edge] so that the tree separates into two parts, and keep only the part containing the root. Continue recursively until the root is cut [only the root is left]. Then the (random) number of cuts made is $X_v(T)$ [$X_e(T)$]. (More precisely, these random variables have the same distribution.) This equivalence is shown in [5], where the asymptotic distributions are found for the random trees that can be constructed as conditioned Galton–Watson trees, for example random labelled trees and random binary trees. See also [6, 1, 7, 8] for earlier results.

We will in this paper study the case of a complete binary tree.

The *complete binary tree* with n vertices has height $m = \lfloor \lg n \rfloor$; it has 2^k vertices of height k , $0 \leq k < m$, and $n - 2^m + 1$ vertices of height m , moreover, the vertices of height m have the leftmost positions among the 2^m possible ones, see e.g. [4, page 401]. We denote this rooted tree by T_n , and denote its root by o .

Let $\{x\} := x - \lfloor x \rfloor$ denote the fractional part of a real number x . Further, for a vertex v in a rooted tree, let $h(v)$ be its height (also known as depth), with the root having height 0.

Theorem 1.1. *Suppose that $n \rightarrow \infty$ such that $\{\lg n - \lg \lg n\} \rightarrow \gamma \in [0, 1]$. Then*

$$\left(X_v(T_n) - \frac{n}{\lg n} - \frac{n \lg \lg n}{\lg^2 n} \right) / \frac{n}{\lg^2 n} \xrightarrow{d} -W_\gamma \quad (1.1)$$

where W_γ has an infinitely divisible distribution with characteristic function

$$\mathbb{E} e^{itW_\gamma} = \exp\left(if(\gamma)t + \int_0^\infty (e^{itx} - 1 - itx\mathbf{1}[x < 1]) d\nu_\gamma(x) \right), \quad (1.2)$$

where $f(\gamma) := 2^\gamma - 1 - \gamma$ and the Lévy measure ν_γ is supported on $(0, \infty)$ and has density

$$\frac{d\nu_\gamma}{dx} = 2^{\lfloor \lg x + \gamma \rfloor} x^{-2}.$$

The same result holds for $X_e(T_n)$.

We prove Theorem 1.1 in Section 2. The strategy of the proof is to approximate $X_v(T_n)$ and $X_e(T_n)$ by a sum of independent random variables derived from $\{\lambda_v\}$, see Lemma 2.4. We will then apply a classical limit theorem for triangular arrays.

Remark 1.2. Let \tilde{X}_n denote the left hand side of (1.1). Instead of stating the result for suitable subsequences, we may say that \tilde{X}_n has approximately the same distribution as $-W_{\{\lg n - \lg \lg n\}}$ for large n ; more precisely, the distance between the two distributions (in for example the Lévy metric) tends to 0 as $n \rightarrow \infty$.

Remark 1.3. Most records occur at height close to the maximum $m \approx \lg n$, simply because almost all vertices are there. On the other hand, it follows from the proof below, see Lemma 2.4 and the proof of Lemma 2.5, that most of the random fluctuations of $X_v(T_n)$ or $X_e(T_n)$ can be explained by the values at heights close to $\lg \lg n$. The explanation is that a few values λ_v at these heights will be so small that they significantly reduce the number of records among their descendants. Vertices of smaller height are too few, and there will usually not be any sufficiently small value among them, while vertices of larger height affect only a small proportion of the tree each, and the random effect caused by their values will be wiped out by the law of large numbers.

Remark 1.4. It is easy to see [5] that $\mathbb{E} X_v(T_n) = \sum_v 1/(h(v) + 1)$ and $\mathbb{E} X_e(T_n) = \sum_{v \neq o} 1/h(v)$; both sums are easily evaluated as $n/m + O(n/m^2) =$

$n/\lg n + O(n/\lg^2 n)$. We see thus from Theorem 1.1 that $X_v(T_n)$ and $X_e(T_n)$ are concentrated well above their means (at a distance of about $n \lg \lg n / \lg^2 n$), so that e.g. $\mathbb{P}(X_v(T_n) \leq \mathbb{E} X_v(T_n)) \rightarrow 0$. This is connected to the fact that the limit W_γ has infinite mean.

Note also that $\mathbb{E} X_e(T_n) - \mathbb{E} X_v(T_n) = \sum_{v \neq o} (h(v)(h(v) + 1))^{-1} - 1 \sim n/\lg^2 n$, while there is no similar difference in the limit distribution in Theorem 1.1.

An explanation of these facts is that the mean is affected by the unlikely event that a vertex close to the root has an extremely small value λ_v , which would reduce the number of records by a large amount.

We see that this behaviour makes it impossible to use the method of moments to find the asymptotic distribution in Theorem 1.1, as we did for other trees in [5].

Remark 1.5. Recall that the Lévy measure $cx^{-2} dx$ gives a (weakly) 1-stable distribution, see e.g. [2, XVII.3]; the measure ν_γ is a version of this with periodic fluctuations, so the distribution of W_γ is roughly similar to a 1-stable distribution. More precisely, we have that if W_γ and W'_γ are independent with the same distribution, then $W_\gamma + W'_\gamma \stackrel{d}{=} 2W_\gamma + 2$, as is easily checked from (1.2), but the corresponding statement for a sum of three copies of W_γ is false.

If we write (1.2) as $\mathbb{E} e^{itW_\gamma} = e^{\psi_\gamma(t)}$, it is possible to compute the Fourier coefficients of $\psi_\gamma(t)$ as a function of γ by integrations, using Fubini and some Gamma integrals, and obtain

$$\begin{aligned} \psi_\gamma(t) &= -\frac{\pi}{2}|t| - (\gamma^* - \frac{1}{2})it \\ &\quad - it \sum_{n \neq 0} \frac{\Gamma(2\pi in / \ln 2 - 1)}{\ln 2 - 2\pi in} e^{-\pi^2 n \operatorname{sign} t / \ln 2} |t|^{-2\pi in / \ln 2} e^{2\pi ni\gamma}, \end{aligned}$$

where γ^* is Euler's constant. We omit the details. This, again, shows the affinity with stable distributions.

The complete binary tree T_n has minimal height among all binary trees with n vertices, but among binary trees with this height, it is maximally unbalanced. The other extreme is the *balanced binary tree* T_n^* , where at each vertex, the two subtrees emanating from it differ in size by at most 1. This tree too has height $m = \lceil \lg n \rceil$, and the same number of vertices at each level as T_n . As a companion to Theorem 1.1, we give a similar theorem for T_n^* ; note that the results are similar but not identical, which shows that the details of the structure of the tree are important. (If we consider only n of the form $2^k - 1$, $T_n = T$ is a *full binary tree*. Indeed, Theorems 1.1 and 1.6 yield the same result in this case.) In contrast, note that the means of X_v and X_e are the same for T_n and T_n^* , see Remark 1.4.

Theorem 1.6. *Suppose that $n \rightarrow \infty$ such that $\{\lg \lg n\} \rightarrow \beta \in [0, 1]$. Then*

$$\left(X_v(T_n^*) - \frac{n}{\lg n} - \frac{n \lg \lg n}{\lg^2 n} \right) / \frac{n}{\lg^2 n} \xrightarrow{d} -W_{1-\beta}$$

where $W_{1-\beta}$ is as in Theorem 1.1. The same result holds for $X_e(T_n^*)$.

The method used below applies also to other binary trees with minimal height, but we leave the details to the reader. Presumably, the method can be used also for a larger class of binary trees, but we have not explored this. In particular, we do not know whether our methods can be used to solve the following problem.

Problem 1.7. *What is the asymptotic distribution of X_v and X_e for a (random) binary search tree?*

2. PROOFS

We first treat the case X_v of Theorem 1.1 in detail, and then indicate the small modifications needed for X_e and for T_n^* .

Let $X_n := X_v(T_n)$, and let, for $y > 0$, $X_{n,y}$ be $X_v(T_n) - 1$ conditioned on the root label $\lambda_o = y$, i.e. the number of records in the rest of the tree if we fix the root label (which always is a record).

We will use the notations $m := \lfloor \lg n \rfloor$ (as above) and $l := \lfloor \lg \lg n \rfloor$; we also let $L := \lfloor \frac{3}{2} \lg \lg n \rfloor \approx 3l/2$. We assume that n is so large that $0 < l < L < m$. If a_n are positive numbers and Z_n random variables such that $Z_n/a_n \xrightarrow{p} 0$ as $n \rightarrow \infty$, we write $Z_n = o_p(a_n)$.

In the sequel, we will write T instead of T_n . For a vertex $v \in T$, we let T_v be the subtree of T rooted at v , and let n_v be the number of vertices in T_v .

For later use we note that if we fix $j < m$ and consider the 2^j vertices of height j , labelling them v_1, \dots, v_{2^j} from left to right, then, with $q_j := \lfloor (n - 2^m + 1)/2^{m-j} \rfloor$,

$$n_{v_i} = \begin{cases} 2^{m+1-j} - 1, & 1 \leq i \leq q_j, \\ 2^{m-j} - 1 + 2^{m-j} \{(n - 2^m + 1)/2^{m-j}\}, & i = q_j + 1, \\ 2^{m-j} - 1, & q_j + 1 < i \leq 2^j. \end{cases} \quad (2.1)$$

We will further assume that the labels λ_v have an exponential distribution $\text{Exp}(1)$ with mean 1; as remarked above, this does not affect the distribution of X_n .

Lemma 2.1. *We have*

$$\mathbb{E} X_{n,y} = \frac{n - 2^m + 1}{m} (1 - e^{-my}) + \sum_{k=1}^{m-1} \frac{2^{m-k}}{m-k} (1 - e^{-(m-k)y}) \quad (2.2)$$

and, uniformly in n and $y > 0$,

$$\text{Var} X_{n,y} = O(m^{-3}n^2).$$

Proof. Fix $y > 0$, and let, for each vertex $v \in T$, I_v be the indicator that λ_v is a minimum, given that $\lambda_o = y$. Thus, $X_{n,y} = \sum_{v \neq o} I_v$. If $h(v) = j$, let $o, v_1, \dots, v_j = v$ be the vertices on the path from the root o to v . Then $I_v = 1$ if and only if $\lambda_{v_j} < y$ and $\lambda_{v_i} > \lambda_{v_j}$ for $i = 1, \dots, j-1$. Hence, since $\lambda_{v_i} \sim \text{Exp}(1)$ are independent,

$$\mathbb{E} I_v = \int_0^y \prod_{i=1}^{j-1} \mathbb{P}(\lambda_{v_i} > x) e^{-x} dx = \int_0^y e^{-jx} dx = \frac{1 - e^{-jy}}{j}. \quad (2.3)$$

Consequently,

$$\mathbb{E} X_{n,y} = \sum_{j=1}^{m-1} 2^j \frac{1 - e^{-jy}}{j} + (n - 2^m + 1) \frac{1 - e^{-my}}{m},$$

proving (2.2) by letting $j = m - k$.

To estimate the variance, assume that v and w are two vertices in T of heights $j = h(v)$ and $k = h(w)$, and with their last common ancestor u at height i .

Suppose first $i < j$ and $i < k$. Let $u_0 = o, u_1, \dots, u_i = u$ be the vertices on the path from o to u , and let $Z := \min\{\lambda_{u_s} : 1 \leq s \leq i\}$. Conditioned on Z , I_v and I_w are independent. Further, since v has height $j - i$ above u , (2.3) yields

$$\mathbb{E}(I_v | Z) = \frac{1 - e^{-(j-i)(Z \wedge y)}}{j - i},$$

and similarly for I_w . Consequently, since $Z \sim \text{Exp}(i^{-1})$, being the minimum of i independent $\text{Exp}(1)$ variables,

$$\begin{aligned} \mathbb{E}(I_v I_w) &= \mathbb{E}\left(\frac{1 - e^{-(j-i)(Z \wedge y)}}{j - i} \cdot \frac{1 - e^{-(k-i)(Z \wedge y)}}{k - i}\right) \\ &= \frac{1}{j - i} \frac{1}{k - i} \left(\int_0^y (1 - e^{-(j-i)z})(1 - e^{-(k-i)z}) i e^{-iz} dz \right. \\ &\quad \left. + e^{-iy} (1 - e^{-(j-i)y})(1 - e^{-(k-i)y}) \right) \\ &= \frac{1}{j - i} \frac{1}{k - i} \left(1 - e^{-iy} - \frac{i}{j} (1 - e^{-jy}) - \frac{i}{k} (1 - e^{-ky}) \right. \\ &\quad \left. + \frac{i}{j + k - i} (1 - e^{-(j+k-i)y}) + e^{-iy} - e^{-jy} - e^{-ky} + e^{-(j+k-i)y} \right). \end{aligned} \quad (2.4)$$

Say that the pair (v, w) is *good* if $i \leq m/3$ and $j, k \geq 2m/3$, and *bad* otherwise. For a good pair (v, w) we have, by (2.4) and (2.3),

$$\begin{aligned} \text{Cov}(I_v, I_w) &= \mathbb{E} I_v I_w - \mathbb{E} I_v \mathbb{E} I_w \\ &= \frac{1 + O(i/m)}{jk} \left(1 - e^{-jy} - e^{-ky} + e^{-(j+k-i)y} + O(i/m) \right) \\ &\quad - \frac{1}{jk} (1 - e^{-jy})(1 - e^{-ky}) \\ &= \frac{1}{jk} e^{-(j+k-i)y} (1 - e^{-iy}) + O(i/m^3) \\ &= O(m^{-2} e^{-my} iy) + O(i/m^3) = O(i/m^3). \end{aligned} \quad (2.5)$$

For given i, j, k , there are at most 2^i choices of u and then at most 2^{j-i} choices of v and 2^{k-i} of w ; thus the total number of such pairs is at most 2^{j+k-i} . Hence (2.5) yields

$$\sum_{\text{good } (v,w)} \text{Cov}(I_v, I_w) = O\left(\sum_{i=1}^m \sum_{j=1}^m \sum_{k=1}^m 2^{j+k-i} i m^{-3} \right) = O(2^{2m} m^{-3}). \quad (2.6)$$

The total number of bad pairs is at most

$$\sum_{i > m/3, j, k \leq m} 2^{j+k-i} + 2 \sum_{i \geq 0, j < 2m/3, k \leq m} 2^{j+k-i} = O(2^{2m-m/3}). \quad (2.7)$$

For the bad pairs we simply use $\text{Cov}(I_v, I_w) \leq \mathbb{E} I_v I_w \leq 1$, and obtain from (2.6) and (2.7)

$$\text{Var } X_{n,y} = \sum_{v,w} \text{Cov}(I_v, I_w) = O(2^{2m} m^{-3}) = O(m^{-3} n^2). \quad \square$$

Let $\varphi(n, y) := \mathbb{E} X_{n,y}$, given by (2.2). In the next lemma, we find it useful to be slightly more general than simply requiring $\bar{m} = m$.

Lemma 2.2. *If $2^{\bar{m}} - 1 \leq n \leq 2^{\bar{m}+1} - 1$, then*

$$\varphi(n, y) = \frac{n}{\bar{m}} (1 - e^{-\bar{m}y}) + \frac{2^{\bar{m}+1}}{\bar{m}^2} + O(\bar{m}^{-2} e^{-\bar{m}y/4} n + \bar{m}^{-3} n).$$

Proof. Let

$$a_k = \frac{2^{m-k}}{m-k} (1 - e^{-(m-k)y}) = \frac{2^{m-k}}{m} \left(1 - \frac{k}{m}\right)^{-1} \left(1 - e^{-my} - e^{(k-m)y} (1 - e^{-ky})\right).$$

For $m/2 < k < m$ we use $a_k = O(2^{m/2})$, and for $k \leq m/2$

$$a_k = \frac{2^{m-k}}{m} \left(1 + \frac{k}{m} + O\left(\frac{k^2}{m^2}\right)\right) \left(1 - e^{-my} + O(ky e^{(k-m)y})\right).$$

Summing over k , we see that (2.2) yields, using $mye^{-my/4} = O(1)$,

$$\begin{aligned}\varphi(n, y) &= \frac{n - 2^m}{m} (1 - e^{-my}) + \frac{2^m}{m} \left(1 + \frac{2}{m} + O(m^{-2})\right) (1 - e^{-my}) \\ &\quad + O\left(m2^{m/2} + \frac{2^m}{m} ye^{-my/2}\right) \\ &= \frac{n}{m} (1 - e^{-my}) + \frac{2^{m+1}}{m^2} + O(m^{-3}n) + O(m^{-2}ne^{-my/4}).\end{aligned}$$

This proves the result when $\bar{m} = m$. The only remaining case is $n = 2^{\bar{m}} - 1$ and $m = \bar{m} - 1$; the result follows easily in this case too, for example by adding a vertex v at height \bar{m} , using the case just considered, and subtracting $\mathbb{E}I_v = (1 - e^{-\bar{m}y})/\bar{m}$ from (2.3). \square

Recall that $L = \lfloor \frac{3}{2} \lg \lg n \rfloor \approx \frac{3}{2}l \approx \frac{3}{2} \lg m$. Let v_i , $1 \leq i \leq 2^L$, be the 2^L vertices of height L , and let $n_i := n_{v_i}$. Note that $n_i = \Theta(n/2^L)$. Further, let Y_i be the minimum of λ_v along the path $P(v_i) = o \dots v_i$ from the root to v_i .

Lemma 2.3. *With notations as above,*

$$X_n = \sum_{i=1}^{2^L} \varphi(n_i, Y_i) + o_p(m^{-2}n). \quad (2.8)$$

Proof. We write the number of records X_n as $V^* + V_1 + \dots + V_{2^L}$, where V^* is the number of records with height $\leq L$ and V_i are the number of records in $T_{v_i} \setminus \{v_i\}$.

If we condition on $\{\lambda_v : h(v) \leq L\}$, then V^* and all Y_i become fixed, while V_i , $1 \leq i \leq 2^L$, become independent random variables with $V_i \stackrel{d}{=} X_{n_i, Y_i}$.

Let \mathcal{F}_L be the σ -field generated by $\{\lambda_v : h(v) \leq L\}$. Then, by the comments just made, $\mathbb{E}(V_i | \mathcal{F}_L) = \mathbb{E}(X_{n_i, Y_i} | Y_i) = \varphi(n_i, Y_i)$ and, with $m_i := \lfloor \lg n_i \rfloor = \lg n - L + O(1) \sim m$,

$$\begin{aligned}\mathbb{E}\left(\left(X_n - V^* - \sum_{i=1}^{2^L} \varphi(n_i, Y_i)\right)^2 \mid \mathcal{F}_L\right) &= \mathbb{E}\left(\left(\sum_{i=1}^{2^L} (V_i - \varphi(n_i, Y_i))\right)^2 \mid \mathcal{F}_L\right) \\ &= \sum_{i=1}^{2^L} \mathbb{E}\left((V_i - \varphi(n_i, Y_i))^2 \mid \mathcal{F}_L\right) = \sum_{i=1}^{2^L} \text{Var}(X_{n_i, Y_i} \mid Y_i) \\ &= \sum_{i=1}^{2^L} O(m_i^{-3}n_i^2) = O(2^L m^{-3} 2^{-2L} n^2) = O(m^{-3} 2^{-L} n^2) \\ &= O(m^{-9/2} n^2).\end{aligned}$$

Taking the expectation, we find

$$\mathbb{E}\left(X_n - V^* - \sum_{i=1}^{2^L} \varphi(n_i, Y_i)\right)^2 = O(m^{-9/2} n^2) = o(m^{-4} n^2)$$

and thus

$$X_n - V^* - \sum_{i=1}^{2^L} \varphi(n_i, Y_i) = o_p(m^{-2}n).$$

The result follows because also

$$0 \leq V^* < 2^{L+1} = O(m^{3/2}) = o(m^{-2}n). \quad \square$$

Next, let $\bar{m} := m - L \sim m$. By (2.1) (with $j = L$), we can apply Lemma 2.2 to each n_i and \bar{m} ; this yields

$$\varphi(n_i, Y_i) = \frac{n_i}{\bar{m}}(1 - e^{-\bar{m}Y_i}) + \frac{2^{\bar{m}+1}}{\bar{m}^2} + O(m^{-2}e^{-\bar{m}Y_i/4}n_i + m^{-3}n_i). \quad (2.9)$$

Since $Y_i \sim \text{Exp}(1/(L+1))$, for every $a > 0$,

$$\mathbb{E} e^{-aY_i} = \int_0^\infty (L+1)e^{-ay-(L+1)y} dy = \frac{L+1}{L+1+a} = O\left(\frac{L}{a}\right). \quad (2.10)$$

Hence

$$\begin{aligned} \mathbb{E}|e^{-\bar{m}Y_i} - e^{-mY_i}| &= \mathbb{E} e^{-\bar{m}Y_i} - \mathbb{E} e^{-mY_i} = \frac{L+1}{L+1+\bar{m}} - \frac{L+1}{L+1+m} \\ &= O\left(L\frac{m-\bar{m}}{m^2}\right) = O\left(\frac{L^2}{m^2}\right). \end{aligned} \quad (2.11)$$

It follows easily from (2.9), (2.10) and (2.11) that

$$\mathbb{E}\left|\varphi(n_i, Y_i) - \frac{n_i}{\bar{m}} + \frac{n_i}{m}e^{-mY_i} - \frac{2^{\bar{m}+1}}{\bar{m}^2}\right| = O(L^2m^{-3}n_i) = o(m^{-2}n_i).$$

Summing over i we find, using Lemma 2.3, since $\sum_i n_i = n - (2^L - 1) = n - O(m^{3/2})$,

$$\begin{aligned} X_n &= \sum_{i=1}^{2^L} \left(\frac{n_i}{\bar{m}} - \frac{n_i}{m}e^{-mY_i} + \frac{2^{\bar{m}+1}}{\bar{m}^2} \right) + o_p(m^{-2}n) \\ &= \frac{n}{m-L} - \frac{1}{m} \sum_{i=1}^{2^L} n_i e^{-mY_i} + 2^L \frac{2^{m-L+1}}{(m-L)^2} + o_p(m^{-2}n) \\ &= \frac{n}{m} + L\frac{n}{m^2} - \frac{1}{m} \sum_{i=1}^{2^L} n_i e^{-mY_i} + \frac{2^{m+1}}{m^2} + o_p(m^{-2}n). \end{aligned} \quad (2.12)$$

We transform this once more.

Lemma 2.4.

$$X_n = \frac{n}{m} + L\frac{n}{m^2} - \frac{1}{m} \sum_{h(v) \leq L} n_v e^{-m\lambda_v} + \frac{2^{m+1}}{m^2} + o_p(m^{-2}n). \quad (2.13)$$

Proof. We recall that each Y_i is the minimum of the $L + 1$ independent variables λ_v , $v \in P(v_i)$; thus e^{-mY_i} is the maximum of the corresponding $e^{-m\lambda_v}$. Let $a = 2 \ln m/m$. The probability that at least two λ_v , $v \in P(v_i)$, are less than a is $O(L^2 a^2) = O(\ln^4 m/m^2)$; hence the probability that this happens for some i is $O(2^L \ln^4 m/m^2) = o(1)$. With probability tending to 1, there is thus at most one λ_v less than a in each $P(v_i)$, and in this case,

$$0 \leq \sum_{v \in P(v_i)} e^{-m\lambda_v} - e^{-mY_i} \leq L e^{-ma} = L/m^2,$$

and thus,

$$\begin{aligned} \sum_{i=1}^{2^L} n_i e^{-mY_i} &= \sum_{i=1}^{2^L} n_i \sum_{v \in P(v_i)} e^{-m\lambda_v} + O(nL/m^2) \\ &= \sum_{h(v) \leq L} e^{-m\lambda_v} \sum_{i: v \in P(v_i)} n_i + O(nL/m^2) \\ &= \sum_{h(v) \leq L} e^{-m\lambda_v} n_v + O(nL/m^2), \end{aligned}$$

because $n_v - 2^L \leq \sum_{i: v \in P(v_i)} n_i \leq n_v$. Hence,

$$\sum_{i=1}^{2^L} n_i e^{-mY_i} = \sum_{h(v) \leq L} e^{-m\lambda_v} n_v + o_p(n/m),$$

and the result follows from (2.12). \square

The sum in (2.13) is a sum of independent random variables. The proof will be completed by a classical result on convergence of such sums for triangular arrays to infinitely divisible distributions, see e.g. [3, Theorem 15.28].

We write, for convenience, $\xi_v := \frac{mn_v}{n} e^{-m\lambda_v}$. We further write $\alpha_n := \{\lg n\}$ and $\beta_n = \{\lg \lg n\}$; Thus $\lg n = m + \alpha_n$ and $\lg m = \lg \lg n + o(1) = l + \beta + o(1)$. We then have, by Lemma 2.4,

$$\begin{aligned} &\frac{m^2}{n} \left(X_n - \frac{n}{\lg n} - \frac{n \lg \lg n}{\lg^2 n} \right) \\ &= m^2 \left(\frac{1}{m} - \frac{1}{\lg n} \right) + L - m \sum_{h(v) \leq L} \frac{n_v}{n} e^{-m\lambda_v} + \frac{2^{m+1}}{n} - \lg \lg n + o_p(1) \\ &= \alpha_n + L - l - \beta_n + 2^{1-\alpha_n} - \sum_{h(v) \leq L} \xi_v + o_p(1). \end{aligned} \tag{2.14}$$

Since $m/\lg n \rightarrow 1$, it is thus enough to show that this converges in distribution to $-W_\gamma$ as $n \rightarrow \infty$ with $\{\lg n - \lg \lg n\} \rightarrow \gamma$.

By considering subsequences, we may assume that the limits $\alpha := \lim \alpha_n$ and $\beta := \lim \beta_n$ exist. Thus $\lg n = m + \alpha + o(1)$ and $\lg m = \lg \lg n + o(1) =$

$l + \beta + o(1)$. Note that $\lg n - \lg \lg n = m - l + \alpha - \beta + o(1)$; thus $\gamma \equiv \alpha - \beta \pmod{1}$ and more precisely,

$$\gamma = \begin{cases} \alpha - \beta & \text{if } \alpha > \beta; \\ \alpha - \beta + 1 & \text{if } \alpha < \beta; \\ 0 \text{ or } 1 & \text{if } \alpha = \beta. \end{cases} \quad (2.15)$$

Lemma 2.5. *Suppose that $n \rightarrow \infty$ such that $\alpha_n \rightarrow \alpha$ and $\beta_n \rightarrow \beta$ for some α and β in $[0, 1]$, and let $h := 2^{\beta - \alpha}$. Then*

- (i) $\sup_v \mathbb{P}(\xi_v > x) \rightarrow 0$ for every $x > 0$. (I.e., $\{\xi_v\}$ form a null array.)
- (ii) $\sum_{h(v) \leq L} \mathbb{P}(\xi_v > x) \rightarrow \nu_\gamma(x, \infty)$ for every $x > 0$.
- (iii) $\sum_{h(v) \leq L} \mathbb{E}(\xi_v \mathbf{1}[\xi_v \leq h]) - (L - l + 2^{1-\alpha} + \alpha - \beta) \rightarrow \beta - \alpha$.
- (iv) $\sum_{h(v) \leq L} \text{Var}(\xi_v \mathbf{1}[\xi_v \leq h]) \rightarrow 3h/2$.

Before proving this lemma, we show how it implies Theorem 1.1. Let $C := L - l + 2^{1-\alpha} + \alpha - \beta$. We apply [3, Theorem 15.28] with $a = 0$ and $b = f(\gamma)$ to $\sum_{h(v) \leq L} \xi_v + \sum_{i=1}^n \xi'_i$, with $\xi'_i = -C/n$ deterministic. (Note that $C/n \rightarrow 0$; thus $\{\xi_v\} \cup \{\xi'_i\}$ is a null array.) We have $d\nu_\gamma/dx = 2^{\lfloor \lg x + \alpha - \beta \rfloor} x^{-2} = 2^{-i+\alpha-\beta} x^{-1}$ when $2^i h < x < 2^{i+1} h$, and thus

$$\int_0^h x^2 d\nu_\gamma(x) = \sum_{i=-\infty}^{-1} \int_{2^i h}^{2^{i+1} h} 2^{-i+\alpha-\beta} x dx = \sum_{i=-\infty}^{-1} \frac{3}{2} 2^i h = \frac{3h}{2}.$$

Similarly, if $\beta \leq \alpha$ so $1/2 \leq h \leq 1$, then

$$\int_h^1 x d\nu_\gamma(x) = \int_h^1 2^{\alpha-\beta} dx = 2^{\alpha-\beta} - 1,$$

while if $\beta \geq \alpha$ so $1 \leq h \leq 2$, then

$$\int_h^1 x d\nu_\gamma(x) = - \int_1^h x d\nu_\gamma(x) = 2^{1+\alpha-\beta} (1 - h) = 2(2^{\alpha-\beta} - 1).$$

It follows, using (2.15) and $f(0) = f(1)$, that in both cases

$$f(\gamma) - \int_h^1 x d\nu_\gamma(x) = 2^\gamma - 1 - \gamma - \int_h^1 x d\nu_\gamma(x) = \beta - \alpha.$$

It is now easy to see from Lemma 2.5 that the conditions of [3, Theorem 15.28] are satisfied, and consequently

$$\sum_{h(v) \leq L} \xi_v - (L - l + 2^{1-\alpha} + \alpha - \beta) = \sum_{h(v) \leq L} \xi_v + \sum_{i=1}^n \xi'_i \xrightarrow{d} W_\gamma.$$

Theorem 1.1 now follows by (2.14).

Proof of Lemma 2.5. For any $x > 0$,

$$\begin{aligned} \mathbb{P}(\xi_v > x) &= \mathbb{P}\left(e^{-m\lambda_v} > \frac{nx}{mn_v}\right) = \mathbb{P}\left(m\lambda_v < \ln \frac{mn_v}{nx}\right) \\ &= 1 - \exp\left(-\frac{1}{m} \ln_+ \frac{mn_v}{nx}\right). \end{aligned} \quad (2.16)$$

This shows first that for every $x > 0$,

$$\mathbb{P}(\xi_v > x) < \frac{1}{m} \ln_+ \frac{mn_v}{nx} \leq \frac{1}{m} \ln_+ \frac{m}{x} \rightarrow 0, \quad (2.17)$$

which proves (i).

On a given level $j < m$ there are, by (2.1), $q_j = n2^{j-m} - 2^j + O(1) = (2^{\alpha_n} - 1)2^j + O(1)$ vertices with $n_v = 2^{m+1-j} - 1$, and $2^j - q_j - 1 = (2 - 2^{\alpha_n})2^j + O(1)$ vertices with $n_v = 2^{m-j} - 1$. There is one additional vertex with an intermediate n_v (which could coincide with one of the two main values); for convenience we call such a vertex *bad*. We also call a vertex v with $n_v \geq 2^{m-l/2}$ (which requires $j \leq l/2$) *bad*. All other vertices v with $h(v) \leq L$ are *good*. The good vertices thus have $n_v = 2^{m-k} - 1$ for some k with $l/2 \leq k \leq L$. For $l/2 \leq k < L$, there are $(2 - 2^{\alpha_n})2^k + O(1)$ such vertices with $h(v) = k$ and $(2^{\alpha_n} - 1)2^{k+1} + O(1)$ with $h(v) = k + 1$; thus together $2^{k+\alpha_n} + O(1)$. For $k = L$, there are only $(2 - 2^{\alpha_n})2^k + O(1)$ such vertices, since we require $h(v) \leq L$. In other words,

$$\#\{v \text{ good} : n_v = 2^{m-k} - 1\} = \begin{cases} 2^{k+\alpha_n} + O(1), & l/2 \leq k < L, \\ (2 - 2^{\alpha_n})2^L + O(1), & k = L. \end{cases} \quad (2.18)$$

The number of bad vertices is $O(L + 2^{l/2}) = O(m^{1/2})$. By (2.17), $\mathbb{P}(\xi_v > x) = O(\ln m/m)$ for every fixed $x > 0$. Hence the sum over bad vertices in (ii) is $O(m^{-1/2} \ln m) = o(1)$.

Similarly, using (2.17) again,

$$\mathbb{E}(\xi_v \mathbf{1}[\xi_v \leq h]) \leq \frac{1}{m} + h \mathbb{P}\left(\xi_v > \frac{1}{m}\right) \leq \frac{1}{m} + h \frac{2 \ln m}{m} = O\left(\frac{\ln m}{m}\right). \quad (2.19)$$

and

$$\text{Var}(\xi_v \mathbf{1}[\xi_v \leq h]) \leq \mathbb{E}(\xi_v^2 \mathbf{1}[\xi_v \leq h]) \leq h \mathbb{E}(\xi_v \mathbf{1}[\xi_v \leq h]) = O\left(\frac{\ln m}{m}\right). \quad (2.20)$$

Consequently, the sum over bad v is $o(1)$ in (ii), (iii) and (iv), so we may in the sequel ignore them and consider only good vertices.

Fix $x > 0$. Then, by (2.16) and (2.17),

$$\mathbb{P}(\xi_v > x) = \frac{1}{m} \ln_+ \left(\frac{mn_v}{nx} \right) \left(1 + O\left(\frac{\ln m}{m} \right) \right) \quad (2.21)$$

If $k \geq L$, then $m2^{m-k} \leq 2^{l+1+m-L} < nx$, provided n is large enough. Thus, for large n , by (2.18) and (2.21), with all $o(1)$ uniform in k for fixed x ,

$$\begin{aligned} \sum_{v \text{ good}} \mathbb{P}(\xi_v > x) &= (1 + o(1)) \sum_{k=l/2}^{\infty} (2^{k+\alpha_n} + O(1)) \frac{1}{m} \ln_+ \left(\frac{(2^{m-k} - 1)m}{nx} \right) \\ &= (1 + o(1)) \sum_{k \geq l/2} 2^{k+\alpha_n-l-\beta_n} \ln_+ (2^{-k-\alpha_n+l+\beta_n+o(1)} x^{-1}) + o(1) \\ &= (1 + o(1)) \sum_{i \leq l/2} 2^{-i+\alpha-\beta} \ln_+ (2^{i-\alpha+\beta+o(1)} x^{-1}) + o(1) \\ &\rightarrow F(x) := \sum_{-\infty}^{\infty} 2^{-i+\alpha-\beta} \ln_+ (2^{i-\alpha+\beta} x^{-1}). \end{aligned}$$

Let $j := \lfloor \lg x + \alpha - \beta \rfloor$; thus $2^{j+\beta-\alpha} \leq x < 2^{k+\beta-\alpha+1}$ and

$$\begin{aligned} F(x) &= \sum_{i=j+1}^{\infty} 2^{-i+\alpha-\beta} \ln(2^{i-\alpha+\beta} x^{-1}) \\ &= \sum_{k=1}^{\infty} 2^{-k-j+\alpha-\beta} (k \ln 2 + \ln(2^{j-\alpha+\beta} x^{-1})) \\ &= 2^{-j+\alpha-\beta} (2 + \lg(2^{j-\alpha+\beta} x^{-1})) \ln 2 \\ &= 2^{\alpha-\beta-\lfloor \lg x + \alpha - \beta \rfloor} (2 - \{\lg x + \alpha - \beta\}) \ln 2 \\ &= 2^{\gamma - \lfloor \lg x + \gamma \rfloor} (2 - \{\lg x + \gamma\}) \ln 2. \end{aligned}$$

Note that $F(x)$ is continuous and decreasing with $F(x) \rightarrow 0$ as $x \rightarrow \infty$. The derivative is

$$\frac{dF(x)}{dx} = -\frac{1}{x} 2^{\gamma - \lfloor \lg x + \gamma \rfloor} = -x^{-2} 2^{\{\lg x + \gamma\}}.$$

Thus $F(x) = \nu_{\gamma}(x, \infty)$, which proves (ii).

For (iii) and (iv) we calculate, for $s > 0$,

$$\begin{aligned} \mathbb{E}(e^{-m\lambda_v} \mathbf{1}[e^{-m\lambda_v} \leq s^{-1}]) &= \int_{m^{-1} \ln_+ s}^{\infty} e^{-mx} e^{-x} dx \\ &= \frac{1}{m+1} e^{-(m+1)\frac{1}{m} \ln_+ s} = \frac{1}{m+1} 2^{-(1+1/m)\lg_+ s} \end{aligned} \quad (2.22)$$

and, similarly,

$$\mathbb{E}(e^{-2m\lambda_v} \mathbf{1}[e^{-m\lambda_v} \leq s^{-1}]) = \frac{1}{2m+1} 2^{-(2+1/m)\lg_+ s},$$

which gives

$$\text{Var}(e^{-m\lambda_v} \mathbf{1}[e^{-m\lambda_v} \leq s^{-1}]) = \frac{1}{2m} 2^{-(2+1/m)\lg_+ s} (1 + O(m^{-1})). \quad (2.23)$$

If v is a good vertex with $n_v = 2^{m-k} - 1 = 2^{m-k+o(1)}$, then

$$\frac{mn_v}{nh} = 2^{(l+\beta)+(m-k)-(m+\alpha)-(\beta-\alpha)+o(1)} = 2^{l-k+o(1)},$$

and thus, by (2.22),

$$\begin{aligned}\mathbb{E}(\xi_v \mathbf{1}[\xi_v \leq h]) &= \frac{mn_v}{n} \mathbb{E}\left(e^{-m\lambda_v} \mathbf{1}\left[e^{-m\lambda_v} \leq \frac{nh}{mn_v}\right]\right) \\ &= \frac{mn_v}{(m+1)n} 2^{-(1+1/m)(l-k)_++o(1)} = 2^{-k-\alpha-(l-k)_++o(1)}.\end{aligned}$$

Note in particular that if $k \geq l+1$, then $\frac{mn_v}{nh} < 1$ for large n , and thus $\xi_v \leq h$ and

$$\mathbb{E}(\xi_v \mathbf{1}[\xi_v \leq h]) = \mathbb{E} \xi_v = \frac{mn_v}{(m+1)n} = 2^{-k-\alpha_n} \left(1 + O\left(\frac{1}{m}\right)\right). \quad (2.24)$$

It follows from (2.18), (2.22) and (2.24) that, with o and O uniform in k ,

$$\begin{aligned}\sum_{v \text{ good}} \mathbb{E}(\xi_v \mathbf{1}[\xi_v \leq h]) &= \sum_{k=l/2}^l 2^{k+\alpha+o(1)} 2^{-k-\alpha-(l-k)_++o(1)} + \sum_{k=l+1}^{L-1} (2^{k+\alpha_n} + O(1)) 2^{-k-\alpha_n} (1 + O(m^{-1})) \\ &\quad + ((2 - 2^{\alpha_n}) 2^L + O(1)) 2^{-L-\alpha+o(1)} \\ &= \sum_{k=l/2}^l 2^{-(l-k)_++o(1)} + \sum_{k=l+1}^{L-1} (1 + O(m^{-1})) + 2^{1-\alpha} - 1 + o(1) \\ &= 2 + L - 1 - l + 2^{1-\alpha} - 1 + o(1) = L - l + 2^{1-\alpha} + o(1).\end{aligned}$$

Similarly, using (2.23),

$$\begin{aligned}\text{Var}(\xi_v \mathbf{1}[\xi_v \leq h]) &= \frac{m^2 n_v^2}{n^2} \text{Var}\left(e^{-m\lambda_v} \mathbf{1}\left[e^{-m\lambda_v} \leq \frac{nh}{mn_v}\right]\right) \\ &= \frac{m^2 n_v^2}{2mn^2} 2^{-(2+1/m)(l-k)_++o(1)} = 2^{l+\beta-1-2k-2\alpha-2(l-k)_++o(1)}\end{aligned}$$

and

$$\begin{aligned}\sum_{v \text{ good}} \text{Var}(\xi_v \mathbf{1}[\xi_v \leq h]) &= \sum_{k=l/2}^L 2^{k+\alpha+o(1)} 2^{l+\beta-1-2k-2\alpha-2(l-k)_++o(1)} + o(1) \\ &= \sum_{k=-\infty}^{\infty} 2^{l-k-2(l-k)_++\beta-\alpha-1+o(1)} + o(1) \\ &= 3 \cdot 2^{\beta-\alpha-1} + o(1) = 3h/2 + o(1).\end{aligned}$$

This completes the proof of Lemma 2.5. \square

We have proved Theorem 1.1 for X_v . For X_e , the only difference is that λ_o is ignored, and thus $Y_i \sim \text{Exp}(1/L)$. The estimates in (2.10) and (2.11) remain valid, and thus (2.13) and (2.14) still hold, summing over $v \neq o$ only. Since $\xi_0 = me^{-m\lambda_o} \xrightarrow{P} 0$ by Lemma 2.5(i), this makes no difference for

the asymptotics of the distribution. (But note that $\mathbb{E} \xi_0 \rightarrow 1$, and that the means differ correspondingly, see Remark 1.4.)

For the completely balanced tree (Theorem 1.6), every vertex v with $h(v) = k$ has $2^{-k}n - 2 < n_v \leq 2^{-k}n$. We call all vertices with $l/2 \leq h(v) \leq L$ good, and replace (2.18) by

$$\#\{v \in T_n^* \text{ good} : n_v = 2^{-k}n + O(1)\} = 2^k, \quad l \leq k \leq L. \quad (2.25)$$

The remaining calculations hold as above, provided we replace α_n and α by 0 and thus γ by $1 - \beta$.

REFERENCES

- [1] P. Chassaing & R. Marchand. In preparation.
- [2] W. Feller, *An Introduction to Probability Theory and Its Applications. Vol. II*. Second edition, Wiley, New York 1971.
- [3] O. Kallenberg, *Foundations of Modern Probability*. 2nd ed., Springer-Verlag, New York, 2002.
- [4] D.E. Knuth, *The Art of Computer Programming. Vol. 1: Fundamental Algorithms*. 3rd ed., Addison-Wesley, Reading, Mass., 1997.
- [5] S. Janson, Random cutting and records in deterministic and random trees. Preprint, 2003. Available from <http://www.math.uu.se/~svante/papers>
- [6] A. Meir & J.W. Moon, Cutting down random trees. *J. Australian Math. Soc.* **11** (1970), 313–324.
- [7] A. Panholzer, Cutting down very simple trees. Preprint, 2003.
- [8] A. Panholzer, Non-crossing trees revisited: cutting down and spanning subtrees. *Proceedings, Discrete Random Walks 2003*, Cyril Banderier and Christian Krattenthaler, Eds., *Discr. Math. Theor. Comput. Sci.* **AC** (2003), 265–276.
- [9] A. Rényi, (1962). On the extreme elements of observations. *MTA III, Oszk. Közl.* **12** (1962) 105–121. Reprinted in *Collected Works*, Vol III, pp. 50-66, Akadémiai Kiadó, Budapest, 1976.

DEPARTMENT OF MATHEMATICS, UPPSALA UNIVERSITY, PO BOX 480, S-751 06 UPPSALA, SWEDEN

E-mail address: svante.janson@math.uu.se

URL: <http://www.math.uu.se/~svante/>