



**HAL**  
open science

## The Moran forest

François Bienvenu, Jean-Jil Duchamps, Félix Foutel-Rodier

► **To cite this version:**

François Bienvenu, Jean-Jil Duchamps, Félix Foutel-Rodier. The Moran forest. *Random Structures and Algorithms*, 2021, 59 (2), pp.155-188. 10.1002/rsa.20997 . hal-02165001v3

**HAL Id: hal-02165001**

**<https://hal.sorbonne-universite.fr/hal-02165001v3>**

Submitted on 16 Feb 2024

**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 Moran forest

François Bienvenu<sup>1,2</sup>, Jean-Jil Duchamps<sup>3</sup>, and Félix  
Foutel-Rodier<sup>1,2</sup>

<sup>1</sup>*Center for Interdisciplinary Research in Biology (CIRB), CNRS UMR 7241,  
Collège de France, PSL Research University, Paris, France*

<sup>2</sup>*Laboratoire de Probabilités, Statistique et Modélisation (LPSM), CNRS UMR 8001,  
Sorbonne Université, Paris, France*

<sup>3</sup>*Laboratoire de mathématiques de Besançon (LmB) UMR 6623, Université Bourgogne  
Franche-Comté, CNRS, F-25000 Besançon, France*

December 16, 2020

## Abstract

Starting from any graph on  $\{1, \dots, n\}$ , consider the Markov chain where at each time-step a uniformly chosen vertex is disconnected from all of its neighbors and reconnected to another uniformly chosen vertex. This Markov chain has a stationary distribution whose support is the set of non-empty forests on  $\{1, \dots, n\}$ . The random forest corresponding to this stationary distribution has interesting connections with the uniform rooted labeled tree and the uniform attachment tree. We fully characterize its degree distribution, the distribution of its number of trees, and the limit distribution of the size of a tree sampled uniformly. We also show that the size of the largest tree is asymptotically  $\alpha \log n$ , where  $\alpha = (1 - \log(e - 1))^{-1} \approx 2.18$ , and that the degree of the most connected vertex is asymptotically  $\log n / \log \log n$ .

## Contents

<b>1</b>	<b>Introduction</b>	<b>2</b>
1.1	The model . . . . .	2
1.2	Main results . . . . .	3
<b>2</b>	<b>Sampling of the stationary distribution</b>	<b>4</b>
2.1	Backward construction . . . . .	4
2.2	Uniform attachment construction . . . . .	5
<b>3</b>	<b>Number of trees</b>	<b>8</b>
3.1	Law of the number of trees . . . . .	8
3.2	Link with uniform labeled trees . . . . .	9
<b>4</b>	<b>Degrees</b>	<b>11</b>
4.1	Degree of a fixed vertex . . . . .	11
4.2	Largest degree . . . . .	14

<b>5</b>	<b>Tree sizes</b>	<b>18</b>
5.1	A discrete-time Yule process . . . . .	18
5.2	Size of some random trees . . . . .	22
5.3	Size of the largest tree . . . . .	24
<b>6</b>	<b>Concluding comments</b>	<b>26</b>
6.1	Aldous's construction . . . . .	26
6.2	A local limit . . . . .	26
6.3	Possible extension . . . . .	27
	<b>References</b>	<b>28</b>
<b>A</b>	<b>Appendix</b>	<b>30</b>
A.1	Proof of point (ii) of Proposition 4.4 . . . . .	30
A.2	Technical lemmas used in the proof of Theorem 1.3 . . . . .	31

# 1 Introduction

## 1.1 The model

Consider a Markov chain on the space of directed graphs on  $\{1, \dots, n\}$ , for a fixed  $n \geq 2$ , whose transition probabilities are defined as follows: at each time-step,

1. Choose an ordered pair of distinct vertices  $(u, v)$  uniformly at random.
2. Disconnect  $v$  from all of its neighbors, then add the edge  $u\vec{v}$ .

Note that if  $u\vec{v}$  is already the only edge attached to  $v$  at time  $t$ , then the graph is unchanged at time  $t + 1$ . A simple example illustrating the dynamics of this Markov chain is depicted in Figure 1.

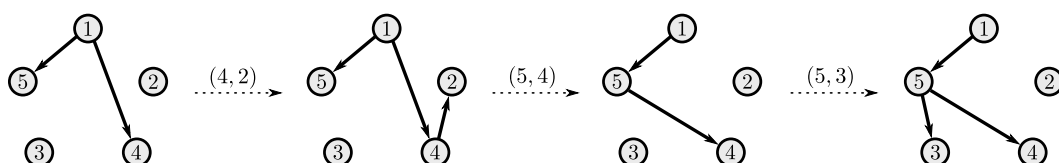


Figure 1: Example of four successive transitions of the Markov chain. Starting from the left-most graph, transitions are represented by dashed arrows decorated with the pair  $(u, v)$  that is chosen uniformly at each step.

This Markov chain has a stationary distribution whose support is the set of non-empty rooted forests on  $\{1, \dots, n\}$ . By *rooted forest* we mean a disjoint union of rooted trees – or, equivalently, a directed graph where each vertex has at most one incoming edge; any vertex  $\rho$  with no incoming edge can then be seen as the root of a tree consisting of all vertices accessible from  $\rho$ . To see why the stationary graph is a rooted forest, note that:

- The graph cannot be empty because there is always an edge between the two vertices involved in the last transition.

- Starting from any graph, the chain will eventually reach a forest (for instance, the sequence of transitions  $(1, 2), (1, 3), \dots, (1, n)$  will at some point turn the graph into the star graph rooted on vertex 1).
- The chain cannot leave the set of forests because its transitions cannot make a vertex have two incoming edges.
- Any non-empty rooted forest is accessible from any other graph (if not clear, this will become apparent in Section 2).
- The chain is aperiodic because it can stay in the same state.

The stationary distribution of this chain is the random forest model that we study in this paper. We call it the *Moran forest*, and use the notation  $\mathcal{F}_n$  to denote a random forest having this distribution. Our interest in this object lies in its connection with the Moran model of population genetics [20]. The Moran model describes the dynamics of a population of constant size  $n$  where, at each time step, two distinct individuals are sampled uniformly at random, and the second one is replaced by a copy of the first one. The Markov chain that we consider thus corresponds to the family structure of extant individuals in a Moran model. The Moran model is a central object in mathematical population genetics [10, 12], whose extensions have been used in a variety of other contexts, including diversification [15, 21] and evolutionary game theory [22].

## 1.2 Main results

Our first result, which we detail in Section 2, is that there is a simple way to sample  $\mathcal{F}_n$ . This construction enables us to study several of its statistics, such as its number of trees (Section 3.1), its degree distribution (Section 4.1), and the typical size of its trees (Section 5.2). Some of these results are presented in Table 1.

Notation	Variable	Distribution
$N_n$	Number of trees	$\sum_{\ell=1}^n I_\ell$ , where $I_\ell \sim \text{Ber}\left(\frac{\ell-1}{n-1}\right)$
$D$	Asymptotic degree distribution	$\text{Ber}(1 - U) + \text{Poisson}(U)$ , where $U \sim \text{Unif}([0, 1])$
$T^U$	Asymptotic size of a uniform tree	$\text{Geometric}(e^{-X})$ , where $X \sim 2xdx$ on $[0, 1]$

Table 1: Some statistics of the Moran forest, for fixed  $n$  in the case of the number of trees, and as  $n \rightarrow \infty$  for the degree and the size of a uniform tree. Note that the degree also has a simple, explicit distribution for fixed  $n$  (see Proposition 4.1); also, the Bernoulli and the Poisson r.v. appearing in the sum correspond respectively to the in- and out-degrees. The Bernoulli variables  $I_\ell$  used to describe the distribution of  $N_n$  are independent and, conditional on  $U$ , so are the Bernoulli and Poisson variables used for the distribution of  $D$ .

In Section 3.2, we show that the Moran forest is closely linked to the uniform rooted labeled tree. Specifically, we prove the following theorem.

**Theorem 1.1.** Let  $\mathcal{T}$  be a uniform rooted tree on  $\{1, \dots, n-1\}$ . From this tree, build a forest  $\mathcal{F}$  on  $\{1, \dots, n\}$  according to the following procedure:

1. Remove all decreasing edges from  $\mathcal{T}$  (that is, edges  $\vec{uv}$  pointing away from the root such that  $u > v$ ).
2. Add a vertex labeled  $n$  and connect it to a uniformly chosen vertex of  $\mathcal{T}$ .
3. Relabel vertices according to a uniform permutation of  $\{1, \dots, n\}$ .

Then, the resulting forest  $\mathcal{F}$  has the law of the Moran forest  $\mathcal{F}_n$ .

Finally, we study the asymptotic concentration of the largest degree and of the size of the largest tree of  $\mathcal{F}_n$ . The following theorems are proved in Sections 4.2 and 5.3, respectively.

**Theorem 1.2.** Let  $D_n^{\max}$  denote the largest degree of  $\mathcal{F}_n$ . Then,

$$D_n^{\max} = \frac{\log n}{\log \log n} + (1 + o_p(1)) \frac{\log n \log \log \log n}{(\log \log n)^2},$$

where  $o_p(1)$  denotes a sequence of random variables that goes to 0 in probability.

**Theorem 1.3.** Let  $T_n^{\max}$  denote the size of the largest tree of  $\mathcal{F}_n$ . Then,

$$T_n^{\max} = \alpha (\log n - (1 + o_p(1)) \log \log n),$$

where  $\alpha = (1 - \log(e-1))^{-1} \approx 2.18019$ .

## 2 Sampling of the stationary distribution

### 2.1 Backward construction

Consider an i.i.d. sequence  $((V_t, W_t), t \in \mathbb{Z})$ , where  $(V_t, W_t)$  is uniformly distributed on the set of ordered pairs of distinct elements of  $\{1, \dots, n\}$ . These variables are meant to encode the transitions of the chain:  $W_t$  represents the vertex that is disconnected at step  $t$ , and  $V_t$  the vertex to which  $W_t$  is then connected. We now explain how to construct a chain  $(\mathcal{F}_n(t), t \in \mathbb{Z})$  of forests by looking at the sequence  $((V_t, W_t), t \in \mathbb{Z})$  backwards in time.

Fix a focal time  $t \in \mathbb{Z}$ . For each vertex  $w$ , let us denote by

$$\tau_t(w) := \max\{s \leq t : W_s = w\}$$

the last time before  $t$  that  $w$  was chosen to be disconnected, and define

$$m_t(w) := V_{\tau_t(w)}$$

to be the vertex to which it was then reconnected. We refer to the time  $\tau_t(w)$  as the *birth time* of  $w$ , and to the vertex  $m_t(w)$  as its *mother*. Note that the variables  $(\tau_t(w), 1 \leq w \leq n)$  are independent of  $(m_t(w), 1 \leq w \leq n)$ .

Now, for each  $s \leq t$ , let the vertices be in one of two states, *active* or *inactive*, as follows: vertex  $w$  is active at times  $s$  such that  $\tau_t(w) \leq s \leq t$ , and inactive at times

$s < \tau_t(w)$ . Finally, let  $\mathcal{F}_n(t)$  be the forest obtained by connecting each vertex  $w$  to its mother if the mother is active at the time of birth of  $w$ , that is,

$$\text{there is an edge from } m_t(w) \text{ to } w \iff \tau_t(m_t(w)) < \tau_t(w).$$

This procedure is illustrated in Figure 2.

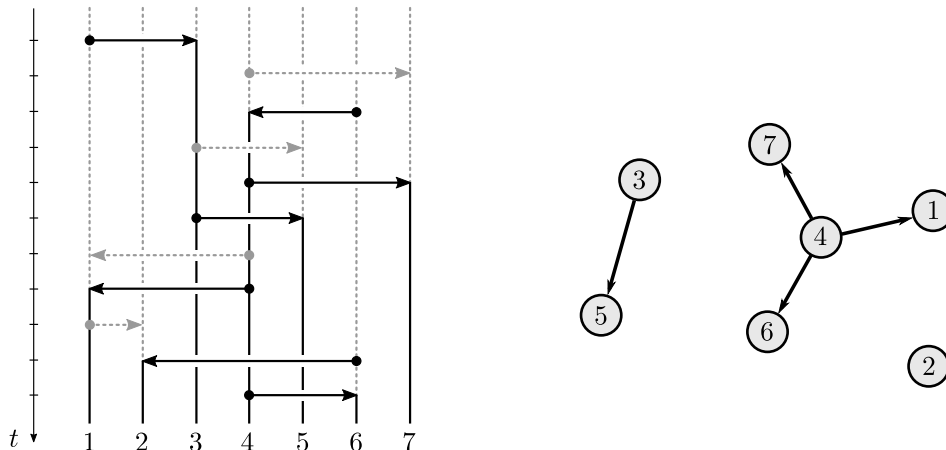


Figure 2: Illustration of the backward construction. Each vertex corresponds to a vertical line. A pair  $(V_t, W_t)$  is represented by an arrow  $V_t \rightarrow W_t$ . The line representing a vertex is solid black when that vertex is active, and dashed grey when it is inactive. Arrows pointing to inactive vertices are represented in dashed grey because they have no impact on the state of the graph at the focal time: their effect has been erased by subsequent arrows.

Let us show that the chain  $(\mathcal{F}_n(t), t \in \mathbb{Z})$  has the same transitions as the chain described in the introduction. First, note that for  $w \neq W_t$  we have  $\tau_t(w) = \tau_{t-1}(w)$ , and thus  $m_t(w) = m_{t-1}(w)$ . As a result, edges that do not involve  $W_t$  are the same in  $\mathcal{F}_n(t)$  and in  $\mathcal{F}_n(t-1)$ . Now,  $\tau_t(W_t) = t$ , so that  $W_t$  is always inactive as a mother in the construction of  $\mathcal{F}_n(t)$ , and  $m_t(W_t) = V_t$  with  $\tau_t(V_t) < t$ , so that  $W_t$  is linked to  $V_t$  in  $\mathcal{F}_n(t)$ . In other words,  $\mathcal{F}_n(t)$  is obtained from  $\mathcal{F}_n(t-1)$  by disconnecting  $W_t$  from its neighbors, and then connecting it to  $V_t$ . This corresponds to the transitions of the chain described in the introduction.

Finally,  $(\mathcal{F}_n(t), t \in \mathbb{Z})$  is stationary by construction, and thus  $\mathcal{F}_n(t)$  is distributed as the Moran forest for all time  $t \in \mathbb{Z}$ .

## 2.2 Uniform attachment construction

We now give a forward-in-time variant of the construction described in the previous section. This forward-in-time procedure, which we call the *uniform attachment construction* (UA construction for short), is our main tool to study  $\mathcal{F}_n$  and will be used throughout the rest of the paper.

In the following, we fix  $n \geq 2$ , since the forest  $\mathcal{F}_n$  is not defined for  $n = 1$ . Let  $(U_n(\ell), 1 \leq \ell \leq n)$  be a vector of independent variables such that  $U_n(\ell)$  is uniformly distributed on  $\{1, \dots, n\} \setminus \{\ell\}$ . Consider the forest  $\mathcal{F}_n^*$  on  $\{1, \dots, n\}$  obtained by setting

$$\text{there is an edge from } k \text{ to } \ell, \text{ with } k < \ell \iff U_n(\ell) = k.$$

We will show that, after relabeling the vertices of  $\mathcal{F}_n^*$  according to a uniform permutation of  $\{1, \dots, n\}$ , we obtain the Moran forest. Before this let us make a few remarks.

First, it will be helpful to think of the construction of  $\mathcal{F}_n^*$  as a sequential process where, starting from a single vertex labeled 1, for  $\ell = 2, \dots, n$  we add a new vertex labeled  $\ell$  and connect it to  $U_n(\ell)$  if  $U_n(\ell) < \ell$ . See Figure 3. This will make the link with some well-known stochastic processes more intuitive. This also explains that we speak of the  $\ell$ -th vertex in the UA construction to refer to vertex  $\ell$  in  $\mathcal{F}_n^*$ .

Second, the edges of  $\mathcal{F}_n^*$  are by construction increasing, in the sense that every edge  $\vec{uv}$  in the graph is such that  $u < v$ .

Rooted trees that have only increasing edges are known as *recursive trees* [9], and forests of recursive trees have been called *recursive forests* [5]. Recursive trees have been studied extensively. In particular, the uniform attachment tree, which corresponds to the uniform distribution over the set of recursive trees, has received much attention [6, 18, 19]. However, the random forest  $\mathcal{F}_n^*$  does not seem to correspond to any previously studied model of random recursive forest (in particular, it is not uniformly distributed over the set of recursive forests).

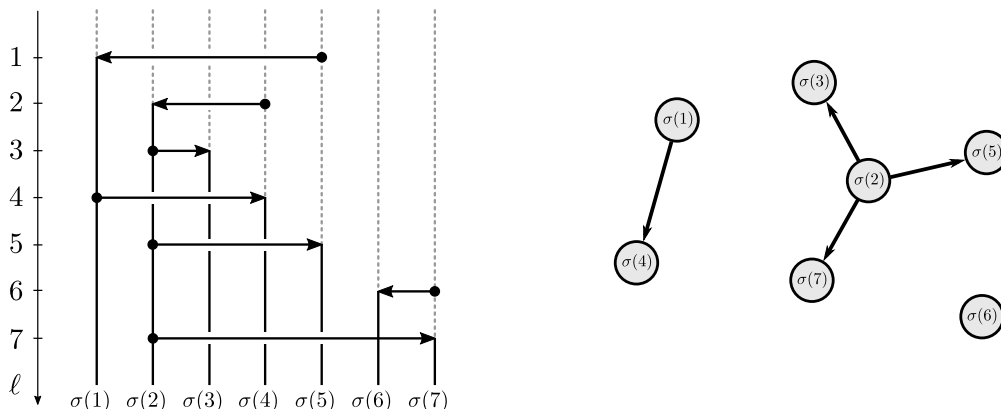


Figure 3: Illustration of the uniform attachment construction for  $n = 7$  and the vector  $(U_n(1), \dots, U_n(n)) = (5, 4, 2, 1, 2, 7, 2)$ . The  $\ell$ -th vertical line from the left corresponds to vertex  $\sigma(\ell)$  (i.e. in the sequential vision, to the  $\ell$ -th vertex that is added).  $U_n(\ell)$  is represented by the arrow pointing from the  $U_n(\ell)$ -th line to the  $\ell$ -th one at time  $\ell$ . Compare this with Figure 2: the vertical lines corresponding to the vertices have been reordered in increasing order of their birth time, and the grey arrows that left no trace on the graph at the focal time have been removed.

**Proposition 2.1.** *The random forest obtained by relabeling the vertices of  $\mathcal{F}_n^*$  according to a uniform permutation of  $\{1, \dots, n\}$  is distributed as the Moran forest.*

*Proof.* Consider the forest  $\mathcal{F}_n(0)$  built from the variables  $((V_t, W_t), t \in \mathbb{Z})$  in the previous section. To ease notation, we will omit the subscript in  $\tau_0$  and  $m_0$ . The proof hinges on the fact that there is a natural coupling of  $\mathcal{F}_n(0)$  with a forest  $\mathcal{F}_n^*$  having the aforementioned distribution, in a way such that conditional on  $\mathcal{F}_n^* = \mathcal{F}$ ,  $\mathcal{F}_n(0)$  is a uniform relabeling of  $\mathcal{F}$ .

Let us relabel the vertices in increasing order of their birth time: since the variables  $(\tau(v), 1 \leq v \leq n)$  are all distinct, there exists a unique permutation  $\sigma$  of  $\{1, \dots, n\}$  such that

$$\tau(\sigma(1)) < \dots < \tau(\sigma(n)).$$

In words,  $\sigma(\ell)$  is the  $\ell$ -th vertex that was born in the construction of  $\mathcal{F}_n(0)$ . Using the new labeling, let us denote its birth time by  $\tau^*(\ell) = \tau(\sigma(\ell))$  and its mother by  $m^*(\ell) = \sigma^{-1}(m(\sigma(\ell)))$ .

Now, for every vertex  $v = \sigma(\ell)$ ,

$$\begin{aligned} \text{there is an edge from } m(v) \text{ to } v \text{ in } \mathcal{F}_n(0) &\iff \tau(m(v)) < \tau(v) \\ &\iff \tau^*(m^*(\ell)) < \tau^*(\ell) \\ &\iff m^*(\ell) < \ell. \end{aligned}$$

Thus, if we set  $U_n(\ell) = m^*(\ell)$  in the construction of  $\mathcal{F}_n^*$  then  $\ell$  is connected to  $m^*(\ell)$  if and only if  $v = \sigma(\ell)$  is connected to  $m(v) = \sigma(m^*(\ell))$  in  $\mathcal{F}_n(0)$ . Therefore, to finish the proof we have to show that:

- (i) The variables  $(m^*(\ell), 1 \leq \ell \leq n)$  are independent and such that  $m^*(\ell)$  is uniformly distributed on  $\{1, \dots, n\} \setminus \{\ell\}$ .
- (ii) The permutation  $\sigma$  is uniform and independent of  $(m^*(\ell), 1 \leq \ell \leq n)$ .

First, note that by construction the variables  $(m(v), 1 \leq v \leq n)$  are independent and that for each  $v$ ,  $m(v)$  is uniformly distributed on  $\{1, \dots, n\} \setminus \{v\}$ . Since the permutation  $\sigma$  depends only on the variables  $(\tau(v), 1 \leq v \leq n)$ , which are independent of  $(m(v), 1 \leq v \leq n)$ , we see that  $\sigma$  is independent of  $(m(v), 1 \leq v \leq n)$ . Moreover, the variables  $(\tau(v), 1 \leq v \leq n)$  are exchangeable so the permutation  $\sigma$  is uniform. Now, for any fixed permutation  $\pi$  of  $\{1, \dots, n\}$  and any fixed map  $f : \{1, \dots, n\} \rightarrow \{1, \dots, n\}$  such that  $f(\ell) \neq \ell$  for all  $\ell$ , we have

$$\begin{aligned} \mathbb{P}(\sigma = \pi, m^* = f) &= \mathbb{P}(\sigma = \pi, m = \pi \circ f \circ \pi^{-1}) \\ &= \frac{1}{n!} \frac{1}{(n-1)^n}, \end{aligned}$$

concluding the proof. □

One might hope to obtain a complete description of the distribution of the Moran forest from the UA construction. Indeed, let us denote by  $(S_1, \dots, S_{N_n})$  the sizes of the trees in the Moran forest, labeled in decreasing order of their sizes, and where  $N_n$  is the number of trees. Then it is clear from the UA construction, that conditional on the vector  $(S_1, \dots, S_{N_n})$ , the trees of the Moran forest are independent uniform attachment trees. Therefore, the study of the Moran forest reduces to that of the distribution of  $(S_1, \dots, S_{N_n})$ . In the terminology of exchangeable partitions, we need to study the *exchangeable partition probability function* (EPPF) of the Moran forest [23]. However, we could not find any closed expression for this EPPF. Note that the Moran forest is *not* sampling consistent, i.e., the restriction of  $\mathcal{F}_{n+1}$  to  $\{1, \dots, n\}$  is not distributed as  $\mathcal{F}_n$ . Therefore the EPPF of the Moran forest cannot be obtained through a Chinese restaurant process.



### 3 Number of trees

#### 3.1 Law of the number of trees

In the UA construction, let  $I_\ell = \mathbb{1}_{\{U_n(\ell) < \ell\}}$  be the indicator variable of the event “the  $\ell$ -th vertex was linked to a previously added vertex”. The variables  $(I_1, \dots, I_n)$  are thus independent Bernoulli variables such that

$$I_\ell \sim \text{Bernoulli}\left(\frac{\ell-1}{n-1}\right).$$

With this notation, the number of edges  $|E_n|$  and the number of trees  $N_n$  are

$$|E_n| = \sum_{\ell=1}^n I_\ell \quad \text{and} \quad N_n = \sum_{\ell=1}^n (1 - I_\ell).$$

Moreover, since  $I_\ell \stackrel{d}{=} 1 - I_{n-\ell+1}$ , we see that

$$\mathbb{P}(N_n = k) = \mathbb{P}(N_n = n - k) = \mathbb{P}(|E_n| = k),$$

that is, the number of trees and the number of edges have the same, symmetric distribution. In consequence, from now on we only use the notation  $N_n$  and refer to it as the number of trees of  $\mathcal{F}_n$  when stating our results – even though we sometimes work with the number of edges in the proofs.

From the representation of  $N_n$  as a sum of independent Bernoulli variables, we immediately get the following result.

**Proposition 3.1.** *Let  $N_n$  denote the number of trees of  $\mathcal{F}_n$ .*

- (i)  $\mathbb{E}(N_n) = \frac{n}{2}$ .
- (ii)  $\text{Var}(N_n) = \frac{n(n-2)}{6(n-1)}$ .
- (iii)  $G_{N_n}(z) := \mathbb{E}(z^{N_n}) = \prod_{k=1}^{n-1} \left(1 + \frac{k}{n-1}(z-1)\right)$ .

The representation of  $N_n$  as a sum of independent Bernoulli variables also makes it straightforward to get the following central limit theorem.

**Proposition 3.2.** *Let  $N_n$  denote the number of trees of  $\mathcal{F}_n$ . Then,*

$$\frac{N_n - n/2}{\sqrt{n/6}} \xrightarrow[n \rightarrow \infty]{d} \mathcal{N}(0, 1).$$

*Proof.* This is an immediate consequence of the Lyapunov CLT for triangular arrays of independent random variables. Indeed,  $\mathbb{E}(|I_\ell - \mathbb{E}(I_\ell)|^3) \leq 1$ . Therefore,

$$\frac{1}{n^{3/2}} \sum_{\ell=1}^n \mathbb{E}(|I_\ell - \mathbb{E}(I_\ell)|^3) \leq \frac{1}{\sqrt{n}} \xrightarrow[n \rightarrow \infty]{} 0,$$

and the result follows, e.g., from Corollary 11.1.4 in [4]. □

### 3.2 Link with uniform labeled trees

As announced in the introduction, there is a strong connection between the Moran forest and uniform labeled trees. Our starting point is the following observation about the probability generating function of  $N_n$ . First, we can rewrite point (iii) of Proposition 3.1 as

$$\begin{aligned} G_{N_n}(z) &= \frac{z}{(n-1)^{n-2}} \prod_{k=1}^{n-2} (n-1-k+kz) \\ &= \sum_{k=0}^{n-2} \frac{a(n-1, k)}{(n-1)^{n-2}} z^{k+1}, \end{aligned}$$

where

$$\sum_{k=0}^{n-2} a(n-1, k) z^k = \prod_{k=1}^{n-2} (n-1-k+kz).$$

Second, the coefficients of this polynomial have a simple combinatorial interpretation:  $a(n-1, k)$  is the number of rooted trees on  $\{1, \dots, n-1\}$  with  $k$  increasing edges, where an edge  $u\vec{v}$  pointing away from the root is said to be increasing if  $u < v$ . This fact is known in the literature as a consequence of the more general Theorem 1.1 of [11] (see also Example 1.7.2 in [8] and Theorem 9.1 in [13]).

This simple observation already gives us the following proposition.

**Proposition 3.3.** *The probability mass function of the number of trees of  $\mathcal{F}_n$  is*

$$\mathbb{P}(N_n = k) = \frac{a(n-1, k-1)}{(n-1)^{n-2}},$$

where  $a(n, k)$  is the number of rooted trees on  $\{1, \dots, n\}$  with  $k$  increasing edges (sequence [A067948](#) of the *On-Line Encyclopedia of Integer Sequences* [1]).

Looking for a bijective proof of Proposition 3.3 naturally leads to the more general Theorem 1.1, which states that the Moran forest  $\mathcal{F}_n$  can be obtained from a uniform rooted tree on  $\{1, \dots, n-1\}$ , denoted by  $\mathcal{T}$ , using the following procedure:

1. Remove all decreasing edges from  $\mathcal{T}$  (that is, edges  $u\vec{v}$  pointing away from the root such that  $u > v$ ).
2. Add a vertex labeled  $n$  and connect it to a uniformly chosen vertex of  $\mathcal{T}$ .
3. Relabel vertices according to a uniform permutation of  $\{1, \dots, n\}$ .

*Proof of Theorem 1.1.* In the UA construction, let  $F|_{n-1}$  denote the forest obtained after the addition of  $n-1$  vertices, before their relabeling. After this, the  $n$ -th vertex will be linked to a uniformly chosen vertex of  $F|_{n-1}$ . As a result, to prove the theorem it suffices to show that  $F|_{n-1}$  has the same law as the forest obtained from  $\mathcal{T}$  by removing its decreasing edges.

To do so, we couple  $F|_{n-1}$  and  $\mathcal{T}$  in such a way that the edges of  $F|_{n-1}$  are exactly the increasing edges of  $\mathcal{T}$ . Formally,  $F|_{n-1}$  is a deterministic function of the random vector  $\mathbf{U} = (U_n(2), \dots, U_n(n-1))$ . Moreover,  $\mathbf{U}$  is uniform on the set

$$\mathcal{S}_{n-1}^* = \left\{ \mathbf{u} \in \{1, \dots, n\}^{\{2, \dots, n-1\}} : u_\ell \neq \ell \right\}.$$

Thus, to end the proof it is sufficient to find a bijection  $\Phi$  from  $\mathcal{S}_{n-1}^*$  to the set of rooted trees on  $\{1, \dots, n-1\}$  and such that

$$k\ell \in F_{|n-1}(\mathbf{u}) \iff k\ell \text{ is an increasing edge of } \Phi(\mathbf{u}).$$

First, let

$$\mathcal{S}_{n-1} = \{1, \dots, n-1\}^{\{2, \dots, n-1\}}$$

and consider the bijection  $\Theta : \mathcal{S}_{n-1}^* \rightarrow \mathcal{S}_{n-1}$  defined by

$$\Theta \mathbf{u} : \ell \mapsto u_\ell - \mathbf{1}_{\{u_\ell > \ell\}}.$$

Importantly, note that  $\Theta$  does not modify the entries of  $\mathbf{u}$  that correspond to edges of  $F_{|n-1}(\mathbf{u})$ , that is, for all  $k < \ell$ ,

$$k\ell \in F_{|n-1}(\mathbf{u}) \iff u_\ell = k \iff (\Theta \mathbf{u})(\ell) = k.$$

As a result, it remains to find a bijection  $\Psi$  from  $\mathcal{S}_{n-1}$  to the set of rooted trees on  $\{1, \dots, n-1\}$  such that

$$u_\ell < \ell \iff u_\ell \text{ and } \ell \text{ are linked by an increasing edge in } \Psi(\mathbf{u}).$$

This bijection will essentially be that used in [11], which can itself be seen as a variant of Joyal's bijection [2, 17].

Let  $\mathcal{G}_{\mathbf{u}}$  be the directed graph on  $\{1, \dots, n-1\}$  obtained by putting a directed edge going from  $u_\ell$  to  $\ell$  for all  $\ell \geq 2$ .

If  $\mathcal{G}_{\mathbf{u}}$  has no cycle or self-loop, then it is a tree. Moreover, the orientation of its edges uniquely identify vertex 1 as its root. Thus we set  $\Psi(\mathbf{u}) = \mathcal{G}_{\mathbf{u}}$ .

If  $\mathcal{G}_{\mathbf{u}}$  is not a tree, set  $\mathcal{C}_0 = \{1\}$  and let  $\mathcal{C}_1, \dots, \mathcal{C}_k$  denote the cycles of  $\mathcal{G}_{\mathbf{u}}$ , taken in increasing order of their largest element and treating self-loops as cycles of length 1. Note that because each vertex has exactly one incoming edge, except for vertex 1 which has none, these cycles are vertex-disjoint and directed.

To turn  $\mathcal{G}_{\mathbf{u}}$  into a tree, set  $s_0 = 1$  and for  $i \geq 1$  let  $m_i$  denote the largest element of  $\mathcal{C}_i$  and  $m_i \vec{s}_i$  its out-going edge in  $\mathcal{C}_i$ . With this notation, for  $i = 1, \dots, k$  remove the edge  $m_i \vec{s}_i$  from  $\mathcal{G}_{\mathbf{u}}$  and replace it by  $m_i \vec{s}_{i-1}$ . Note that

- This turns  $\mathcal{C}_0 \sqcup \dots \sqcup \mathcal{C}_k$  into a directed path  $\mathcal{P}$  going from  $s_k$  to 1.
- Because  $m_i = \max \mathcal{C}_i$  and that  $1 < m_1 < \dots < m_k$ , every edge  $m_i \vec{s}_i$  was non-increasing and has been replaced by the decreasing edge  $m_i \vec{s}_{i-1}$ .

Therefore, this procedure turns  $\mathcal{G}_{\mathbf{u}}$  into a tree  $\Psi(\mathbf{u})$  rooted in  $s_k$ , without modifying its increasing edges. Consequently, the increasing edges of  $\Psi(\mathbf{u})$  are exactly the pairs  $k\ell$  for which  $k = u_\ell < \ell$ .

To see that  $\Psi$  is a bijection, it suffices to note that the cycles  $\mathcal{C}_0, \dots, \mathcal{C}_k$  can be recovered unambiguously from the path  $\mathcal{P}$  going from the root to vertex 1. Indeed, writing this path as the word  $1m_1 \dots s_1 m_2 \dots s_k$ , the  $m_i$  are exactly the left-to-right maxima of that word.

Setting  $\Phi = \Psi \circ \Theta$  thus gives us the bijection that we were looking for, concluding the proof.  $\square$

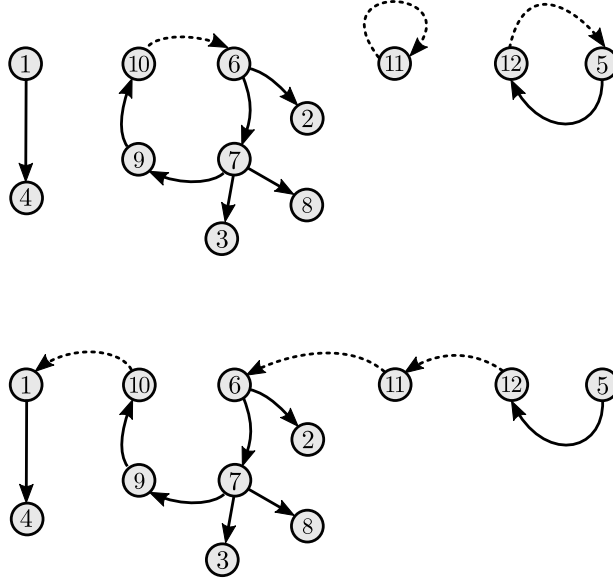


Figure 4: Example of construction of  $\Phi(\mathbf{u})$ , for  $\mathbf{u} = (7, 8, 1, 13, 11, 6, 7, 7, 9, 12, 5)$ . Applying  $\Theta$  yields  $\mathbf{u}' = \Theta\mathbf{u} = (6, 7, 1, 12, 10, 6, 7, 7, 9, 11, 5)$ . The directed graph  $\mathcal{G}_{\mathbf{u}'}$  encoding  $\mathbf{u}'$  is represented on top. Its cycles are  $\mathcal{C}_1 = (10, 6, 7, 9)$ ,  $\mathcal{C}_2 = (11)$  and  $\mathcal{C}_3 = (12, 5)$ , and we set  $\mathcal{C}_0 = (1)$ . The edges  $m_i^s_i$  are dashed. Rewiring them as described in the main text turns  $\mathcal{G}_{\mathbf{u}'}$  into the rooted tree  $\Psi(\mathbf{u}')$  represented on bottom. No information is lost when turning the cycles  $(1)(10, 6, 7, 9)(11)(12, 5)$  into the path going from 5 to 1 encoded by the word  $(1, 10, 6, 7, 9, 11, 12, 5)$ , because the left-to-right maxima of that word – here 1, 10, 11 and 12 – each mark the start of a new cycle.

## 4 Degrees

### 4.1 Degree of a fixed vertex

Using the UA construction and the notation from Section 2.2, let us denote by

- $I_\ell = \mathbb{1}_{\{U_n(\ell) < \ell\}}$  the indicator variable of the event “the  $\ell$ -th vertex has an incoming edge linking it to a previously added vertex”.
- $X_\ell^{(v)} = \mathbb{1}_{\{U_n(\ell) = \sigma^{-1}(v)\}}$  the indicator variable of the event “the  $\ell$ -th vertex is linked to vertex  $v$ ”.
- $B_v = \sigma^{-1}(v)$  the step of the construction at which vertex  $v$  is added.

With this notation, the degree of vertex  $v$  is

$$D_n^{(v)} = I_{B_v} + \sum_{\ell=B_v+1}^n X_\ell^{(v)},$$

where  $I_{B_v}$  is the in-degree, and  $\sum_{\ell=B_v+1}^n X_\ell^{(v)}$  is the out-degree of vertex  $v$ . Moreover, conditional on  $\{B_v = b\}$ ,  $(X_{b+1}^{(v)}, \dots, X_n^{(v)})$  are i.i.d. Bernoulli variables with parameter  $1/(n-1)$  and  $I_b$  is a Bernoulli variable with parameter  $\frac{b-1}{n-1}$  that is independent of  $(X_{b+1}^{(v)}, \dots, X_n^{(v)})$ . As a result, conditional on  $B_v$  and writing  $L_v$  for  $n - B_v$ ,

$$D_n^{(v)} \stackrel{d}{=} \text{Ber}\left(1 - \frac{L_v}{n-1}\right) + \text{Bin}\left(L_v, \frac{1}{n-1}\right),$$

where the Bernoulli and the binomial variables are independent conditional on  $L_v$  (here, as in similar expressions in the rest of this document, the Ber and Bin

notation refers to random variables with the corresponding distribution, not the distributions themselves). Using that  $L_v$  is uniformly distributed on  $\{0, \dots, n-1\}$ , the mean, variance and probability generating function of  $D_n^{(v)}$  are obtained by routine calculations.

**Proposition 4.1.** *Let  $D_n$  be the degree of a fixed vertex of  $\mathcal{F}_n$ . Then,*

- (i)  $\mathbb{E}(D_n) = 1.$
- (ii)  $\text{Var}(D_n) = \frac{2(n-2)}{3(n-1)}.$
- (iii)  $G_{D_n}(z) := \mathbb{E}(z^{D_n}) = \frac{1}{n} \sum_{\ell=0}^{n-1} \left(1 + (1 - \frac{\ell}{n-1})(z-1)\right) \left(1 + \frac{1}{n-1}(z-1)\right)^\ell.$
- (iii')  $G_{D_n}(z) = 2 \left(1 - \frac{1}{n}\right) \frac{\left(1 + \frac{z-1}{n-1}\right)^n - 1}{z-1} - 1.$

**Remark 4.2.** Note that, conditional on  $L_v$ , the probability that  $v$  has a (unique) incoming edge is  $1 - \frac{L_v}{n-1}$ , and its mean out-degree is  $\frac{L_v}{n-1}$ . Therefore, summing the two we have  $\mathbb{E}(D_n^{(v)} | L_v) = 1$ , that is, the average degree of a vertex is independent of the step at which it was added in the UA construction.  $\diamond$

**Proposition 4.3.** *The degree  $D_n$  of a fixed vertex of  $\mathcal{F}_n$  converges in distribution to the variable  $D$  satisfying:*

- (i)  $D \sim \text{Ber}(1-U) + \text{Poisson}(U)$ , where  $U$  is uniform on  $[0, 1]$  and the Bernoulli and Poisson variables are independent conditional on  $U$ .
- (ii)  $G_D(z) := \mathbb{E}(z^D) = \int_0^1 \left(1 + (1-x)(z-1)\right) e^{x(z-1)} dx = 2 \frac{e^{z-1} - 1}{z-1} - 1.$
- (iii) For all  $p \geq 1$ ,  $\mathbb{E}(D(D-1)\cdots(D-p+1)) = \frac{2}{p+1}.$
- (iv)  $\mathbb{P}(D=0) = 1 - 2/e$  and, for  $k \geq 1$ ,

$$\mathbb{P}(D=k) = \frac{2}{e} \sum_{j>k} \frac{1}{j!}.$$

*Proof.* First, for all  $z \in \mathbb{C} \setminus \{1\}$ ,

$$G_{D_n}(z) = 2 \left(1 - \frac{1}{n}\right) \frac{\left(1 + \frac{z-1}{n-1}\right)^n - 1}{z-1} - 1 \xrightarrow{n \rightarrow \infty} 2 \frac{e^{z-1} - 1}{z-1} - 1.$$

This pointwise convergence of the probability generating function of  $D_n$  proves the convergence in distribution of  $D_n$  to a random variable  $D$  satisfying (ii). Point (i) then follows immediately from the integral expression of  $G_D$ .

To compute the factorial moments of  $D$ , note that

$$G_D(z) = 2 \sum_{k \geq 0} \frac{(z-1)^k}{(k+1)!} - 1.$$

As a result, for  $p \geq 1$  the  $p$ -th derivative of  $G_D$  is

$$G_D^{(p)}(z) = 2 \sum_{k \geq 0} \frac{(z-1)^k}{(k+1+p)k!},$$

and, in particular,  $\mathbb{E}(D(D-1)\cdots(D-p+1)) = G_D^{(p)}(1) = \frac{2}{p+1}$ , proving (iii).

Finally, to prove (iv), using (i) we see that

$$\mathbb{P}(D=0) = \int_0^1 x e^{-x} dx = 1 - \frac{2}{e}$$

and that, for  $k \geq 1$ ,

$$\mathbb{P}(D=k) = \frac{1}{k!} \int_0^1 (kx^{k-1} - kx^k + x^{k+1}) e^{-x} dx.$$

Noting that  $(kx^{k-1} - kx^k + x^{k+1})e^{-x} = 2x^k e^{-x} + \frac{d}{dx}((x^k - x^{k+1})e^{-x})$ , we get

$$\mathbb{P}(D=k) = \frac{2}{k!} \int_0^1 x^k e^{-x} dx,$$

and an easy integration by parts yields

$$\mathbb{P}(D=k+1) = \mathbb{P}(D=k) - \frac{2}{e(k+1)!},$$

from which (iv) follows by induction.  $\square$

Before closing this section, let us give an asymptotic equivalent of the tail of  $D_n$ . We will need it in the proof of Theorem 1.2 on the largest degree.

**Proposition 4.4.** *Let  $D_n$  be the degree of a fixed vertex of  $\mathcal{F}_n$  and let  $D$  have the asymptotic distribution of  $D_n$ .*

(i) *For all  $k \geq 1$ ,*

$$\frac{2/e}{(k+1)!} \leq \mathbb{P}(D \geq k) \leq \left(1 + \frac{1}{k}\right)^2 \frac{2/e}{(k+1)!}.$$

(ii) *For all  $K_n = o(\sqrt{n})$ , there exists  $\varepsilon_n = o(1)$  such that, for all  $k \leq K_n$ ,*

$$|\mathbb{P}(D_n \geq k) - \mathbb{P}(D \geq k)| \leq \varepsilon_n \mathbb{P}(D \geq k).$$

(iii) *For all  $k_n \rightarrow +\infty$  and  $K_n \geq k_n$  such that  $K_n = o(\sqrt{n})$ ,*

$$\mathbb{P}(D_n \geq k) \sim \frac{2/e}{(k+1)!},$$

*uniformly in  $k$  such that  $k_n \leq k \leq K_n$ .*

*Proof.* First, observe that

$$\frac{1}{(\ell+1)!} \leq \frac{1}{\ell \cdot \ell!} - \frac{1}{(\ell+1) \cdot (\ell+1)!},$$

so that

$$\sum_{\ell>i} \frac{1}{\ell!} \leq \frac{1}{i \cdot i!} = \left(1 + \frac{1}{i}\right) \frac{1}{(i+1)!}.$$

Recalling from Proposition 4.3 that

$$\mathbb{P}(D \geq k) = \frac{2}{e} \sum_{i \geq k} \sum_{\ell > i} \frac{1}{\ell!},$$

point (i) follows readily.

The proof of (ii) is somewhat technical so we only outline it here and refer the reader to Section A.1 of the Appendix for the detailed calculations.

Consider the function

$$\Delta_n(z) = \sum_{i \geq 0} \left( \mathbb{P}(D \geq i) - \mathbb{P}(D_n \geq i) \right) z^i.$$

With this function, (ii) can be re-expressed as

$$\Delta_n^{(k)}(0) = \frac{\varepsilon_n}{k+1} \quad \text{for all } k \leq K_n = o(\sqrt{n}),$$

where  $\Delta_n^{(k)}$  denotes the  $k$ -th derivative of  $\Delta_n$ . But  $\Delta_n$  can be expressed in terms of the generating functions of  $D$  and  $D_n$ , namely as

$$\Delta_n(z) = \left(1 + \frac{1}{z-1}\right) (G_D(z) - G_{D_n}(z)).$$

The expressions of  $G_D$  and  $G_{D_n}$  obtained in Propositions 4.1 and 4.3 thus make it straightforward to obtain a power series expansion of  $\Delta_n$  at  $z = 1$ , and this expansion can be used to bound  $\Delta_n^{(k)}(0)$  and conclude the proof.

Finally, (iii) is a direct consequence of (i) and (ii).  $\square$

## 4.2 Largest degree

The aim of this section is to prove Theorem 1.2 concerning the largest degree of  $\mathcal{F}_n$ , i.e. to show that

$$D_n^{\max} = \frac{\log n}{\log \log n} + \left(1 + o_p(1)\right) \frac{\log n \log \log \log n}{(\log \log n)^2},$$

where  $o_p(1)$  denotes a sequence of random variables that goes to 0 in probability and  $D_n^{\max} = \max_v D_n^{(v)}$ .

Our proof is a standard application of the first and second moment method and it implies that, asymptotically,  $D_n^{\max}$  behaves like the maximum of  $n$  independent random variables distributed as  $D$ . As is typically the case with this method, the main difficulty – and therefore the bulk of the proof – consists in bounding the asymptotic dependency between two fixed vertices of  $\mathcal{F}_n$ .

Because the first and second moment method part of our reasoning will also be used in the proof of Theorem 1.3 concerning the size of the largest tree, we isolate it as a lemma, whose proof we recall for the sake of completeness.

**Lemma 4.5.** For all integers  $n$ , let  $(X_n^{(1)}, \dots, X_n^{(n)})$  be a vector of exchangeable random variables and

$$X_n^{\max} = \max\{X_n^{(i)} : i = 1, \dots, n\}.$$

Write  $p_n(k)$  for  $\mathbb{P}(X_n^{(i)} \geq k)$ , and suppose that there exists a sequence  $(m_n)$  and a constant  $\beta$  such that, for all  $\varepsilon > 0$ , as  $n \rightarrow \infty$ ,

- (i)  $np_n((\beta + \varepsilon)m_n) \rightarrow 0$ .
- (ii)  $np_n((\beta - \varepsilon)m_n) \rightarrow +\infty$ .
- (iii)  $\mathbb{P}(X_n^{(1)} \geq (\beta - \varepsilon)m_n, X_n^{(2)} \geq (\beta - \varepsilon)m_n) \sim p_n((\beta - \varepsilon)m_n)^2$ .

Then for all  $\varepsilon > 0$ ,

$$\mathbb{P}(X_n^{\max} \geq (\beta + \varepsilon)m_n) \rightarrow 0 \quad \text{and} \quad \mathbb{P}(X_n^{\max} \geq (\beta - \varepsilon)m_n) \rightarrow 1,$$

which can also be written

$$X_n^{\max} = (\beta + o_p(1))m_n,$$

where  $o_p(1)$  denotes a sequence of random variables that goes to 0 in probability.

*Proof.* First,

$$\begin{aligned} \mathbb{P}(X_n^{\max} \geq (\beta + \varepsilon)m_n) &= \mathbb{P}\left(\bigcup_{i=1}^n \{X_n^{(i)} \geq (\beta + \varepsilon)m_n\}\right) \\ &\leq np_n((\beta + \varepsilon)m_n), \end{aligned}$$

which goes to zero by (i). Now, denote by

$$Z_n = \sum_{i=1}^n \mathbb{1}_{\{X_n^{(i)} \geq (\beta - \varepsilon)m_n\}}$$

the number of variables  $X_n^{(i)}$  that are greater than or equal to  $(\beta - \varepsilon)m_n$ . Using the Cauchy–Schwartz inequality, we have  $\mathbb{E}(Z_n)^2 = \mathbb{E}(Z_n \mathbb{1}_{\{Z_n > 0\}})^2 \leq \mathbb{E}(Z_n^2) \mathbb{P}(Z_n > 0)$  so

$$\mathbb{P}(X_n^{\max} \geq (\beta - \varepsilon)m_n) = \mathbb{P}(Z_n > 0) \geq \frac{\mathbb{E}(Z_n)^2}{\mathbb{E}(Z_n^2)}.$$

Moreover,

$$\begin{aligned} \mathbb{E}(Z_n^2) &= np_n((\beta - \varepsilon)m_n) \\ &\quad + n(n-1)\mathbb{P}(X_n^{(1)} \geq (\beta - \varepsilon)m_n, X_n^{(2)} \geq (\beta - \varepsilon)m_n), \end{aligned}$$

and so, by (ii) and (iii),  $\mathbb{E}(Z_n)^2/\mathbb{E}(Z_n^2) \rightarrow 1$  as  $n \rightarrow \infty$ . □

**Remark 4.6.** Note that under assumption (ii) of this lemma, for any  $\varepsilon > 0$ , letting  $n \rightarrow \infty$  in  $\mathbb{E}(Z_n)^2/\mathbb{E}(Z_n^2) \leq 1$  shows that

$$p_n((\beta - \varepsilon)m_n)^2 \leq \mathbb{P}(X_n^{(1)} \geq (\beta - \varepsilon)m_n, X_n^{(2)} \geq (\beta - \varepsilon)m_n)(1 + o(1)).$$

Therefore, to prove (iii) it suffices to show



$$(iii') \quad \mathbb{P}\left(X_n^{(1)} \geq (\beta - \varepsilon)m_n, X_n^{(2)} \geq (\beta - \varepsilon)m_n\right) \leq p_n((\beta - \varepsilon)m_n)^2(1 + o(1)). \quad \diamond$$

We now turn to the proof of Theorem 1.2.

*Proof of Theorem 1.2.* Instead of proving the theorem directly for the degree random variables  $(D_n^{(1)}, \dots, D_n^{(n)})$ , we prove it for the out-degrees, which we denote by  $(\tilde{D}_n^{(1)}, \dots, \tilde{D}_n^{(n)})$  and whose maximum has the same asymptotic behavior as  $D_n^{\max}$ . The point in doing this is that the tails of the out-degrees are less correlated than those of the variables  $D_n^{(v)}$ , making it easier to study their maximum by the first and second moment method.

Remember from Section 4.1 that, in the UA construction,

$$D_n^{(v)} = I_{B_v} + \sum_{\ell=B_v+1}^n X_\ell^{(v)},$$

where  $B_v$  is the step at which vertex  $v$  was added,  $X_\ell^{(v)}$  is the indicator of “the  $\ell$ -th vertex is linked to vertex  $v$ ”, and  $I_\ell$  is the indicator of “the  $\ell$ -th vertex is linked to a previously added vertex”. With this notation, let

$$\tilde{D}_n^{(v)} = \sum_{\ell=B_v+1}^n X_\ell^{(v)},$$

denote the out-degree of vertex  $v$ , and set  $\tilde{D}_n^{\max} = \max\{\tilde{D}_n^{(v)} : v = 1, \dots, n\}$ . Since  $\tilde{D}_n^{\max}$  and  $D_n^{\max}$  differ by at most 1, for any  $m_n \rightarrow +\infty$ ,

$$D_n^{\max} - \tilde{D}_n^{\max} = o_p(m_n),$$

i.e.  $(D_n^{\max} - \tilde{D}_n^{\max})/m_n$  goes to 0 in probability. Thus, to prove the theorem we apply Lemma 4.5 to the variables

$$\left(\tilde{D}_n^{(1)} - \frac{\log n}{\log \log n}, \dots, \tilde{D}_n^{(n)} - \frac{\log n}{\log \log n}\right),$$

with  $m_n = (\log n)(\log \log \log n)/(\log \log n)^2$  and  $\beta = 1$ .

Using Proposition 4.4 and Stirling’s formula, we see that for any  $k_n = o(\sqrt{n})$ ,

$$\log\left(\mathbb{P}(D_n \geq k_n)\right) = -k_n \log k_n + k_n + O(\log k_n).$$

Writing  $\tilde{D}_n$  to refer to the common distribution of the variables  $\tilde{D}_n^{(v)}$ , since

$$\mathbb{P}(D_n \geq k_n + 1) \leq \mathbb{P}(\tilde{D}_n \geq k_n) \leq \mathbb{P}(D_n \geq k_n),$$

we also have

$$\log\left(\mathbb{P}(\tilde{D}_n \geq k_n)\right) = -k_n \log k_n + k_n + O(\log k_n).$$

In particular, for  $k_n = (\log n)/(\log \log n) + \gamma m_n$  with

$$m_n = \frac{\log n \log \log \log n}{(\log \log n)^2},$$

this gives

$$\log\left(\mathbb{P}(\tilde{D}_n \geq k_n)\right) = -\log n - (\gamma - 1) \frac{\log n \log \log \log n}{\log \log n} + O\left(\frac{\log n}{\log \log n}\right) \quad (1)$$

As a result, for all  $\varepsilon > 0$ ,

- (i)  $n \mathbb{P} \left( \tilde{D}_n - \frac{\log n}{\log \log n} \geq (1 + \varepsilon)m_n \right) \rightarrow 0.$
- (ii)  $n \mathbb{P} \left( \tilde{D}_n - \frac{\log n}{\log \log n} \geq (1 - \varepsilon)m_n \right) \rightarrow +\infty.$

Thus, to apply Lemma 4.5 and finish the proof it suffices to show that

$$\mathbb{P} \left( \tilde{D}_n^{(1)} \geq k_n, \tilde{D}_n^{(2)} \geq k_n \right) \sim \mathbb{P} \left( \tilde{D}_n \geq k_n \right)^2$$

whenever  $k_n = (\log n)/(\log \log n) + (1 - \varepsilon)m_n$ . More precisely, using Remark 4.6 it is sufficient to show that

$$\mathbb{P} \left( \tilde{D}_n^{(1)} \geq k_n, \tilde{D}_n^{(2)} \geq k_n \right) \leq \mathbb{P} \left( \tilde{D}_n \geq k_n \right)^2 + o \left( \mathbb{P} \left( \tilde{D}_n \geq k_n \right)^2 \right).$$

First let us fix  $b_1 \neq b_2 \in \{1, \dots, n\}$ . Conditional on  $\{B_1 = b_1, B_2 = b_2\}$ , recall that the variables  $(X_\ell^{(2)}, b_2 + 1 \leq \ell \leq n)$  are independent Bernoulli variables with parameter  $1/(n - 1)$ . By further conditioning on the variables  $X_\ell^{(1)}$ , the independence of  $(X_\ell^{(2)}, b_2 + 1 \leq \ell \leq n)$  still holds but their distribution is changed. Indeed, choose  $(x_\ell, \ell \neq b_1) \in \{0, 1\}^{n-1}$  and consider the event

$$A := \left\{ B_1 = b_1, B_2 = b_2, \forall \ell \neq b_1, X_\ell^{(1)} = x_\ell \right\}.$$

Then by construction, for all  $\ell \notin \{b_1, b_2\}$ , we have

$$\mathbb{P} \left( X_\ell^{(2)} = 1 \mid A \right) = \begin{cases} 0 & \text{if } x_\ell = 1 \\ \frac{1}{n-2} & \text{if } x_\ell = 0. \end{cases}$$

Consequently  $X_\ell^{(2)}$  is always stochastically dominated by a Bernoulli( $\frac{1}{n-2}$ ) random variable, and so we bound the distribution of  $\tilde{D}_n^{(2)} = \sum_{\ell > b_2} X_\ell^{(2)}$  conditional on  $A$  by

$$\left( \tilde{D}_n^{(2)} \mid A \right) \stackrel{d}{\leq} \text{Binomial} \left( n - b_2, \frac{1}{n - 2} \right).$$

To get a bound on the distribution of  $\tilde{D}_n^{(2)}$  conditional on  $\tilde{D}_n^{(1)} = i$  for some  $i$ , first note that summing over all configurations  $b_1, b_2, (x_\ell, \ell \neq b_1)$  such that  $\sum_{\ell > b_1} x_\ell = i$  gives

$$\left( \tilde{D}_n^{(2)} \mid B_1 = b_1, B_2 = b_2, \tilde{D}_n^{(1)} = i \right) \stackrel{d}{\leq} \text{Binomial} \left( n - b_2, \frac{1}{n - 2} \right).$$

Let us now write for conciseness  $L_1 = n - B_1$  and  $L_2 = n - B_2$ . Note that  $L_2$  is not independent of  $\{\tilde{D}_n^{(1)} = i\}$  because they are linked by  $L_1$ . Indeed,  $L_1$  is positively correlated to  $\tilde{D}_n^{(1)}$  and we always have  $L_2 \neq L_1$ . Nevertheless, since conditional on  $L_1$ ,  $L_2$  is independent of  $\tilde{D}_n^{(1)}$  and uniform on  $\{0, \dots, n - 1\} \setminus L_1$ , we have the following stochastic ordering:

$$\left( L_2 \mid B_1 = b_1, \tilde{D}_n^{(1)} = i \right) \stackrel{d}{\leq} \bar{L}_2,$$

where  $\bar{L}_2$  is uniformly distributed on  $\{1, \dots, n - 1\}$ . Summing over  $b_1$  and  $b_2$ , we thus get

$$\left( \tilde{D}_n^{(2)} \mid \tilde{D}_n^{(1)} = i \right) \stackrel{d}{\leq} \text{Binomial} \left( \bar{L}_2, \frac{1}{n - 2} \right).$$

Let us define a random variable  $M_n \sim \text{Bin}\left(\bar{L}_2, \frac{1}{n-2}\right)$ . As the previous bound is uniform in  $i$ , we have

$$\mathbb{P}\left(\tilde{D}_n^{(2)} \geq k_n \mid \tilde{D}_n^{(1)} \geq k_n\right) \leq \mathbb{P}(M_n \geq k_n).$$

To conclude, it is sufficient to show that  $\mathbb{P}(M_n \geq k_n) \sim \mathbb{P}(\tilde{D}_n \geq k_n)$  since this implies

$$\mathbb{P}\left(\tilde{D}_n^{(1)} \geq k_n, \tilde{D}_n^{(2)} \geq k_n\right) \leq \mathbb{P}(\tilde{D}_n \geq k_n)\mathbb{P}(M_n \geq k_n) \sim \mathbb{P}(\tilde{D}_n \geq k_n)^2.$$

For this, define on the same probability space as the variables  $\bar{L}_2$  and  $M_n$  the variable

$$\underline{L}_2 := \bar{L}_2 \mathbb{1}_{\{\bar{L}_2 \leq n-2\}}.$$

$\underline{L}_2$  is then uniformly distributed on  $\{0, \dots, n-2\}$ , and we have the equality in distribution

$$M_n \mathbb{1}_{\{\bar{L}_2 \leq n-2\}} \stackrel{d}{=} \tilde{D}_{n-1} \sim \text{Binomial}\left(\underline{L}_2, \frac{1}{n-2}\right).$$

As the two variables  $M_n$  and  $M_n \mathbb{1}_{\{\bar{L}_2 \leq n-2\}}$  differ on an event of probability no greater than  $1/(n-1)$ , we have

$$\mathbb{P}(M_n \geq k_n) = \mathbb{P}(\tilde{D}_{n-1} \geq k_n) + O\left(\frac{1}{n}\right),$$

and finally (1) with  $\gamma = (1-\varepsilon)$  allows us to conclude that this expression is indeed equivalent to  $\mathbb{P}(\tilde{D}_n \geq k_n)$ .  $\square$

## 5 Tree sizes

In this section, we study the size of the trees composing the Moran forest. Section 5.2 is concerned with the typical size of these trees, while Section 5.3 focuses on the asymptotics of the size of the largest tree. But before going any further we need to introduce a process that will play a central role throughout the rest of this paper.

### 5.1 A discrete-time Yule process

Let  $\Upsilon_n = (\Upsilon_n(\ell), \ell \geq 0)$  be the pure birth Markov chain defined by  $\Upsilon_n(0) = 1$  and the following transition probabilities:

$$\mathbb{P}\left(\Upsilon_n(\ell+1) = j \mid \Upsilon_n(\ell) = i\right) = \begin{cases} \frac{i}{n-1} & \text{if } j = i+1 \\ 1 - \frac{i}{n-1} & \text{if } j = i, \end{cases}$$

and stopped when reaching  $n$ .

The reason why this process will play an important role when studying the trees of  $\mathcal{F}_n$  is the following: let  $\mathcal{T}_n^{(v)}$  denote the tree containing  $v$ , and  $\tilde{\mathcal{T}}_n^{(v)}$  the subtree descending from  $v$  in the UA construction – that is, letting  $m(v)$  denote the mother of  $v$  and  $\mathcal{T}_n^{(v)} \setminus \{vm(v)\}$  the forest obtained by removing the edge between  $v$  and  $m(v)$  from  $\mathcal{T}_n^{(v)}$  (if that edge existed),  $\tilde{\mathcal{T}}_n^{(v)}$  is the tree of  $\mathcal{T}_n^{(v)} \setminus \{vm(v)\}$  containing  $v$ .

Recalling that  $L_v$  denotes the number of steps after vertex  $v$  was added in the UA construction and letting  $\tilde{T}_n^{(v)} = |\tilde{\mathcal{T}}_n^{(v)}|$  be the size of  $\tilde{\mathcal{T}}_n^{(v)}$ , we have

$$\tilde{T}_n^{(v)} \stackrel{d}{=} \Upsilon_n(L_v),$$

where  $\Upsilon_n$  is independent of  $L_v$ . In particular, the size of a tree created at step  $n - h$  of the UA construction is distributed as  $\Upsilon_n(h)$ .

In the rest of this section, we list a few basic properties of  $\Upsilon_n$  that will be used in subsequent proofs.

**Lemma 5.1.** *For all  $0 \leq \ell \leq n - 1$ ,*

$$\mathbb{E}(\Upsilon_n(\ell)) = \left(1 + \frac{1}{n-1}\right)^\ell.$$

*Proof.* For  $0 \leq \ell < n - 1$ , we have  $\Upsilon_n(\ell) < n$  almost surely, therefore we can write

$$\begin{aligned} \mathbb{E}(\Upsilon_n(\ell + 1) \mid \Upsilon_n(\ell)) &= \frac{\Upsilon_n(\ell)}{n-1}(\Upsilon_n(\ell) + 1) + \left(1 - \frac{\Upsilon_n(\ell)}{n-1}\right)\Upsilon_n(\ell) \\ &= \Upsilon_n(\ell) \left(1 + \frac{1}{n-1}\right), \end{aligned}$$

and the result follows by induction.  $\square$

We now compare the discrete-time process  $\Upsilon_n$  to the Yule process. By Yule process, we refer to the continuous-time Markov chain  $(Y(t), t \geq 0)$  that jumps from  $i$  to  $i + 1$  at rate  $i$  (see e.g. [24], Section 5.3).

**Lemma 5.2.** *As  $n \rightarrow \infty$ ,*

$$\left(\Upsilon_n(\lfloor tn \rfloor), t \geq 0\right) \Longrightarrow \left(Y(t), t \geq 0\right),$$

where “ $\Longrightarrow$ ” denotes convergence in distribution in the Skorokhod space [7], and  $(Y(t), t \geq 0)$  is a Yule process.

*Proof.* Since both processes only have increments of  $+1$ , it suffices to prove that the sequence of jump times of  $(\Upsilon_n(\lfloor tn \rfloor), t \geq 0)$  converges in distribution to that of the Yule process. For  $1 \leq i \leq n$ , let

$$t_n(i) = \inf\{\ell \geq 0 : \Upsilon_n(\ell) = i\}$$

be the jump times of the chain  $\Upsilon_n$ . By the strong Markov property, the variables  $(t_n(i+1) - t_n(i), 1 \leq i \leq n-1)$  are independent, and  $t_n(i+1) - t_n(i) \sim \text{Geometric}(\frac{i}{n-1})$ . Therefore,

$$\left(\frac{1}{n}(t_n(i+1) - t_n(i)), 1 \leq i \leq n-1\right) \xrightarrow[n \rightarrow \infty]{d} (\mathcal{E}(i), i \geq 1),$$

where the variables  $(\mathcal{E}(i), i \geq 1)$  are independent and  $\mathcal{E}(i) \sim \text{Exponential}(i)$ . This concludes the proof.  $\square$

**Lemma 5.3.** For all integers  $0 \leq k \leq \ell \leq n - 1$ ,

$$\mathbb{P}\left(Y\left(\frac{\ell-k+1}{n-1}\right) > k\right) \leq \mathbb{P}(\Upsilon_n(\ell) > k) \leq \mathbb{P}\left(Y\left(\lambda_n(k)\frac{\ell}{n-1}\right) > k\right),$$

where

$$\lambda_n(k) = -\frac{n-1}{k} \log\left(1 - \frac{k}{n-1}\right).$$

*Proof.* Let us start with the upper bound, and write  $\lambda := \lambda_n(k)$  for simplicity. Note that, for all  $t \geq 0$  and  $i \geq 1$ ,

$$\mathbb{P}\left(Y\left(t + \frac{\lambda}{n-1}\right) = i \mid Y(t) = i\right) = e^{-\frac{i\lambda}{n-1}},$$

and that we have chosen  $\lambda$  such that if  $i \leq k$  then

$$e^{-\frac{i\lambda}{n-1}} \leq 1 - \frac{i}{n-1} = \mathbb{P}\left(\Upsilon_n(\ell+1) = i \mid \Upsilon_n(\ell) = i\right).$$

Thus, until it reaches  $k+1$  individuals, the process  $\Upsilon_n$  is dominated by the Markov chain  $(Y(\frac{\lambda\ell}{n-1}), 0 \leq \ell \leq n-1)$ . This shows that

$$\mathbb{P}(\Upsilon_n(\ell) > k) \leq \mathbb{P}\left(Y\left(\frac{\lambda\ell}{n-1}\right) > k\right),$$

proving the second inequality of the lemma.

To prove the first inequality, we couple  $\Upsilon_n$  with a ‘‘censored’’ Yule process  $Y_c$ . Intuitively, this censoring consists in ignoring births that occur less than  $1/(n-1)$  unit of time after another birth.

Formally, we define  $Y_c$  by specifying the sequence  $t_0 = 0 < t_1 < t_2 < \dots$  of times corresponding to births in the population. Let  $(\mathcal{E}_i, i \geq 1)$  be an independent sequence of exponential random variables where  $\mathcal{E}_i \sim \text{Exponential}(i)$ . Set  $t_0 = 0$  and, for each  $i \geq 1$ ,

$$t_i := \mathcal{E}_1 + \sum_{j=2}^i \left(\frac{1}{n-1} + \mathcal{E}_j\right) = \frac{i-1}{n-1} + \sum_{j=1}^i \mathcal{E}_j. \quad (2)$$

We now define, for all  $t \geq 0$ ,

$$Y_c(t) := 1 + \sum_{i \geq 1} \mathbf{1}_{\{t_i \leq t\}} = \sum_{i \geq 1} i \mathbf{1}_{\{t_{i-1} \leq t < t_i\}}.$$

The censoring of the Yule process after birth events implies that for any time  $t \geq 0$ , the random variable  $Y_c(t + \frac{1}{n-1}) - Y_c(t)$  takes values in  $\{0, 1\}$ . Furthermore, for any  $i \in \mathbb{N}$ ,

$$\mathbb{P}\left(Y_c\left(t + \frac{1}{n-1}\right) = i+1 \mid Y_c(t) = i\right) \leq 1 - e^{-\frac{i}{n-1}} \leq \frac{i}{n-1}.$$

Therefore, we can couple  $(\Upsilon_n(\ell), 0 \leq \ell \leq n-1)$  and  $(Y_c(t), t \geq 0)$  in such a way that, for all  $0 \leq \ell \leq n-1$ ,

$$Y_c\left(\frac{\ell}{n-1}\right) \leq \Upsilon_n(\ell).$$

Now, by construction, the sequence  $(t_i - \frac{i-1}{n-1}, i \geq 1)$  has the distribution of the sequence of jump times of a Yule process. Therefore,

$$\begin{aligned} \mathbb{P}(\Upsilon_n(\ell) > k) &\geq \mathbb{P}\left(Y_c\left(\frac{\ell}{n-1}\right) > k\right) \\ &= \mathbb{P}\left(t_k \leq \frac{\ell}{n-1}\right) \\ &= \mathbb{P}\left(t_k - \frac{k-1}{n-1} \leq \frac{\ell-k+1}{n-1}\right) \\ &= \mathbb{P}\left(Y\left(\frac{\ell-k+1}{n-1}\right) > k\right), \end{aligned}$$

which yields the lower bound of the lemma.  $\square$

We now use the previous lemma to obtain the following result, which will be used to derive asymptotics for the tail probability of the size of a tree in the Moran forest.

**Proposition 5.4.** *Let  $L$  be a uniform random variable on  $\{0, \dots, n-1\}$ , independent of the process  $\Upsilon_n$ . Then for any sequence of integers  $k_n \rightarrow \infty$  with  $k_n = o(\sqrt{n})$ ,*

$$\mathbb{P}(\Upsilon_n(L) > k_n) \sim \frac{e}{k_n}(1 - e^{-1})^{k_n+1}.$$

*Proof.* Using the upper bound in Lemma 5.3, we have

$$\begin{aligned} \mathbb{P}(\Upsilon_n(L) > k_n) &\leq \frac{1}{n} \sum_{\ell=0}^{n-1} \mathbb{P}\left(Y\left(\lambda_n(k) \frac{\ell}{n-1}\right) > k_n\right) \\ &= \frac{n-1}{n} \int_0^1 \mathbb{P}\left(Y\left(\lambda_n(k_n) \frac{\lfloor x(n-1) \rfloor}{n-1}\right) > k_n\right) dx + \frac{1}{n} \mathbb{P}(Y(\lambda_n(k_n)) > k_n) \\ &\leq \int_0^1 \mathbb{P}(Y(\lambda_n(k_n)x) > k_n) dx + \frac{1}{n}(1 - e^{-\lambda_n(k_n)})^{k_n} \\ &= \int_0^1 \left(1 - e^{-\lambda_n(k_n)x}\right)^{k_n} dx + \frac{1}{n}(1 - e^{-\lambda_n(k_n)})^{k_n}. \end{aligned}$$

Now recall that  $\lambda_n(k_n) = -\frac{n-1}{k_n} \log\left(1 - \frac{k_n}{n-1}\right) = 1 + O\left(\frac{k_n}{n}\right)$ , so uniformly in  $x \in [0, 1]$ ,

$$e^{-\lambda_n(k_n)x} = e^{-x} + O\left(\frac{k_n}{n}\right).$$

Since  $k_n = o(\sqrt{n})$ , we have  $k_n/n = o(1/k_n)$  and thus Lemma A.1 from the Appendix gives

$$\int_0^1 \left(1 - e^{-\lambda_n(k_n)x}\right)^{k_n} dx \sim \frac{e}{k_n}(1 - e^{-1})^{k_n+1}.$$

Elementary calculations show that when  $k_n = o(\sqrt{n})$ , we also have

$$\frac{1}{n}(1 - e^{-\lambda_n(k_n)})^{k_n} \sim \frac{1}{n}(1 - e^{-1})^{k_n} = o\left(\frac{(1 - e^{-1})^{k_n}}{k_n}\right).$$

It remains to examine the lower bound in Lemma 5.3. As above, we get an integral

$$\begin{aligned} \mathbb{P}(\Upsilon_n(L) > k_n) &\geq \frac{1}{n} \sum_{\ell=0}^{n-1} \mathbb{P}\left(Y\left(\frac{\ell-k_n+1}{n-1}\right) > k_n\right) \\ &\geq \frac{n-1}{n} \int_0^1 \mathbb{P}\left(Y\left(\frac{\lfloor x(n-1) \rfloor - k_n}{n-1}\right) > k_n\right) dx \\ &\geq \frac{n-1}{n} \int_0^1 \mathbb{P}\left(Y\left(x - \frac{k_n}{n-1}\right) > k_n\right) dx. \end{aligned}$$

Since

$$\mathbb{P}\left(Y\left(x - \frac{k_n}{n-1}\right) > k_n\right) = \left(1 - \exp\left(-x + \frac{k_n}{n-1}\right)\right)^{k_n} = \left(1 - e^{-x} + O(k_n/n)\right)^{k_n},$$

using Lemma A.1 again, we get

$$\mathbb{P}(\Upsilon_n(L) > k_n) \geq \frac{n-1}{n} \int_0^1 \mathbb{P}\left(Y\left(x - \frac{k_n}{n-1}\right) > k_n\right) dx \sim \frac{e}{k_n} (1 - e^{-1})^{k_n+1},$$

which completes the proof.  $\square$

## 5.2 Size of some random trees

In this section, we study the size of some typical trees of  $\mathcal{F}_n$ . In particular, we study the asymptotics of the size  $T_n^{(1)}$  of the tree containing vertex 1 and of the size  $T_n^U$  of a tree sampled uniformly at random among the trees composing  $\mathcal{F}_n$ . Our main result is the following theorem.

### Theorem 5.5.

(i) Let  $T_n^U$  be the size of a uniform tree of  $\mathcal{F}_n$ . Then,

$$\mathbb{P}\left(T_n^U = k\right) \xrightarrow{n \rightarrow \infty} 2 \int_0^1 x e^{-x} (1 - e^{-x})^{k-1} dx,$$

that is,  $T_n^U \xrightarrow{d} T^U$  where  $T^U \sim \text{Geometric}(e^{-X})$ , and  $X \sim 2xdx$  on  $[0, 1]$ .

(ii) Let  $T_n^{(1)}$  be the size of the tree containing vertex 1 in  $\mathcal{F}_n$ . Then,

$$\mathbb{P}\left(T_n^{(1)} = k\right) \xrightarrow{n \rightarrow \infty} k \int_0^1 x e^{-x} (1 - e^{-x})^{k-1} dx,$$

that is,  $T_n^{(1)}$  converges in distribution to the size-biasing of  $T^U$ .

**Remark 5.6.** Note that even though the limit distribution of  $T_n^{(1)}$  is the size-biased limit distribution of  $T_n^U$ , for finite  $n$  the distribution of  $T_n^{(1)}$  is *not* the size-biased distribution of  $T_n^U$ . Indeed, note that

$$\mathbb{P}\left(T_n^{(1)} = k\right) = \mathbb{E}\left(\sum_{\mathcal{T} \in \mathcal{F}_n} \mathbb{1}_{\{|\mathcal{T}|=k\}} \mathbb{1}_{\{1 \in \mathcal{T}\}}\right) = \frac{k}{n} \mathbb{E}\left(\sum_{\mathcal{T} \in \mathcal{F}_n} \mathbb{1}_{\{|\mathcal{T}|=k\}}\right),$$

while

$$\mathbb{P}\left(T_n^U = k\right) = \mathbb{E}\left(\frac{1}{N_n} \sum_{\mathcal{T} \in \mathcal{F}_n} \mathbb{1}_{\{|\mathcal{T}|=k\}}\right),$$

where, as in Section 3,  $N_n$  denotes the number of trees in  $\mathcal{F}_n$ . However, these computations are enough to show that, in the limit, the size-biasing holds: indeed, note that  $n/N_n \rightarrow 2$  in probability by Proposition 3.1. Furthermore, using for instance Hoeffding's inequality [16] to control the deviation of  $N_n$  from its mean, it is easy to show that  $n/N_n \rightarrow 2$  in  $L^1$  as well, so that

$$\left|\mathbb{P}\left(T_n^{(1)} = k\right) - \frac{k}{2} \mathbb{P}\left(T_n^U = k\right)\right| \leq \frac{k}{2} \mathbb{E}\left(\left|\frac{n}{N_n} - 2\right|\right) \rightarrow 0.$$

This shows that points (i) and (ii) of Theorem 5.5 are equivalent.  $\diamond$

We start by giving the distribution of  $T_n^{(1)}$  in terms of the process  $\Upsilon_n$  defined in Section 5.1. For this, we first need to introduce some notation. Let  $\mathcal{F}_n^{(v)}$  be the tree containing vertex  $v$  in  $\mathcal{F}_n$ . We denote by  $H_n^{(v)}$  the number of steps after the root of  $\mathcal{F}_n^{(v)}$  was added in the UA construction. Recalling the notation from Section 2.2, where  $\sigma^{-1}(v) \in \{1, \dots, n\}$  denotes the step of the UA construction at which vertex  $v$  was added, we thus have

$$H_n^{(v)} = n - \min\{\sigma^{-1}(u) : u \in \mathcal{F}_n^{(v)}\}.$$

**Proposition 5.7.** *Let  $T_n^{(1)}$  be the size of the tree containing vertex 1 in  $\mathcal{F}_n$ , and denote by  $H_n^{(1)}$  the number of steps after the root of that tree was added in the UA construction. Then,*

(i) For  $0 \leq h \leq n - 1$ ,  $\mathbb{P}(H_n^{(1)} = h) = \frac{h}{n(n-1)} \left(1 + \frac{1}{n-1}\right)^h$ .

(ii) Conditional on  $\{H_n^{(1)} = h\}$ ,  $T_n^{(1)}$  is distributed as the size-biasing of  $\Upsilon_n(h)$ .

**Remark 5.8.** The size-biasing of  $\Upsilon_n(h)$  can be easily represented as follows. Consider the Markov chain  $\Upsilon_n^* = (\Upsilon_n^*(\ell), 0 \leq \ell \leq n - 1)$  defined by  $\Upsilon_n^*(0) = 1$  and the following transition probabilities:

$$\mathbb{P}(\Upsilon_n^*(\ell + 1) = j \mid \Upsilon_n^*(\ell) = i) = \begin{cases} \frac{i+1}{n} & \text{if } j = i + 1 \\ 1 - \frac{i+1}{n} & \text{if } j = i. \end{cases}$$

A straightforward induction on  $\ell$  shows that  $\Upsilon_n^*(\ell)$  is distributed as the size-biasing of  $\Upsilon_n(\ell)$ .  $\diamond$

*Proof.* First, note that  $H_n^{(1)} = h$  if and only if a new tree is created at step  $n - h$ , and vertex 1 belongs to this tree. Now, the probability that a new tree is created at step  $n - h$  is  $\frac{h}{n-1}$ , and the size of this tree is then distributed as  $\Upsilon_n(h)$ . Moreover, at the end of the UA construction the labels are assigned to the vertices uniformly. As a result, conditional on a tree having size  $i$ , the probability that it contains vertex 1 is  $i/n$ . We thus have

$$\mathbb{P}(H_n^{(1)} = h, T_n^{(1)} = i) = \frac{h}{n-1} \cdot \frac{i}{n} \mathbb{P}(\Upsilon_n(h) = i).$$

Summing over  $i$  and using Lemma 5.1 yields

$$\mathbb{P}(H_n^{(1)} = h) = \frac{h}{n(n-1)} \left(1 + \frac{1}{n-1}\right)^h.$$

Finally,

$$\mathbb{P}(T_n^{(1)} = i \mid H_n^{(1)} = h) = i \mathbb{P}(\Upsilon_n(h) = i) \left(1 + \frac{1}{n-1}\right)^{-h},$$

which concludes the proof.  $\square$

We can now turn to the proof of our main result.



*Proof of Theorem 5.5.* By Remark 5.8, it is sufficient to prove (ii). Now from Proposition 5.7, we can write

$$\begin{aligned}\mathbb{P}(T_n^{(1)} = k) &= \frac{1}{n} \sum_{h=0}^{n-1} \frac{h}{n-1} \mathbb{E}(\Upsilon_n(h) \mathbf{1}_{\{\Upsilon_n(h)=k\}}) \\ &= \frac{k}{n} \sum_{h=0}^{n-1} \frac{h}{n-1} \mathbb{P}(\Upsilon_n(h) = k).\end{aligned}$$

Therefore using Lemma 5.2, we get

$$\frac{k}{n} \sum_{h=0}^{n-1} \frac{h}{n-1} \mathbb{P}(\Upsilon_n(h) = k) \xrightarrow{n \rightarrow \infty} k \int_0^1 x \mathbb{P}(Y(x) = k) dx,$$

and recalling the well-known fact that  $Y(x)$  has a Geometric( $e^{-x}$ ) distribution (see for instance Section 5.3 in [24]) concludes the proof.  $\square$

**Remark 5.9.** If  $H_n^U$  denotes the number of steps in the UA construction after the root of a uniformly chosen tree was added, one could give an alternative proof of Theorem 5.5 by showing that  $H_n^U/n \rightarrow X$ , and then using Lemma 5.2.  $\diamond$

### 5.3 Size of the largest tree

The goal of this section is to derive asymptotics for  $T_n^{\max} := \max_v T_n^{(v)}$ , the size of the largest tree in the Moran forest on  $n$  vertices, when  $n \rightarrow \infty$ . Namely, we show that

$$T_n^{\max} = \alpha \left( \log n - (1 + o_p(1)) \log \log n \right),$$

where  $\alpha = (1 - \log(e-1))^{-1}$ . Similarly to Theorem 1.2 concerning the largest degree, this corresponds to the maximum of  $n/2$  independent  $T_n^U$ -distributed trees. Again the key element to the proof is to control the asymptotic independence of two distinct trees of  $\mathcal{F}_n$ .

As in Section 5.1, for any vertex  $v$  let us define  $\tilde{\mathcal{T}}_n^{(v)} \subset \mathcal{T}_n^{(v)}$  as the subtree descending from  $v$  in the UA construction. For our purpose, it will be sufficient to study the size  $\tilde{T}_n^{(v)} := |\tilde{\mathcal{T}}_n^{(v)}|$  of those subtrees instead of that of the trees  $\mathcal{T}_n^{(v)}$ . Indeed, observe that

$$T_n^{\max} = \max_v \tilde{T}_n^{(v)},$$

so that applying Lemma 4.5 with  $m_n = \alpha \log \log n$  and  $\beta = -1$  to the exchangeable variables  $(\tilde{T}_n^{(1)} - \alpha \log n, \dots, \tilde{T}_n^{(n)} - \alpha \log n)$  will prove the theorem. Again, we omit the superscript and denote by  $\tilde{T}_n$  a random variable with distribution equal to that of  $\tilde{T}_n^{(1)}$ .

For the rest of the section, we thus study the tail probabilities of the variable  $\tilde{T}_n$ . Recall from the UA construction that the number  $L$  of steps after a fixed vertex was added is uniformly distributed on  $\{0, \dots, n-1\}$ , and from Section 5.1 that, conditional on  $\{L = \ell\}$ ,

$$\tilde{T}_n \stackrel{d}{=} \Upsilon_n(\ell).$$

As a consequence, applying directly Proposition 5.4 yields that for any sequence of integers  $k_n \rightarrow \infty$  with  $k_n = o(\sqrt{n})$ ,

$$\mathbb{P}(\tilde{T}_n > k_n) \sim \frac{e}{k_n} (1 - e^{-1})^{k_n+1}. \quad (3)$$

Note that if  $k_n$  is not integer-valued, then

$$\mathbb{P}(\tilde{T}_n > k_n) = \mathbb{P}(\tilde{T}_n > \lfloor k_n \rfloor) \sim \frac{e}{k_n} (1 - e^{-1})^{\lfloor k_n \rfloor + 1},$$

which is not necessarily equivalent to  $\frac{e}{k_n} (1 - e^{-1})^{k_n + 1}$  since  $k_n - \lfloor k_n \rfloor$  may oscillate between 0 and 1. However, we do have  $\mathbb{P}(\tilde{T}_n > k_n) = \Theta((1 - e^{-1})^{k_n} / k_n)$ , where the Bachmann–Landau notation  $u_n = \Theta(v_n)$  indicates that there exist two positive constants  $c$  and  $C$  such that  $cv_n \leq u_n \leq Cv_n$  for  $n$  large enough. This approximation is sufficient for our purpose.

We may now prove Theorem 1.3 using the first and second moment method that we already used for the largest degree.

*Proof of Theorem 1.3.* We apply Lemma 4.5 to the exchangeable variables

$$(X_n^{(1)}, \dots, X_n^{(n)}) = (\tilde{T}_n^{(1)} - \alpha \log n, \dots, \tilde{T}_n^{(n)} - \alpha \log n),$$

with  $m_n = \alpha \log \log n$  and  $\beta = -1$ . The first two points of the lemma are readily checked, since (3) tells us that for

$$\alpha = (1 - \log(e - 1))^{-1} = -(\log(1 - e^{-1}))^{-1}$$

and any  $\gamma > 0$ , we have for  $k_n := \alpha(\log n - \gamma \log \log n)$

$$\mathbb{P}(\tilde{T}_n - \alpha \log n \geq -\gamma \alpha \log \log n) = \mathbb{P}(\tilde{T}_n \geq k_n) = \Theta\left(\frac{(\log n)^{\gamma-1}}{n}\right). \quad (4)$$

Thus, for all  $\varepsilon > 0$ ,

$$(i) \quad n\mathbb{P}(\tilde{T}_n - \alpha \log n \geq (-1 + \varepsilon)\alpha \log \log n) \rightarrow 0.$$

$$(ii) \quad n\mathbb{P}(\tilde{T}_n - \alpha \log n \geq (-1 - \varepsilon)\alpha \log \log n) \rightarrow +\infty.$$

All that remains to check is the third point of the lemma. From now on, we fix  $k_n = \alpha(\log n - (1 + \varepsilon) \log \log n)$  for some  $\varepsilon > 0$ , and for the sake of readability, we set  $R_n := \mathbb{P}(\tilde{T}_n \geq k_n)$ . With this notation, given Remark 4.6 we need to show

$$\mathbb{P}(\tilde{T}_n^{(1)} \geq k_n, \tilde{T}_n^{(2)} \geq k_n) \leq R_n^2 + o(R_n^2). \quad (5)$$

Since this is rather technical, we defer the complete proof to Lemma A.2 in Appendix A.2, and only outline the main ideas of the proof here. As in the study of the largest degree, we prove this by showing that the law of  $\tilde{T}_n^{(2)}$  conditional on  $\{\tilde{T}_n^{(1)} \geq k_n\}$  is close to its unconditional law. We first prove that

$$\mathbb{P}(A_n, \tilde{T}_n^{(1)} \geq k_n, \tilde{T}_n^{(2)} \geq k_n) = o(R_n^2),$$

where  $A_n := \{\tilde{\mathcal{G}}_n^{(2)} \subset \tilde{\mathcal{G}}_n^{(1)}\} \sqcup \{\tilde{\mathcal{G}}_n^{(1)} \subset \tilde{\mathcal{G}}_n^{(2)}\}$  is the event that one of the two vertices 1 and 2 is an ancestor of the other in the UA construction. We then show

$$\mathbb{P}(A_n^c, \tilde{T}_n^{(1)} \geq k_n, \tilde{T}_n^{(2)} \geq k_n) \leq R_n^2 + o(R_n^2),$$

where  $A_n^c$  denotes the complement of  $A_n$ . This is done by showing that, conditional on  $\{\tilde{T}_n^{(1)} = i\}$ , on the event  $A_n^c$  the process counting the number of vertices of the

tree  $\mathcal{F}_n^{(2)}$  in the UA construction behaves as a modified  $\Upsilon_n$  process, which we essentially bound from above by  $\Upsilon_{n-i}$ . Therefore,  $\tilde{T}_n^{(2)}$  can be compared with an independent variable with distribution  $\tilde{T}_{n-i}$ . Finally, we show that

$$\sum_{i \geq k_n} \mathbb{P}(\tilde{T}_n^{(1)} = i) \mathbb{P}(\tilde{T}_{n-i} \geq k_n) \leq R_n^2 + o(R_n^2),$$

thereby proving (5) and concluding the proof of Theorem 1.3.  $\square$

## 6 Concluding comments

### 6.1 Aldous's construction

The UA construction described in Section 2.2 is reminiscent of Aldous's construction [3, Algorithm 2] of a uniform rooted labeled tree as a function of  $n$  uniform random variables on  $\{1, \dots, n\}$  – an algorithm arguably simpler than the variant of Joyal's bijection used in the proof of Theorem 1.1. However, we could not find a simple way to couple both procedures so that the forest obtained from the UA construction coincides with the one obtained by removing decreasing edges in Aldous's uniform tree.

### 6.2 A local limit

The construction of  $\mathcal{F}_n$  from a uniform random tree in Section 3.2 gives us a way to build an infinite forest as a limiting object for the Moran forest, in a weak sense. Let us start by describing the local weak limit of the uniform (rooted) random tree. Recall that the local weak limit of a sequence of random graph  $\mathcal{G}_n$  is a (possibly infinite) random pointed graph  $(\mathcal{G}, u_{\mathcal{G}})$  such that for each finite radius  $r \geq 1$ , the  $r$ -neighborhood of a uniformly chosen vertex  $u_n \in \mathcal{G}_n$  converges in distribution to the ball of radius  $r$  around  $u_{\mathcal{G}}$  in  $\mathcal{G}$ . Consider the following infinite random tree:

- Start from an infinite *spine* of vertices  $u_0, u_1, u_2, \dots$ , with edges  $(u_i \leftarrow u_{i+1})$  between subsequent vertices, directed toward the focal vertex  $u_0$ .
- Let independent Galton-Watson trees with Poisson(1) offspring distribution start from each  $u_i$ , for  $i \geq 0$ , with edges directed from mothers to daughters.

The graph  $\mathcal{T}_{\infty}$  described above is the local weak limit of the random rooted uniform tree on  $n$  labeled vertices [14]. The root is informally placed at the end of the infinite spine.

In order to translate this result to the Moran forest, we need to remove the decreasing edges of  $\mathcal{T}_{\infty}$ . To do so, we equip each vertex  $v$  with an independent uniform variable  $V_v$  on  $[0, 1]$ , that corresponds to the limiting renormalized label of  $v$ . The graph obtained by removing from  $\mathcal{T}_{\infty}$  each edge  $uv$  such that  $V_u > V_v$  can be understood as the local weak limit of the Moran forest.

This construction can be used to derive some limiting results about the local structure of the Moran forest. We can for instance recover the limiting degree of a uniformly chosen vertex: it is clear that the focal vertex  $u_0$  in the construction above has degree

$$D \stackrel{d}{=} \text{Ber}(V_{u_0}) + \text{Poisson}(1 - V_{u_0}),$$

as described in Proposition 4.3. However, global results such as the size of the largest tree or the largest degree cannot be easily derived from this local weak limit.

**Remark 6.1.**

- (i) By definition, the local weak limit should be a.s. connected. Only the random tree that contains  $u_0$  corresponds to the actual local weak limit of the Moran forest. It seems difficult to give meaning to the other trees that are obtained in this procedure. They are connected to the tree containing the focal vertex through the labels  $(V_v)$ . This indicates that the tree that is adjacent to the focal tree is not simply an independent “local limit tree” for another uniformly chosen vertex.
- (ii) Note that the local weak limit can be obtained in a simpler, more direct way. Indeed, Theorem 5.5 gives us the limiting size of the tree containing a uniformly chosen vertex, and we know from the UA construction that, conditional on the number of trees and their sizes, trees in the Moran forest are uniform attachment trees. ◇

### 6.3 Possible extension

An important property of the Moran model is that it is an exchangeable population model: its distribution is invariant under re-labeling of the vertices. General exchangeable population models are known as *Cannings models*, and just as for the Moran model it is possible to associate a forest-valued process to any Cannings model.

A Cannings model is defined from an exchangeable vector  $(\xi(1), \dots, \xi(n))$  of non-negative integers verifying  $\xi(1) + \dots + \xi(n) = n$ . This vector encodes the offspring distribution of a population labeled by  $\{1, \dots, n\}$ . If  $\xi(v) = 0$ , we say that individual  $v$  is dead. Otherwise, it has  $\xi(v) - 1$  children. It is clear that the number of dead individuals is equal to the total number of children in the population. We can thus assign a mother to each dead individual, in such a way that the number of children of each live individual is  $\xi(v) - 1$ .

Starting from any directed graph, we can now define a transition as follows:

1. Draw a vector distributed as  $(\xi(1), \dots, \xi(n))$ .
2. Disconnect each dead vertex from all its neighbors.
3. Assign a mother to each dead vertex, uniformly among all possibilities. For each dead vertex  $v$ , draw an edge from its mother to  $v$ .

This defines a Markov chain on the set of rooted forests. It is not hard to see that, if  $(\xi(1), \dots, \xi(n))$  is the uniform permutation of the vector  $(2, 0, 1, \dots, 1)$ , we recover the Markov chain leading to the Moran forest described in Section 1.1.

Studying these general exchangeable forest processes is not the aim of the current work. In particular, the techniques we used here rely heavily on the construction of the Moran forest described in Section 2.2, which cannot be easily adapted to more general exchangeable forest processes.

## Acknowledgments

We thank Amaury Lambert for initiating the discussion that led us to consider this model and for comments on the first version of this manuscript, and Justin Salez for interesting discussions on the link between the uniform random rooted tree, its local limit and the Moran forest. We are also grateful to the editor for pointing out the similarity between the UA construction and Aldous's algorithm, and to two anonymous referees for helpful comments, including the possibility of adapting the definition of the Moran forest to Cannings models.

This work was for the most part done while JJD worked in CIRB and LPSM.

## References

- [1] The On-Line Encyclopedia of Integer Sequences, published electronically at <http://oeis.org>, 2019.
- [2] M. Aigner and G. M. Ziegler. *Proofs from THE BOOK*. Springer-Verlag Berlin, 6th edition, 2018. doi:[10.1007/978-3-642-00856-6](https://doi.org/10.1007/978-3-642-00856-6).
- [3] D. J. Aldous. The random walk construction of uniform spanning trees and uniform labelled trees. *SIAM Journal on Discrete Mathematics*, 3(4):450–465, Nov. 1990. doi:[10.1137/0403039](https://doi.org/10.1137/0403039).
- [4] K. B. Athreya and S. N. Lahiri. *Measure Theory and Probability Theory*. Springer Science+Business Media, 2006. doi:[10.1007/978-0-387-35434-7](https://doi.org/10.1007/978-0-387-35434-7).
- [5] K. T. Balińska, L. V. Quintas, and J. Szymański. Random recursive forests. *Random Structures & Algorithms*, 5(1):3–12, 1994. doi:[10.1002/rsa.3240050103](https://doi.org/10.1002/rsa.3240050103).
- [6] F. Bergeron, P. Flajolet, and B. Salvy. Varieties of increasing trees. In *CAAP '92*, pages 24–48, Berlin, Heidelberg, 1992. Springer Berlin Heidelberg. doi:[10.1007/3-540-55251-0\\_2](https://doi.org/10.1007/3-540-55251-0_2).
- [7] P. Billingsley. *Convergence of Probability Measures*. Wiley, 2nd edition, 1999.
- [8] B. Drake. *An Inversion Theorem for Labeled Trees and Some Limits of Areas Under Lattice Paths*. PhD thesis, Brandeis University, 2008.
- [9] M. Drmota. *Random Trees: An Interplay Between Combinatorics and Probability*. Springer-Verlag Vienna, 1st edition, 2009. doi:[10.1007/978-3-211-75357-6](https://doi.org/10.1007/978-3-211-75357-6).
- [10] R. Durrett. *Probability Models for DNA Sequence Evolution*. Springer-Verlag New York, 2nd edition, 2008. doi:[10.1007/978-1-4757-6285-3](https://doi.org/10.1007/978-1-4757-6285-3).
- [11] Ö. Eğecioğlu and J. B. Remmel. Bijections for Cayley trees, spanning trees, and their q-analogues. *Journal of Combinatorial Theory, Series A*, 42(1): 15–30, 1986. doi:[10.1016/0097-3165\(86\)90004-X](https://doi.org/10.1016/0097-3165(86)90004-X).
- [12] A. Etheridge. *Some Mathematical Models From Population Genetics. École d'Été de Probabilités de Saint-Flour XXXIX-2009*, volume 2012. Springer Science+Business Media, 2011. doi:[10.1007/978-3-642-16632-7](https://doi.org/10.1007/978-3-642-16632-7).

- [13] I. M. Gessel and S. Seo. A refinement of Cayley’s formula for trees. *The Electronic Journal of Combinatorics*, 11(2):R27, 2006.
- [14] G. R. Grimmett. Random labelled trees and their branching networks. *Journal of the Australian Mathematical Society*, 30(2):229–237, Dec. 1980. doi:[10.1017/S1446788700016517](https://doi.org/10.1017/S1446788700016517).
- [15] J. Hey. Using phylogenetic trees to study speciation and extinction. *Evolution*, 46(3):627–640, 1992. doi:[10.1111/j.1558-5646.1992.tb02071.x](https://doi.org/10.1111/j.1558-5646.1992.tb02071.x).
- [16] W. Hoeffding. Probability inequalities for sums of bounded random variables. *Journal of the American Statistical Association*, 58(301):13–30, Mar. 1963. doi:[10.1080/01621459.1963.10500830](https://doi.org/10.1080/01621459.1963.10500830).
- [17] A. Joyal. Une théorie combinatoire des séries formelles. *Advances in Mathematics*, 42(1):1–82, 1981. doi:[10.1016/0001-8708\(81\)90052-9](https://doi.org/10.1016/0001-8708(81)90052-9).
- [18] H. M. Mahmoud and R. T. Smythe. On the distribution of leaves in rooted subtrees of recursive trees. *The Annals of Applied Probability*, 1(3):406–418, 1991. doi:[10.1214/aoap/1177005874](https://doi.org/10.1214/aoap/1177005874).
- [19] A. Meir and J. W. Moon. Cutting down recursive trees. *Mathematical Biosciences*, 21(3):173–181, 1974. doi:[10.1016/0025-5564\(74\)90013-3](https://doi.org/10.1016/0025-5564(74)90013-3).
- [20] P. A. P. Moran. Random processes in genetics. *Mathematical Proceedings of the Cambridge Philosophical Society*, 54(1):60–71, 1958. doi:[10.1017/S0305004100033193](https://doi.org/10.1017/S0305004100033193).
- [21] H. Morlon, M. D. Potts, and J. B. Plotkin. Inferring the dynamics of diversification: A coalescent approach. *PLOS Biology*, 8(9):1–13, 2010. doi:[10.1371/journal.pbio.1000493](https://doi.org/10.1371/journal.pbio.1000493).
- [22] M. A. Nowak. *Evolutionary Dynamics: Exploring the Equations of Life*. Harvard University Press, 2006. doi:[10.2307/j.ctvjghw98](https://doi.org/10.2307/j.ctvjghw98).
- [23] J. Pitman. *Combinatorial stochastic processes. École d’Été de Probabilités de Saint-Flour XXXII-2002*, volume 1875. Springer Science+Business Media, 2006. doi:[10.1007/b11601500](https://doi.org/10.1007/b11601500).
- [24] S. M. Ross. *Stochastic Processes*. Wiley, 2nd edition, 1995.

# A Appendix

## A.1 Proof of point (ii) of Proposition 4.4

We want to prove that, for all  $K_n = o(\sqrt{n})$ , there exists  $\varepsilon_n = o(1)$  such that, for all  $k \leq K_n$ ,

$$|\mathbb{P}(D_n \geq k) - \mathbb{P}(D \geq k)| \leq \varepsilon_n \mathbb{P}(D \geq k).$$

Doing this directly from the expressions of  $D_n$  and  $D$  involves unappealing calculations. To somewhat circumvent this, we make use of the simple expressions of the probability generating functions  $G_{D_n}$  and  $G_D$ . For this, let

$$\Delta_n(z) := \sum_{i \geq 0} (\mathbb{P}(D \geq i) - \mathbb{P}(D_n \geq i)) z^i,$$

so that the  $k$ -th derivative of  $\Delta_n$  evaluated at  $z = 0$  is

$$\Delta_n^{(k)}(0) = k! (\mathbb{P}(D \geq k) - \mathbb{P}(D_n \geq k)).$$

Since  $\mathbb{P}(D \geq k) \geq \frac{2/e}{(k+1)!}$ , we have to show that for any given sequence  $K_n = o(\sqrt{n})$ ,

$$\Delta_n^{(k)}(0) = \frac{\varepsilon_n}{k+1}$$

for some  $\varepsilon_n \rightarrow 0$  and all  $k \leq K_n$ . Now, since for any non-negative integer-valued random variable  $X$ ,

$$\sum_{i \geq 0} \mathbb{P}(X \geq i) z^i = \frac{z \mathbb{E}(z^X) - 1}{z - 1},$$

we can express  $\Delta_n$  in terms of the generating functions of  $D$  and  $D_n$ , that is,

$$\Delta_n(z) = \left(1 + \frac{1}{z-1}\right) (G_D(z) - G_{D_n}(z)).$$

Moreover, we know from Proposition 4.3 that

$$G_D(z) = 2 \frac{e^{z-1} - 1}{z-1} - 1 = 2 \sum_{i \geq 0} \frac{(z-1)^i}{(i+1)!} - 1$$

and from Proposition 4.1 that

$$\begin{aligned} G_{D_n}(z) &= 2 \left(1 - \frac{1}{n}\right) \frac{\left(1 + \frac{z-1}{n-1}\right)^n - 1}{z-1} - 1 \\ &= 2 \left(1 - \frac{1}{n}\right) \sum_{i=0}^{n-1} \binom{n}{i+1} \left(\frac{1}{n-1}\right)^{i+1} (z-1)^i - 1 \\ &= 2 \sum_{i=0}^{n-1} \left(\prod_{\ell=1}^i \frac{n-\ell}{n-1}\right) \frac{(z-1)^i}{(i+1)!} - 1, \end{aligned}$$

where the empty product is 1. Therefore,

$$G_D(z) - G_{D_n}(z) = \sum_{i \geq 0} A(n, i) \frac{(z-1)^i}{(i+1)!},$$

where

$$A(n, i) = 2 \left[ 1 - \left( \prod_{\ell=1}^i \frac{n-\ell}{n-1} \right) \mathbf{1}_{\{i \leq n-1\}} \right].$$

Using that  $A(n, 0) = A(n, 1) = 0$  and rearranging a bit, we obtain the following expansion of  $\Delta_n$  at  $z = 1$ :

$$\Delta_n(z) = \sum_{i \geq 1} \left( A(n, i) + \frac{A(n, i+1)}{i+2} \right) \frac{(z-1)^i}{(i+1)!},$$

from which we get

$$\Delta_n^{(k)}(0) = \sum_{i \geq k} \left( A(n, i) + \frac{A(n, i+1)}{i+2} \right) \frac{(-1)^{i-k}}{(i-k)!(i+1)}.$$

Now, pick any  $J_n = o(\sqrt{n})$  such that  $K_n = o(J_n)$ . For all  $i < J_n$ ,

$$\left| A(n, i) + \frac{A(n, i+1)}{i+2} \right| \leq 4 \left( 1 - \prod_{\ell=1}^{J_n} \frac{n-\ell}{n-1} \right) = \varepsilon_n,$$

with  $\varepsilon_n \rightarrow 0$ , since

$$\prod_{\ell=1}^{J_n} \frac{n-\ell}{n-1} \geq \left( \frac{n-J_n}{n-1} \right)^{J_n} = \exp\left(-\frac{J_n^2}{n} + o\left(\frac{J_n^2}{n}\right)\right).$$

For  $i \geq J_n$ , we have

$$\left| A(n, i) + \frac{A(n, i+1)}{i+2} \right| \leq 4.$$

Combining these two upper bounds, we get

$$\begin{aligned} |\Delta_n^{(k)}(0)| &\leq \sum_{i=k}^{J_n-1} \frac{\varepsilon_n}{(i-k)!(i+1)} + \sum_{i \geq J_n} \frac{4}{(i-k)!(i+1)} \\ &\leq \frac{\varepsilon_n C_1}{(k+1)} + \frac{C_2}{(J_n+1)}. \end{aligned}$$

Finally, since  $K_n = o(J_n)$ , we have for all  $k \leq K_n$ ,

$$\frac{1}{J_n+1} \leq \frac{1}{k+1} \cdot \frac{K_n+1}{J_n+1},$$

with  $(K_n+1)/(J_n+1) = o(1)$ . This concludes the proof.

Note that although we have been quite crude in that we have used the triangle inequality on an alternating series, a more careful analysis would show that the  $o(\sqrt{n})$  requirement on  $K_n$  is in fact optimal.

## A.2 Technical lemmas used in the proof of Theorem 1.3

**Lemma A.1.** *For any sequence  $k_n \rightarrow \infty$  and any sequence of measurable maps  $f_n : [0, 1] \rightarrow \mathbb{R}$  such that for all  $x \in [0, 1]$ ,  $(1 - e^{-x} + f_n(x)) \geq 0$  and  $\sup_x |f_n(x)| = o(1/k_n)$ , we have*

$$\int_0^1 (1 - e^{-x} + f_n(x))^{k_n} dx \sim \frac{e}{k_n} (1 - e^{-1})^{k_n+1}.$$



*Proof.* Let us compute

$$\begin{aligned} \int_0^1 \frac{(1 - e^{-x} + f_n(x))^{k_n}}{(1 - e^{-1})^{k_n}} k_n dx &= \int_0^1 \left( 1 - \frac{e^{1-x} - 1}{e - 1} + \frac{e}{e - 1} f_n(x) \right)^{k_n} k_n dx \\ &= \int_0^{k_n} \left( 1 - \frac{y}{k_n} + g_n(y) \right)^{k_n} \frac{e - 1}{1 + (e - 1) \frac{y}{k_n}} dy, \end{aligned}$$

where we used the change of variable  $y = k_n(e^{1-x} - 1)(e - 1)^{-1}$ , and defined the map  $g_n$  as

$$g_n(y) = \frac{e}{e - 1} f_n \left( 1 - \log \left( 1 + \frac{y}{k_n} (e - 1) \right) \right).$$

Since  $(1 - \frac{y}{k_n} + g_n(y))^{k_n} \leq \exp(-y + \frac{e}{e-1} k_n \sup_x f_n(x))$ , it follows from dominated convergence that

$$\int_0^1 \frac{(1 - e^{-x} + f_n(x))^{k_n}}{(1 - e^{-1})^{k_n}} k_n dx \xrightarrow{n \rightarrow \infty} \int_0^\infty e^{-y} (e - 1) dy = e - 1,$$

concluding the proof.  $\square$

**Lemma A.2.** Let  $\tilde{T}_n^{(v)}$  denote the size of the subtree descending from  $v$  in the UA construction of  $\tilde{\mathcal{F}}_n$ . Then, for  $\alpha = -1/\log(1 - e^{-1})$  and any  $\varepsilon > 0$ , letting  $k_n = \alpha(\log n - (1 + \varepsilon)\log \log n)$  and  $R_n = \mathbb{P}(\tilde{T}_n \geq k_n)$ ,

$$\mathbb{P}(\tilde{T}_n^{(1)} \geq k_n, \tilde{T}_n^{(2)} \geq k_n) \leq R_n^2 + o(R_n^2).$$

*Proof.* Let us denote by  $A_n := \{\tilde{\mathcal{G}}_n^{(2)} \subset \tilde{\mathcal{G}}_n^{(1)}\} \sqcup \{\tilde{\mathcal{G}}_n^{(1)} \subset \tilde{\mathcal{G}}_n^{(2)}\}$  the event that one of the vertices 1 and 2 is an ancestor of the other. We start by showing that

$$\mathbb{P}(A_n, \tilde{T}_n^{(1)} \geq k_n, \tilde{T}_n^{(2)} \geq k_n) = o(R_n^2). \quad (6)$$

By exchangeability, we have

$$\begin{aligned} &\mathbb{P}(A_n, \tilde{T}_n^{(1)} \geq k_n, \tilde{T}_n^{(2)} \geq k_n) \\ &= 2 \mathbb{P}(\tilde{\mathcal{G}}_n^{(2)} \subset \tilde{\mathcal{G}}_n^{(1)}, \tilde{T}_n^{(1)} \geq k_n, \tilde{T}_n^{(2)} \geq k_n) \\ &= \sum_{i \geq k_n} \mathbb{P}(\tilde{\mathcal{G}}_n^{(2)} \subset \tilde{\mathcal{G}}_n^{(1)}, \tilde{T}_n^{(2)} \geq k_n \mid \tilde{T}_n^{(1)} = i) \mathbb{P}(\tilde{T}_n = i). \end{aligned}$$

Let us call the *height* of a vertex the number of steps after it was added in the UA construction. Conditional on  $\{\tilde{T}_n^{(1)} = i\}$  and on the heights of the vertices of  $\tilde{\mathcal{G}}_n^{(1)}$  being  $\ell_1 > \dots > \ell_i$ , the height  $L_2$  of vertex 2 is uniformly distributed on  $\{0, \dots, n - 1\} \setminus \{\ell_1\}$ . Moreover, in order to have

$$\{\tilde{\mathcal{G}}_n^{(2)} \subset \tilde{\mathcal{G}}_n^{(1)}, \tilde{T}_n^{(2)} \geq k_n\},$$

the height of vertex 2 must belong to  $\{\ell_2, \dots, \ell_{i-(k_n-1)}\}$ , which happens with probability  $\frac{i-k_n}{n-1}$ . Therefore,

$$\begin{aligned} &\mathbb{P}(A_n, \tilde{T}_n^{(1)} \geq k_n, \tilde{T}_n^{(2)} \geq k_n) \\ &\leq \sum_{i \geq k_n} \mathbb{P}(\tilde{T}_n = i) \frac{i - k_n}{n - 1} \\ &= \frac{1}{n - 1} \sum_{i \geq k_n} \mathbb{P}(\tilde{T}_n \geq i). \end{aligned}$$

To show that this is small enough, we let  $K_n := k_n + \alpha(\log n)^\delta$  with  $0 < \delta < \min(1, \varepsilon)$  and  $K'_n := 2\alpha \log n$ , and crudely bound

$$\sum_{i > k_n} \mathbb{P}(\tilde{T}_n \geq i) \leq (K_n - k_n) \mathbb{P}(\tilde{T}_n \geq k_n) + K'_n \mathbb{P}(\tilde{T}_n \geq K_n) + n \mathbb{P}(\tilde{T}_n \geq K'_n).$$

Now let us show that these three terms are negligible compared to  $nR_n^2$ . Recalling from (4) that  $R_n = \Theta\left(\frac{(\log n)^\varepsilon}{n}\right)$ , we have  $nR_n^2 = \Theta((\log n)^{2\varepsilon}/n)$  and therefore

- $(K_n - k_n) \mathbb{P}(\tilde{T}_n \geq k_n) \sim \alpha(\log n)^\delta R_n = \Theta\left(\frac{(\log n)^{\delta+\varepsilon}}{n}\right) = o(nR_n^2)$ .
- $K'_n \mathbb{P}(\tilde{T}_n \geq K_n) = \Theta(\log n R_n e^{-(\log n)^\delta}) = o(R_n) = o(nR_n^2)$ .
- $n \mathbb{P}(\tilde{T}_n \geq K'_n) = \Theta\left(n \frac{n^{-2}}{\log n}\right) = o(1/n) = o(nR_n^2)$ .

As a result, (6) is proven and it remains to show that

$$\mathbb{P}(A_n^c, \tilde{T}_n^{(1)} \geq k_n, \tilde{T}_n^{(2)} \geq k_n) \leq R_n^2 + o(R_n^2),$$

where  $A_n^c$  denotes the complement of  $A_n$ . We now fix  $n \geq 1$ ,  $i \geq k_n$ , and a finite sequence  $n - 1 \geq \ell_1 > \dots > \ell_i \geq 0$ . Let us write  $B$  for the event that  $\tilde{\mathcal{T}}_n^{(1)}$  contains exactly the vertices with heights  $\ell_1 > \dots > \ell_i$ . Conditional on  $B$ , let us examine the distribution of  $\tilde{\mathcal{T}}_n^{(2)}$ . Recall that the height  $L_2$  of vertex 2 is uniformly distributed on  $\{0, \dots, n - 1\} \setminus \{\ell_1\}$ . In the UA construction, define  $\mathbb{T}$  as the tree obtained by starting from a root arrived at height  $L_2$  and allowing the attachment of a vertex with height  $\ell$  to  $\mathbb{T}$  only if  $\ell \notin \{\ell_1, \dots, \ell_i\}$ . Then, on the event  $A_n^c$ , this tree must coincide with  $\mathcal{T}_n^{(2)}$ , and so

$$\mathbb{P}(A_n^c, \tilde{T}_n^{(2)} \geq k_n \mid B) = \mathbb{P}(A_n^c, |\mathbb{T}| \geq k_n \mid B).$$

From the UA construction, for any  $\ell \notin \{\ell_1, \dots, \ell_i\}$ , conditional on  $B \cap \{L_2 = \ell\}$ , we can describe  $|\mathbb{T}|$  using the process  $(\tilde{\Upsilon}_\ell(m), 0 \leq m \leq \ell)$  defined by

- $\tilde{\Upsilon}_\ell(0) = 1$ .
- For all  $0 < m \leq \ell$ ,  $\tilde{\Upsilon}_\ell(m) - \tilde{\Upsilon}_\ell(m - 1) \in \{0, 1\}$  and, conditional on  $\{\tilde{\Upsilon}_\ell(m - 1) = j\}$ ,  $\tilde{\Upsilon}_\ell(m) = j + 1$  with probability

$$\begin{cases} \frac{j}{n - 1 - J_m} & \text{if } \ell - m \notin \{\ell_1, \dots, \ell_i\} \\ 0 & \text{if } \ell - m \in \{\ell_1, \dots, \ell_i\}, \end{cases}$$

where  $J_m = |\{\ell_1, \dots, \ell_i\} \cap \{\ell - m, \dots, n\}|$  is the number of vertices of  $\tilde{T}_n^{(1)}$  with height greater than  $\ell - m$  in the UA construction.

With this definition, for any  $\ell \notin \{\ell_1, \dots, \ell_i\}$ , conditional on  $B \cap \{L_2 = \ell\}$ , we have by construction  $|\mathbb{T}| \stackrel{d}{=} \tilde{\Upsilon}_\ell(\ell)$ . Now, note that the probability of increasing is always bounded by  $j/(n - 1 - i)$ . Therefore,  $\tilde{\Upsilon}_\ell$  can be coupled with  $\Upsilon_{n-i}$  in such a way that, for all  $0 \leq m \leq \ell < n - i$ ,

$$\tilde{\Upsilon}_\ell(m) \leq \Upsilon_{n-i}(m).$$

For  $\ell \geq n - i$ , we use the crude bound  $\mathbb{P}(\tilde{\Upsilon}_\ell(\ell) \geq k_n) \leq \mathbb{E}(\tilde{\Upsilon}_\ell(\ell))/k_n$ . Using the same reasoning as in Lemma 5.1, we get

$$\mathbb{E}(\tilde{\Upsilon}_\ell(\ell)) \leq \left(1 + \frac{1}{n - i - 1}\right)^{n-i-1} \leq e.$$

We thus have

$$\begin{aligned} \mathbb{P}(A_n^c, |\mathbb{T}| \geq k_n \mid B) &\leq \mathbb{P}(L_2 \notin \{\ell_1, \dots, \ell_i\}, |\mathbb{T}| \geq k_n \mid B) \\ &= \frac{1}{n-1} \sum_{\substack{\ell=0 \\ \ell \notin \{\ell_1, \dots, \ell_i\}}}^{n-1} \mathbb{P}(\tilde{\Upsilon}_\ell(\ell) \geq k_n), \\ &\leq \frac{ei}{k_n(n-1)} + \frac{1}{n-1} \sum_{\ell=0}^{n-i-1} \mathbb{P}(\Upsilon_{n-i}(\ell) \geq k_n) \end{aligned} \quad (7)$$

$$= \frac{ei}{k_n(n-1)} + \frac{n-i}{n-1} \mathbb{P}(\tilde{T}_{n-i} \geq k_n). \quad (8)$$

Since this bound depends on the set  $\{\ell_1, \dots, \ell_i\}$  only via its cardinality  $i$ , one can integrate with respect to the distribution of  $\mathcal{F}_n^{(1)}$  to get

$$\mathbb{P}(A_n^c, |\mathbb{T}| \geq k_n \mid \tilde{T}_n^{(1)} = i) \leq \frac{ei}{k_n(n-1)} + \frac{n-i}{n-1} \mathbb{P}(\tilde{T}_{n-i} \geq k_n).$$

Finally, because  $\Upsilon_{n-(i+1)}(\ell) \stackrel{d}{\geq} \Upsilon_{n-i}(\ell)$ , the expression (7) – and therefore (8) – is nondecreasing in  $i$ , and we have

$$\begin{aligned} &\mathbb{P}(A_n^c, \tilde{T}_n^{(1)} \geq k_n, \tilde{T}_n^{(2)} \geq k_n) \\ &\leq \sum_{i \geq k_n} \mathbb{P}(\tilde{T}_n = i) \left( \frac{ei}{k_n(n-1)} + \frac{n-i}{n-1} \mathbb{P}(\tilde{T}_{n-i} \geq k_n) \right) \\ &\leq \sum_{i=k_n}^{K_n} \mathbb{P}(\tilde{T}_n = i) \left( \frac{eK_n}{k_n(n-1)} + \frac{n-K_n}{n-1} \mathbb{P}(\tilde{T}_{n-K_n} \geq k_n) \right) \end{aligned} \quad (9)$$

$$+ \sum_{i \geq K_n} \mathbb{P}(\tilde{T}_n = i) \left( \frac{ei}{k_n(n-1)} + \frac{n-i}{n-1} \mathbb{P}(\tilde{T}_{n-i} \geq k_n) \right), \quad (10)$$

for any sequence  $K_n \geq k_n$ . Letting  $K_n := \alpha(\log n)^{1+\varepsilon/2}$ , we then show that (9) is asymptotically no greater than  $R_n^2$ , and that (10) is negligible compared to  $R_n^2$ . Indeed, (9) is bounded from above by

$$R_n \left( \frac{eK_n}{k_n(n-1)} + \frac{n-K_n}{n-1} \mathbb{P}(\tilde{T}_{n-K_n} \geq k_n) \right).$$

Now note that  $\frac{eK_n}{k_n(n-1)} = O\left(\frac{(\log n)^{\varepsilon/2}}{n}\right) = o(R_n)$ , and that since  $n - K_n \sim n$ , we have  $k_n = o(\sqrt{n - K_n})$ . Therefore, using (3) we get  $\mathbb{P}(\tilde{T}_{n-K_n} \geq k_n) \sim R_n$ . Finally, up to a multiplicative constant, (10) is bounded from above by

$$\mathbb{P}(\tilde{T}_n \geq K_n) = \Theta\left(\frac{n^{-(\log n)^{\varepsilon/2}}}{K_n}\right) = o(n^{-2}) = o(R_n^2).$$

Putting everything together, we have proved that

$$\mathbb{P}(\tilde{T}_n^{(1)} \geq k_n, \tilde{T}_n^{(2)} \geq k_n) \leq R_n^2 + o(R_n^2),$$

which concludes the proof.  $\square$