



**HAL**  
open science

# Exact simulation of the genealogical tree for a stationary branching population and application to the asymptotics of its total length

Romain Abraham, Jean-François Delmas

## ► To cite this version:

Romain Abraham, Jean-François Delmas. Exact simulation of the genealogical tree for a stationary branching population and application to the asymptotics of its total length. 2018. hal-01413614v2

**HAL Id: hal-01413614**

**<https://hal.science/hal-01413614v2>**

Preprint submitted on 18 Apr 2018 (v2), last revised 30 Jan 2020 (v3)

**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.

# EXACT SIMULATION OF THE GENEALOGICAL TREE FOR A STATIONARY BRANCHING POPULATION AND APPLICATION TO THE ASYMPTOTICS OF ITS TOTAL LENGTH

ROMAIN ABRAHAM AND JEAN-FRANÇOIS DELMAS

ABSTRACT. We consider a model of stationary population with random size given by a continuous state branching process with immigration with a quadratic branching mechanism. We give an exact elementary simulation procedure of the genealogical tree of  $n$  individuals randomly chosen among the extant population at a given time. Then, we prove the convergence of the renormalized total length of this genealogical tree as  $n$  goes to infinity, see also Pfaffelhuber, Wakolbinger and Weisshaupt (2011) in the context of a constant size population. The limit appears already in Bi and Delmas (2016) but with a different approximation of the full genealogical tree. The proof is based on the ancestral process of the extant population at a fixed time which was defined by Aldous and Popovic (2005) in the critical case.

## 1. INTRODUCTION

Continuous state branching (CB) processes are stochastic processes that can be obtained as the scaling limits of sequences of Galton-Watson processes when the initial number of individuals tends to infinity. They hence can be seen as a model for a large branching population. The genealogical structure of a CB process can be described by a continuum random tree introduced first by Aldous [4] for the quadratic critical case, see also Le Gall and Le Jan [19] and Duquesne and Le Gall [12] for the general critical and sub-critical cases. We shall only consider the quadratic case; it is characterized by a branching mechanism  $\psi_\theta$ :

$$\psi_\theta(\lambda) = \beta\lambda^2 + 2\beta\theta\lambda, \quad \lambda \in [0, +\infty),$$

where  $\beta > 0$  and  $\theta \in \mathbb{R}$ . The sub-critical (resp. critical) case corresponds to  $\theta > 0$  (resp.  $\theta = 0$ ). The parameter  $\beta$  can be seen as a time scaling parameter, and  $\theta$  as a population size parameter.

In this model the population dies out a.s. in the critical and sub-critical cases. In order to model branching population with stationary size distribution, which corresponds to what is observed at an ecological equilibrium, one can simply condition a sub-critical or a critical CB to not die out. This gives a Q-process, see Roelly-Coppoleta and Rouault [22], Lambert [18] and Abraham and Delmas [1], which can also be viewed as a CB with a specific immigration. The genealogical structure of the Q-process in the stationary regime is a tree with an infinite spine. This infinite spine has to be removed if one adopts the immigration point of view, in this case the genealogical structure can be seen as a forest of trees. For  $\theta > 0$ , let  $(Z_t, t \in \mathbb{R})$  be this Q-process in the stationary regime, so that  $Z_t$  is the size of the population at time  $t \in \mathbb{R}$ . See Chen and Delmas [9] for studies on this model in a more general framework. Let  $A_t$  be the time to the most recent common ancestor of the population living at time  $t$ , see (16) for a precise definition. According to [9], we have  $\mathbb{E}[Z_t] = 1/\theta$ , and  $\mathbb{E}[A_t] = 3/4\beta\theta$ , so that  $\theta$  is indeed a population size

---

*Date:* April 18, 2018.

*2010 Mathematics Subject Classification.* 60J80,60J85.

*Key words and phrases.* Stationary branching processes, Real trees, Genealogical trees, Ancestral process, Simulation.

parameter and  $\beta$  is a time parameter.

Aldous and Popovic [5], see also Popovic [21], give a description of the genealogical tree of the extant population at a fixed time using the so-called ancestral process which is a point process representation of the lineage in a setting very close to  $\theta = 0$  in the present model. We extend the presentation of [5] to the case  $\theta \geq 0$ , see Propositions 3.6 and 3.7 which can be summarized as follows. According to [9], the size of the population at time 0,  $Z_0$ , is distributed as  $E_g + E_d$ , where  $E_g$  and  $E_d$  are two independent exponential random variables with mean  $1/2\theta$  (take  $E_g = E_d = +\infty$  if  $\theta = 0$ ). Conditionally given  $(E_g, E_d)$ , we describe the genealogy of the extant population by a Poisson point measure (that we call the ancestral process), namely  $\mathcal{A}(du, d\zeta) = \sum_{i \in \mathcal{I}} \delta_{u_i, \zeta_i}(du, d\zeta)$  where  $u_i$  represents the individual and  $\zeta_i$  its 'age'. From this point measure, we construct informally a genealogical tree as follows. We view this process as a sequence of vertical segments in  $\mathbb{R}^2$ , the tops of the segments being the  $u_i$ 's and their lengths being the  $\zeta_i$ 's. We add the half line  $\{0\} \times (-\infty, 0]$  in this collection of segments. We then attach the bottom of each segment such that  $u_i > 0$  (resp.  $u_i < 0$ ) to the first left (resp. first right) longer segment. See Figure 1 for an example.

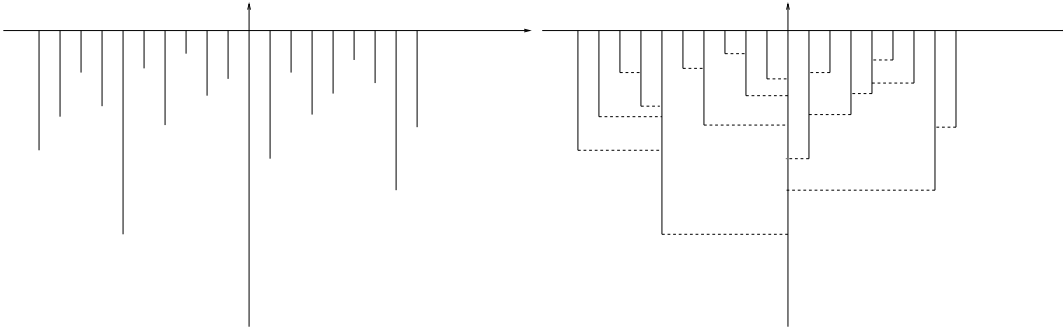


FIGURE 1. An instance of an ancestral process and the corresponding genealogical tree

The ancestral process description allows to give elementary exact simulations of the genealogical tree of  $n$  individuals randomly chosen in the extant population at time 0 (or at some time  $t \in \mathbb{R}$  as the population has a stationary distribution). We give first a static simulation for fixed  $n$  in Subsection 4.1, then two dynamic simulations in Subsections 4.2 and 4.3, where the individuals are taken one by one and the genealogical tree is then updated. Our framework allows also to simulate the genealogical tree of  $n$  extant individuals conditionally given the time  $A_0$  to the most recent common ancestor of the extant population, see Subsection 4.4. Let us stress that the existence of an elementary simulation method is new, and the question goes back to Lambert [17] and Theorem 4.7 in [9]

The ancestral process description allows also to compute the limit distribution of the total length of the genealogical tree of the extant population at time  $t \in \mathbb{R}$ . More precisely, let  $\Lambda_n$  be the total length of the tree of  $n$  individuals randomly chosen in the extant population at time  $t$ , see (18) for a precise definition. Then we prove, see Proposition 5.1 and (34), that  $(\Lambda_n - \mathbb{E}[\Lambda_n | Z_t], n \in \mathbb{N}^*)$  converges a.s. and in  $L^2$  as  $n$  goes down to 0 towards a limit, say  $\mathcal{L}_t$ .

The Laplace transform of the distribution of  $\mathcal{L}_t$  is given by, for  $\lambda > 0$ :

$$\mathbb{E} \left[ e^{-\lambda \mathcal{L}_t} | Z_t \right] = e^{\theta Z_t \varphi(\lambda/(2\beta\theta))}, \quad \text{with} \quad \varphi(\lambda) = \lambda \int_0^1 \frac{1-v^\lambda}{1-v} dv.$$

The proof is based on technical  $L^2$  computations. This result is in the spirit of Pfaffelhuber, Wakolbinger and Weisshaupt [20] on the tree length of the coalescent, which is a model for constant population size. We also prove that  $\mathcal{L}_t$  coincides with the limit of the total length  $L_\varepsilon$  of the genealogical tree up to  $t - \varepsilon$  of the individuals alive at time  $t$  obtained in [6]. More precisely, we have that  $(L_\varepsilon - \mathbb{E}[L_\varepsilon | Z_t], \varepsilon > 0)$  converges a.s. towards  $\mathcal{L}_t$  as  $\varepsilon$  goes down to zero. See [6] for some properties of the process  $(\mathcal{L}_t, t \in \mathbb{R})$ .

The paper is organized as follows. We first introduce in Section 2 the framework of real trees and we define the Brownian CRT that describes the genealogy of the CB in the quadratic case. Section 3 is devoted to the description via a Poisson point measure of the ancestral process of the extant population at time 0 and Section 4 gives the different simulations of the genealogical tree of  $n$  individuals randomly chosen in this population. Then, Section 5 concerns the asymptotic length of the genealogical tree for those  $n$  sampled individuals.

## 2. NOTATIONS

**2.1. Real trees.** The study of real trees has been motivated by algebraic and geometric purposes. See in particular the survey [10]. It has been first used in [15] to study random continuum trees, see also [14].

**Definition 2.1** (Real tree). *A real tree is a metric space  $(\mathbf{t}, d_{\mathbf{t}})$  such that:*

- (i) *For every  $x, y \in \mathbf{t}$ , there is a unique isometric map  $f_{x,y}$  from  $[0, d_{\mathbf{t}}(x, y)]$  to  $\mathbf{t}$  such that  $f_{x,y}(0) = x$  and  $f_{x,y}(d_{\mathbf{t}}(x, y)) = y$ .*
- (ii) *For every  $x, y \in \mathbf{t}$ , if  $\phi$  is a continuous injective map from  $[0, 1]$  to  $\mathbf{t}$  such that  $\phi(0) = x$  and  $\phi(1) = y$ , then  $\phi([0, 1]) = f_{x,y}([0, d_{\mathbf{t}}(x, y)])$ .*

Notice that a real tree is a length space as defined in [8]. We say that  $(\mathbf{t}, d_{\mathbf{t}}, \partial_{\mathbf{t}})$  is a *rooted* real tree, where  $\partial = \partial_{\mathbf{t}}$  is a distinguished vertex of  $\mathbf{t}$ , which will be called the root. Remark that the set  $\{\partial\}$  is a rooted tree that only contains the root.

Let  $\mathbf{t}$  be a compact rooted real tree and let  $x, y \in \mathbf{t}$ . We denote by  $\llbracket x, y \rrbracket$  the range of the map  $f_{x,y}$  described in Definition 2.1. We also set  $\llbracket x, y \llbracket = \llbracket x, y \rrbracket \setminus \{y\}$ . We define the out-degree of  $x$ , denoted by  $k_{\mathbf{t}}(x)$ , as the number of connected components of  $\mathbf{t} \setminus \{x\}$  that do not contain the root. If  $k_{\mathbf{t}}(x) = 0$ , resp.  $k_{\mathbf{t}}(x) > 1$ , then  $x$  is called a leaf, resp. a branching point. A tree is said to be binary if the out-degree of its vertices belongs to  $\{0, 1, 2\}$ . The skeleton of the tree  $\mathbf{t}$  is the set  $\text{sk}(\mathbf{t})$  of points of  $\mathbf{t}$  that are not leaves. Notice that  $\text{cl}(\text{sk}(\mathbf{t})) = \mathbf{t}$ , where  $\text{cl}(A)$  denote the closure of  $A$ .

We denote by  $\mathbf{t}_x$  the sub-tree of  $\mathbf{t}$  above  $x$  i.e.

$$\mathbf{t}_x = \{y \in \mathbf{t}, x \in \llbracket \partial, y \rrbracket\}$$

rooted at  $x$ . We say that  $x$  is an ancestor of  $y$ , which we denote by  $x \preceq y$ , if  $y \in \mathbf{t}_x$ . We write  $x \prec y$  if furthermore  $x \neq y$ . Notice that  $\preceq$  is a partial order on  $\mathbf{t}$ . We denote by  $x \wedge y$  the Most Recent Common Ancestor (MRCA) of  $x$  and  $y$  in  $\mathbf{t}$  i.e. the unique vertex of  $\mathbf{t}$  such that  $\llbracket \partial, x \rrbracket \cap \llbracket \partial, y \rrbracket = \llbracket \partial, x \wedge y \rrbracket$ .

We denote by  $h_{\mathbf{t}}(x) = d_{\mathbf{t}}(\partial, x)$  the height of the vertex  $x$  in the tree  $\mathbf{t}$  and by  $H(\mathbf{t})$  the height of the tree  $\mathbf{t}$ :

$$H(\mathbf{t}) = \max\{h_{\mathbf{t}}(x), x \in \mathbf{t}\}.$$

For  $\varepsilon > 0$ , we define the erased tree  $r_\varepsilon(\mathbf{t})$  (sometimes called in the literature the  $\varepsilon$ -trimming of the tree  $\mathbf{t}$ ) by

$$r_\varepsilon(\mathbf{t}) = \{x \in \mathbf{t} \setminus \{\partial\}, H(\mathbf{t}_x) \geq \varepsilon\} \cup \{\partial\}.$$

For  $\varepsilon > 0$ ,  $r_\varepsilon(\mathbf{t})$  is indeed a tree and  $r_\varepsilon(\mathbf{t}) = \{\partial\}$  for  $\varepsilon > H(\mathbf{t})$ . Notice that

$$(1) \quad \bigcup_{\varepsilon > 0} r_\varepsilon(\mathbf{t}) = \text{sk}(\mathbf{t}).$$

Recall  $\mathbf{t}$  is a compact rooted real tree and let  $(\mathbf{t}_i, i \in I)$  be a family of trees, and  $(x_i, i \in I)$  a family of vertices of  $\mathbf{t}$ . We denote by  $\mathbf{t}_i^\circ = \mathbf{t}_i \setminus \{\partial_{\mathbf{t}_i}\}$ . We define the tree  $\mathbf{t} \otimes_{i \in I} (\mathbf{t}_i, x_i)$  obtained by grafting the trees  $\mathbf{t}_i$  on the tree  $\mathbf{t}$  at points  $x_i$  by

$$\begin{aligned} \mathbf{t} \otimes_{i \in I} (\mathbf{t}_i, x_i) &= \mathbf{t} \sqcup \left( \bigsqcup_{i \in I} \mathbf{t}_i^\circ \right), \\ d_{\mathbf{t} \otimes_{i \in I} (\mathbf{t}_i, x_i)}(y, y') &= \begin{cases} d_{\mathbf{t}}(y, y') & \text{if } y, y' \in \mathbf{t}, \\ d_{\mathbf{t}_i}(y, y') & \text{if } y, y' \in \mathbf{t}_i^\circ, \\ d_{\mathbf{t}}(y, x_i) + d_{\mathbf{t}_i}(\partial_{\mathbf{t}_i}, y') & \text{if } y \in \mathbf{t} \text{ and } y' \in \mathbf{t}_i^\circ, \\ d_{\mathbf{t}_i}(y, \partial_{\mathbf{t}_i}) + d_{\mathbf{t}}(x_i, x_j) + d_{\mathbf{t}_j}(\partial_{\mathbf{t}_j}, y') & \text{if } y \in \mathbf{t}_i^\circ \text{ and } y' \in \mathbf{t}_j^\circ \text{ with } i \neq j, \end{cases} \\ \partial_{\mathbf{t} \otimes_{i \in I} (\mathbf{t}_i, x_i)} &= \partial_{\mathbf{t}}, \end{aligned}$$

where  $A \sqcup B$  denotes the disjoint union of the sets  $A$  and  $B$ . Notice that  $\mathbf{t} \otimes_{i \in I} (\mathbf{t}_i, x_i)$  might not be compact.

**2.2. The Gromov-Hausdorff topology.** In order to define random real trees, we endow the set of (isometry classes of) rooted compact real trees with a metric, the so-called Gromov-Hausdorff metric, which hence defines a Borel  $\sigma$ -algebra on this set.

First, let us recall the definition of the Hausdorff distance between two compact subsets: let  $A, B$  be two compact subsets of a metric space  $(X, d_X)$ . For every  $\varepsilon > 0$ , we set:

$$A^\varepsilon = \{x \in X, d_X(x, A) \leq \varepsilon\}.$$

Then, the Hausdorff distance between  $A$  and  $B$  is defined by:

$$d_{X, \text{Haus}}(A, B) = \inf\{\varepsilon > 0, B \subset A^\varepsilon \text{ and } A \subset B^\varepsilon\}.$$

Now, let  $(\mathbf{t}, d_{\mathbf{t}}, \partial_{\mathbf{t}})$ ,  $(\mathbf{t}', d_{\mathbf{t}'}, \partial_{\mathbf{t}'})$  be two compact rooted real trees. We define the pointed Gromov-Hausdorff distance between them, see [16, 15], by:

$$d_{GH}(\mathbf{t}, \mathbf{t}') = \inf\{d_{Z, \text{Haus}}(\varphi(\mathbf{t}), \varphi'(\mathbf{t}')) \vee d_Z(\varphi(\partial_{\mathbf{t}}), \varphi'(\partial_{\mathbf{t}'}))\},$$

where the infimum is taken over all metric spaces  $(Z, d_Z)$  and all isometric embeddings  $\varphi : \mathbf{t} \rightarrow Z$  and  $\varphi' : \mathbf{t}' \rightarrow Z$ .

Notice that  $d_{GH}$  is only a pseudo-metric. We say that two rooted real trees  $\mathbf{t}$  and  $\mathbf{t}'$  are equivalent (and we note  $\mathbf{t} \sim \mathbf{t}'$ ) if there exists a root-preserving isometry that maps  $\mathbf{t}$  onto  $\mathbf{t}'$ , that is  $d_{GH}(\mathbf{t}, \mathbf{t}') = 0$ . This clearly defines an equivalence relation. We denote by  $\mathbb{T}$  the set of equivalence classes of compact rooted real trees. The Gromov-Hausdorff distance  $d_{GH}$  hence induces a metric on  $\mathbb{T}$  (that is still denoted by  $d_{GH}$ ). Moreover, the metric space  $(\mathbb{T}, d_{GH})$  is complete and separable, see [15]. If  $\mathbf{t}, \mathbf{t}'$  are two compact rooted real trees such that  $\mathbf{t} \sim \mathbf{t}'$ , then, for every  $\varepsilon > 0$ , we have  $r_\varepsilon(\mathbf{t}) \sim r_\varepsilon(\mathbf{t}')$ . Thus, the erasure function  $r_\varepsilon$  is well-defined on  $\mathbb{T}$ . It is easy to check that the functions  $r_\varepsilon$  for  $\varepsilon > 0$  are 1-Lipschitz.

**2.3. Coding a compact real tree by a function and the Brownian CRT.** Let  $\mathcal{E}$  be the set of continuous function  $g : [0, +\infty) \rightarrow [0, +\infty)$  with compact support and such that  $g(0) = 0$ . For  $g \in \mathcal{E}$ , we set  $\sigma(g) = \sup\{x, g(x) > 0\}$ . Let  $g \in \mathcal{E}$ , and assume that  $\sigma(g) > 0$ , that is  $g$  is not identically zero. For every  $s, t \geq 0$ , we set

$$m_g(s, t) = \inf_{r \in [s \wedge t, s \vee t]} g(r),$$

and

$$(2) \quad d_g(s, t) = g(s) + g(t) - 2m_g(s, t).$$

It is easy to check that  $d_g$  is a pseudo-metric on  $[0, +\infty)$ . We then say that  $s$  and  $t$  are equivalent iff  $d_g(s, t) = 0$  and we set  $T_g$  the associated quotient space. We keep the notation  $d_g$  for the induced distance on  $T_g$ . Then the metric space  $(T_g, d_g)$  is a compact real-tree, see [13]. We denote by  $p_g$  the canonical projection from  $[0, +\infty)$  to  $T_g$ . We will view  $(T_g, d_g)$  as a rooted real tree with root  $\partial = p_g(0)$ . We will call  $(T_g, d_g)$  the real tree coded by  $g$ , and conversely that  $g$  is a contour function of the tree  $T_g$ . We denote by  $F$  the application that associates with a function  $g \in \mathcal{E}$  the equivalence class of the tree  $T_g$ .

Conversely every rooted compact real tree  $(T, d)$  can be coded by a continuous function  $g$  (up to a root-preserving isometry), see [11].

Let  $\theta \in \mathbb{R}$ ,  $\beta > 0$  and  $B^{(\theta)} = (B_t^{(\theta)}, t \geq 0)$  be a Brownian motion with drift  $-2\theta$  and scale  $\sqrt{2/\beta}$ : for  $t \geq 0$ ,

$$B_t^{(\theta)} = \sqrt{2/\beta} B_t - 2\theta t,$$

where  $B$  is a standard Brownian motion. For  $\theta \geq 0$ , let  $n^{(\theta)}[de]$  denote the Itô measure on  $\mathcal{E}$  of positive excursions of  $B^{(\theta)}$  normalized such that for  $\lambda \geq 0$ :

$$(3) \quad n^{(\theta)} \left[ 1 - e^{-\lambda\sigma} \right] = \psi_\theta^{-1}(\lambda),$$

where  $\sigma = \sigma(e)$  denotes the duration (or the length) of the excursion  $e$  and for  $\lambda \geq 0$ :

$$(4) \quad \psi_\theta(\lambda) = \beta\lambda^2 + 2\beta\theta\lambda.$$

Let  $\zeta = \zeta(e) = \max_{s \in [0, \sigma]}(e_s)$  be the maximum of the excursion. We set  $c_\theta(h) = n^{(\theta)}[\zeta \geq h]$  for  $h > 0$ , and we recall, see Section 7 in [9] for the case  $\theta > 0$ , that:

$$(5) \quad c_\theta(h) = \begin{cases} (\beta h)^{-1} & \text{if } \theta = 0 \\ 2\theta (e^{2\beta\theta h} - 1)^{-1} & \text{if } \theta > 0. \end{cases}$$

We define the Brownian CRT,  $\tau = F(e)$ , as the (equivalence class of the) tree coded by the positive excursion  $e$  under  $n^{(\theta)}$ . And we define the measure  $\mathbb{N}^{(\theta)}$  on  $\mathbb{T}$  as the “distribution” of  $\tau$ , that is the push-forward of the measure  $n^{(\theta)}$  by the application  $F$ . Notice that  $H(\tau) = \zeta(e)$ .

*Remark 2.2.* If we translate the former construction into the framework of [12], then, for  $\theta \geq 0$ ,  $B^{(\theta)}$  is the height process which codes the Brownian CRT with branching mechanism  $\psi_\theta$  and it is obtained from the underlying Lévy process  $X = (X_t, t \geq 0)$  with  $X_t = \sqrt{2\beta} B_t - 2\beta\theta t$ .

Let  $e$  with “distribution”  $n^{(\theta)}(de)$  and let  $(\Lambda_s^a, s \geq 0, a \geq 0)$  be the local time of  $e$  at time  $s$  and level  $a$ . Then we define the local time measure of  $\tau$  at level  $a \geq 0$ , denoted by  $\ell_a(dx)$ , as the push-forward of the measure  $d\Lambda_s^a$  by the map  $F$ , see Theorem 4.2 in [13]. We shall define  $\ell_a$  for  $a \in \mathbb{R}$  by setting  $\ell_a = 0$  for  $a \in \mathbb{R} \setminus [0, H(\tau)]$ .

**2.4. Trees with one semi-infinite branch.** The goal of this section is to describe the genealogical tree of a stationary CB with immigration (restricted to the population that appeared before time 0). For this purpose, we add an immortal individual living from  $-\infty$  to 0 that will be the spine of the genealogical tree (i.e. the semi-infinite branch) and will be represented by the half straight line  $(-\infty, 0]$ , see Figure 2. Since we are interested in the genealogical tree, we don't record the population generated by the immortal individual after time 0. The distinguished vertex in the tree will be the point 0 and hence would be the root of the tree in the terminology of Section 2.1. We will however speak of the distinguished leaf in what follows in order to fit with the natural intuition. In the same spirit, we will give another definition for the height of a vertex in such a tree in order to allow negative heights.

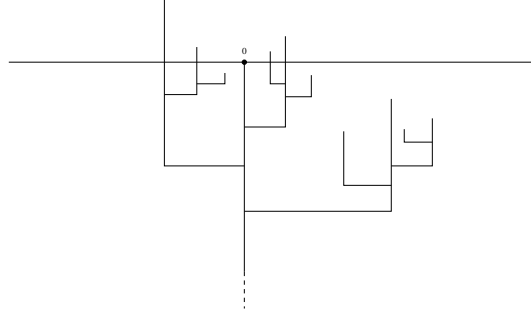


FIGURE 2. An instance of a tree with a semi-infinite branch

**2.4.1. Forests.** A forest  $\mathbf{f}$  is a family  $((h_i, \mathbf{t}_i), i \in I)$  of points of  $\mathbb{R} \times \mathbb{T}$ . Using an immediate extension of the grafting procedure, for an interval  $\mathfrak{J} \subset \mathbb{R}$ , we define the real tree

$$(6) \quad \mathbf{f}_{\mathfrak{J}} = \mathfrak{J} \otimes_{i \in I, h_i \in \mathfrak{J}} (\mathbf{t}_i, h_i).$$

Let us denote, for  $i \in I$ , by  $d_i$  the distance of the tree  $\mathbf{t}_i$  and by  $\mathbf{t}_i^\circ = \mathbf{t}_i \setminus \{\partial_{\mathbf{t}_i}\}$  the tree  $\mathbf{t}_i$  without its root. The distance on  $\mathbf{f}_{\mathfrak{J}}$  is then defined, for  $x, y \in \mathbf{f}_{\mathfrak{J}}$ , by:

$$d_{\mathbf{f}}(x, y) = \begin{cases} d_i(x, y) & \text{if } x, y \in \mathbf{t}_i^\circ, \\ h_{\mathbf{t}_i}(x) + |h_i - h_j| + h_{\mathbf{t}_j}(y) & \text{if } x \in \mathbf{t}_i^\circ, y \in \mathbf{t}_j^\circ \text{ with } i \neq j, \\ |x - h_j| + h_{\mathbf{t}_j}(y) & \text{if } x \notin \bigsqcup_{i \in I} \mathbf{t}_i^\circ, y \in \mathbf{t}_j^\circ \\ |x - y| & \text{if } x, y \notin \bigsqcup_{i \in I} \mathbf{t}_i^\circ. \end{cases}$$

Let us recall the following lemma (see [2]).

**Lemma 2.3.** *Let  $\mathfrak{J} \subset \mathbb{R}$  be a closed interval. If for every  $a, b \in \mathfrak{J}$ , such that  $a < b$ , and every  $\varepsilon > 0$ , the set  $\{i \in I, h_i \in [a, b], H(\mathbf{t}_i) > \varepsilon\}$  is finite, then the tree  $\mathbf{f}_{\mathfrak{J}}$  is a complete locally compact length space.*

**2.4.2. Trees with one semi-infinite branch.**

**Definition 2.4.** *We set  $\mathbb{T}_1$  the set of forests  $\mathbf{f} = ((h_i, \mathbf{t}_i), i \in I)$  such that*

- *for every  $i \in I$ ,  $h_i \leq 0$ ,*
- *for every  $a < b$ , and every  $\varepsilon > 0$ , the set  $\{i \in I, h_i \in [a, b], H(\mathbf{t}_i) > \varepsilon\}$  is finite.*

The following corollary, which is an elementary consequence of Lemma 2.3, associates with a forest  $\mathbf{f} \in \mathbb{T}_1$  a complete and locally compact real tree.

**Corollary 2.5.** *Let  $\mathbf{f} = ((h_i, \mathbf{t}_i), i \in I) \in \mathbb{T}_1$ . Then, the tree  $\mathbf{f}_{(-\infty, 0]}$  defined by (6) is a complete and locally compact real tree.*

Conversely, let  $(\mathbf{t}, d_{\mathbf{t}}, \rho_0)$  be a complete and locally compact rooted real tree. We denote by  $\mathcal{S}(\mathbf{t})$  the set of vertices  $x \in \mathbf{t}$  such that at least one of the connected components of  $\mathbf{t} \setminus \{x\}$  that do not contain  $\rho_0$  is unbounded. If  $\mathcal{S}(\mathbf{t})$  is not empty, then it is a tree which contains  $\rho_0$ . We say that  $\mathbf{t}$  has a unique semi-infinite branch if  $\mathcal{S}(\mathbf{t})$  is non-empty and has no branching point. We set  $(\mathbf{t}_i^\circ, i \in I)$  the connected components of  $\mathbf{t} \setminus \mathcal{S}(\mathbf{t})$ . For every  $i \in I$ , we set  $x_i$  the unique point of  $\mathcal{S}(\mathbf{t})$  such that  $\inf\{d_{\mathbf{t}}(x_i, y), y \in \mathbf{t}_i^\circ\} = 0$ , and:

$$\mathbf{t}_i = \mathbf{t}_i^\circ \cup \{x_i\}, \quad h_i = -d(\rho_0, x_i).$$

We shall say that  $x_i$  is the root of  $\mathbf{t}_i$ . Notice first that  $(\mathbf{t}_i, d_{\mathbf{t}}, x_i)$  is a bounded rooted tree. It is also compact since, according to the Hopf-Rinow theorem (see Theorem 2.5.26 in [8]), it is a bounded closed subset of a complete locally compact length space. Thus it belongs to  $\mathbb{T}$ .

The family  $\mathbf{f} = ((h_i, \mathbf{t}_i), i \in I)$  is therefore a forest with  $h_i < 0$ . To check that it belongs to  $\mathbb{T}_1$ , we need to prove that the second condition in Definition 2.4 is satisfied which is a direct consequence of the fact that the tree  $\mathbf{f}_{[a,b]}$  is locally compact.

We can therefore identify the set  $\mathbb{T}_1$  with the set of (equivalence classes) of complete locally compact rooted real trees with a unique semi-infinite branch. We can follow [3] to endow  $\mathbb{T}_1$  with a Gromov-Hausdorff-type distance for which  $\mathbb{T}_1$  is a Polish space.

We extend the partial order defined for trees in  $\mathbb{T}$  to forests in  $\mathbb{T}_1$ , with the idea that the distinguished leaf  $\rho_0 = 0$  is at the tip of the semi-infinite branch. Let  $\mathbf{f} = (h_i, \mathbf{t}_i)_{i \in I} \in \mathbb{T}_1$  and write  $\mathbf{t} = \mathbf{f}_{(-\infty, 0]}$  viewed as a real tree rooted at  $\rho_0 = 0$  (with a unique semi-infinite branch  $\mathcal{S}(\mathbf{t}) = (-\infty, 0]$ ). For  $x, y \in \mathbf{t}$ , we set  $x \preceq y$  if either  $x, y \in \mathcal{S}(\mathbf{t})$  and  $d_{\mathbf{f}}(x, \rho_0) \geq d_{\mathbf{f}}(y, \rho_0)$ , or  $x, y \in \mathbf{t}_i$  for some  $i \in I$  and  $x \preceq y$  (with the partial order for the rooted compact real tree  $\mathbf{t}_i$ ), or  $x \in \mathcal{S}(\mathbf{t})$  and  $y \in \mathbf{t}_i$  for some  $i \in I$  and  $d_{\mathbf{f}}(x, \rho_0) \geq |h_i|$ . We write  $x \prec y$  if furthermore  $x \neq y$ . We define  $x \wedge y$  the MRCA of  $x, y \in \mathbf{t}$  as  $x$  if  $x \preceq y$ , as  $x \wedge y$  if  $x, y \in \mathbf{t}_i$  for some  $i \in I$  (with the MRCA for the rooted compact real tree  $\mathbf{t}_i$ ), as  $h_i \wedge h_j$  if  $x \in \mathbf{t}_i$  and  $y \in \mathbf{t}_j$  for some  $i \neq j$ . We define the height of a vertex  $x \in \mathbf{t}$  as

$$h_{\mathbf{f}}(x) = d_{\mathbf{f}}(x, \rho_0 \wedge x) - d_{\mathbf{f}}(\rho_0, \rho_0 \wedge x).$$

Notice that the definition of the height function  $h_{\mathbf{f}}$  for a forest  $\mathbf{f} = (h_i, \mathbf{t}_i)_{i \in I} \in \mathbb{T}_1$  is different than the height function of the tree  $\mathbf{t} = \mathbf{f}_{(-\infty, 0]}$  viewed as a tree in  $\mathbb{T}$ , as in the former case the root  $\rho_0$  is viewed as a distinguished vertex above the semi-infinite branch (all elements of this semi-infinite branch have negative heights for  $h_{\mathbf{f}}$  whereas all the heights are nonnegative for  $h_{\mathbf{t}}$ ).

**2.4.3. Coding a forest by a contour function.** We want to extend the construction of a tree of the type  $\mathbf{f}_{(-\infty, 0]}$  via a contour function as in Section 2.3. Let  $\mathcal{E}_1$  be the set of continuous functions  $g$  defined on  $\mathbb{R}$  such that  $g(0) = 0$  and  $\liminf_{x \rightarrow -\infty} g(x) = \liminf_{x \rightarrow +\infty} g(x) = -\infty$ . For such a function, we still consider the pseudo-metric  $d_g$  defined by (2) (but for  $s, t \in \mathbb{R}$ ) and define the tree  $T_g$  as the quotient space on  $\mathbb{R}$  induced by this pseudo-metric. We set  $p_g$  as the canonical projection from  $\mathbb{R}$  onto  $T_g$ .

**Lemma 2.6.** *Let  $g \in \mathcal{E}_1$ . The triplet  $(T_g, d_g, p_g(0))$  is a complete locally compact rooted real tree with a unique semi-infinite branch.*

*Proof.* We define the infimum function  $\underline{g}(x)$  on  $\mathbb{R}$  as the infimum of  $g$  between 0 and  $x$ :  $\underline{g}(x) = \inf_{[x \wedge 0, x \vee 0]} g$ . The function  $g - \underline{g}$  is non-negative and continuous. Let  $((a_i, b_i), i \in I)$  be the excursion intervals of  $g - \underline{g}$  above 0. Because of the hypothesis on  $g$ , the intervals  $(a_i, b_i)$  are bounded. For  $i \in I$ , set  $h_i = g(a_i)$  and  $g_i(x) = g((a_i + x) \wedge b_i) - h_i$  so that  $g_i \in \mathcal{E}$ . Consider the forest  $\mathbf{f} = ((h_i, T_{g_i}), i \in I)$ .



It is elementary to check that  $(\mathbf{f}_{(-\infty, g(0)]}, d_{\mathbf{f}}, g(0))$  and  $(T_g, d_g, p_g(0))$  are root-preserving isometric. To conclude, it is enough to check that  $\mathbf{f} \in \mathbb{T}_1$ . First remark that, by definition,  $h_i \leq 0$  for every  $i \in I$ . Let  $r > 0$  and set  $r_g = \inf\{x, \underline{g}(x) \geq g(0) - r\}$  and  $r_d = \sup\{x, \underline{g}(x) \geq g(0) - r\}$ . Because of the hypothesis on  $g$ , we have that  $r_g$  and  $r_d$  are finite. By continuity of  $g - \underline{g}$  on  $[r_g, r_d]$ , we deduce that for any  $\varepsilon > 0$ , the set  $\{i \in I; (a_i, b_i) \subset [r_g, r_d] \text{ and } \sup_{(a_i, b_i)}(g - \underline{g}) > \varepsilon\}$  is finite. Since this holds for any  $r > 0$  and that  $H(T_{g_i}) = \sup_{(a_i, b_i)}(g - \underline{g})$  for all  $i \in I$ , we deduce that  $\mathbf{f} \in \mathbb{T}_1$ . This concludes the proof.  $\square$

**2.4.4. Genealogical tree of an extant population.** For a tree  $\mathbf{t} \in \mathbb{T}$  or  $\mathbf{t} \in \mathbb{T}_1$  (recall that we identify a forest  $\mathbf{f} \in \mathbb{T}_1$  with the tree  $\mathbf{t} = \mathbf{f}_{(-\infty, 0]}$  with a different definition for the height function) and  $h \geq 0$ , we define  $\mathcal{Z}_h(\mathbf{t}) = \{x \in \mathbf{t}, h_{\mathbf{t}}(x) = h\}$  the set of vertices of  $\mathbf{t}$  at level  $h$  also called the extant population at time  $h$ , and the genealogical tree of the vertices of  $\mathbf{t}$  at level  $h$  by:

$$(7) \quad \mathcal{G}_h(\mathbf{t}) = \{x \in \mathbf{t}; \exists y \in \mathcal{Z}_h(\mathbf{t}) \text{ such that } x \preceq y\}.$$

For  $\mathbf{f} \in \mathbb{T}_1$ , we write  $\mathcal{G}_h(\mathbf{f})$  for  $\mathcal{G}_h(\mathbf{f}_{(-\infty, 0]})$ ;

### 3. ANCESTRAL PROCESS

Usually, the ancestral process records the genealogy of  $n$  extant individuals at time 0 picked at random among the whole population. Using the ideas of [5], we are able to describe in the case of a Brownian forest the genealogy of all extant individuals at time 0 by a simple Poisson point process on  $\mathbb{R}^2$ .

#### 3.1. Construction of a tree from a point measure.

**Definition 3.1.** A point process  $\mathcal{A}(dx, d\zeta) = \sum_{i \in \mathcal{I}} \delta_{(x_i, \zeta_i)}(dx, d\zeta)$  on  $\mathbb{R}^* \times (0, +\infty)$  is said to be an ancestral process if

- (i)  $\forall i, j \in \mathcal{I}, i \neq j \implies x_i \neq x_j$ .
- (ii)  $\forall a, b \in \mathbb{R}, \forall \varepsilon > 0, \mathcal{A}([a, b] \times [\varepsilon, +\infty)) < +\infty$ .
- (iii)  $\sup\{\zeta_i, x_i > 0\} = +\infty$  if  $\sup_{i \in \mathcal{I}} x_i = +\infty$ ; and  $\sup\{\zeta_i, x_i < 0\} = +\infty$  if  $\inf_{i \in \mathcal{I}} x_i = -\infty$ .

Let  $\mathcal{A} = \sum_{i \in \mathcal{I}} \delta_{(x_i, \zeta_i)}$  be a point process on  $\mathbb{R}^* \times [0, +\infty)$  satisfying (i) and (ii) from Definition 3.1. We shall associate with this ancestral process a genealogical tree. Informally the genealogical tree is constructed as follows. We view this process as a sequence of vertical segments in  $\mathbb{R}^2$ , the tips of the segments being the  $x_i$ 's and their lengths being the  $\zeta_i$ 's. We then attach the bottom of each segment such that  $x_i > 0$  (resp.  $x_i < 0$ ) to the first left (resp. first right) longer segment or to the half line  $\{0\} \times (-\infty, 0]$  if such a segment does not exist. This gives a (unrooted, non-compact) real tree that may not be complete. See also Figure 1 for an example.

Let us turn to a more formal definition. Let us denote by  $\mathcal{I}^d = \{i \in \mathcal{I}, x_i > 0\}$  and  $\mathcal{I}^g = \{i \in \mathcal{I}, x_i < 0\} = \mathcal{I} \setminus \mathcal{I}^d$ . We also set  $\mathcal{I}_0 = \mathcal{I} \sqcup \{0\}$ ,  $x_0 = 0$  and  $\zeta_0 = +\infty$ . We set, for every  $i \in \mathcal{I}_0$ ,  $S_i = \{x_i\} \times (-\zeta_i, 0]$  the vertical segment in  $\mathbb{R}^2$  that links the points  $(x_i, 0)$  and  $(x_i, -\zeta_i)$ . Notice that we omit the lowest point of the vertical segments. Finally we define

$$(8) \quad \mathfrak{T} = \bigsqcup_{i \in \mathcal{I}_0} S_i.$$

We now define a distance on  $\mathfrak{T}$ . We first define the distance between leaves of  $\mathfrak{T}$ , i.e. points  $(x_i, 0)$  with  $i \in \mathcal{I}_0$ , then we extend it to every point of  $\mathfrak{T}$ . For  $i, j \in \mathcal{I}_0$  such that  $x_i < x_j$ , we set

$$(9) \quad d((x_i, 0), (x_j, 0)) = 2 \max\{\zeta_k, x_k \in J(x_i, x_j)\},$$

where, for  $x < y$ ,  $J(x, y) = (x, y]$  (resp.  $[x, y)$ , resp.  $[x, y] \setminus \{0\}$ ) if  $x \geq 0$  (resp.  $y \leq 0$ , resp.  $x < 0$  and  $y > 0$ ), with the convention  $\max \emptyset = 0$ . For  $u = (x_i, a) \in S_i$  and  $v = (x_j, b) \in S_j$ , we set, with  $r = \frac{1}{2}d((x_i, 0), (x_j, 0))$ :

$$(10) \quad d(u, v) = |a - b| \mathbf{1}_{\{x_i = x_j\}} + (|a - r| + |b - r|) \mathbf{1}_{\{x_i \neq x_j\}}.$$

See Figure 3 for an example. It is easy to verify that  $d$  is a distance on  $\mathfrak{T}$ . Notice that  $\mathfrak{T}$  is not compact in particular because of the infinite half-line attached to  $(0, 0)$ . In order to stick to the framework of Section 2.4, the origin  $(0, 0)$  will be the distinguished point in  $\mathfrak{T}$  located at height  $h = 0$ .

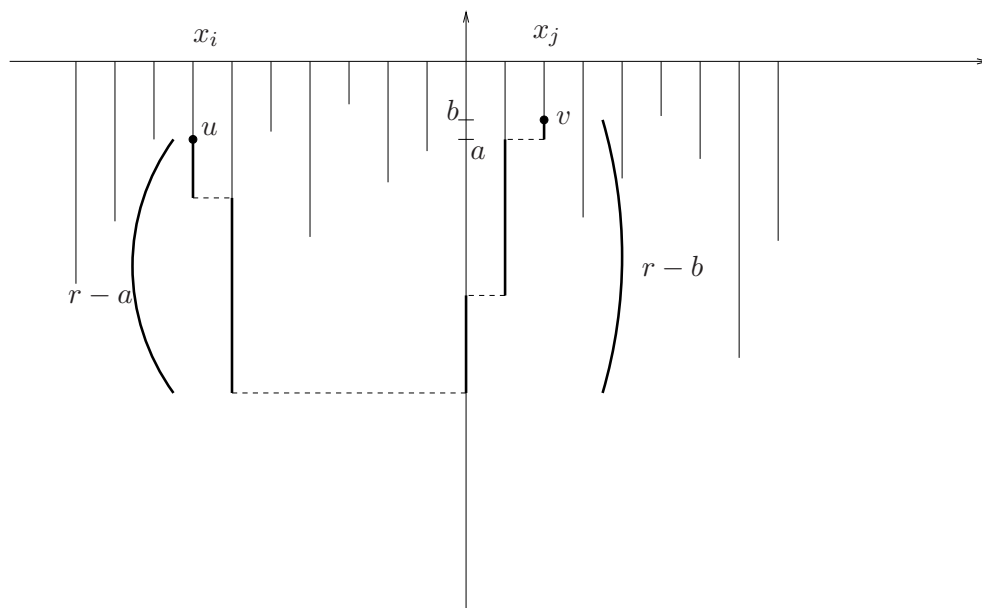


FIGURE 3. An example of the distance  $d(u, v)$  defined in (9)

Finally, we define  $\mathfrak{T}(\mathcal{A})$ , with the metric  $d$ , as the completion of the metric space  $(\mathfrak{T}, d)$ .

*Remark 3.2.* For every  $i \in \mathcal{I}^d$ , we set  $i_\ell$  the index in  $\mathcal{I}_0$  such that

$$x_{i_\ell} = \max\{x_j, 0 \leq x_j < x_i \text{ and } \zeta_j > \zeta_i\}.$$

Remark that  $i_\ell$  is well defined since there are only a finite number of indices  $j \in \mathcal{I}_0$  such that  $x_j \in [0, x_i)$  and  $\zeta_j > \zeta_i$ . Similarly, for  $i \in \mathcal{I}^g$ , we set  $i_r$  the index in  $\mathcal{I}_0$  such that

$$x_{i_r} = \min\{x_j, x_i < x_j \leq 0 \text{ and } \zeta_j > \zeta_i\}.$$

The distance  $d$  identifies the point  $(x_i, -\zeta_i)$  (which does not belong to  $\mathfrak{T}$  by definition) with the point  $(x_{i_\ell}, -\zeta_i)$  if  $x_i > 0$  and with the point  $(x_{i_r}, -\zeta_i)$  if  $x_i < 0$  as illustrated on the right-hand side of Figure 4.

**Proposition 3.3.** *Let  $\mathcal{A}$  be an ancestral process. The tree  $(\mathfrak{T}(\mathcal{A}), d, (0, 0))$  is a complete and locally compact rooted real tree with a unique semi-infinite branch and the associated forest belongs to  $\mathbb{T}_1$ .*

We shall call  $\mathfrak{T}(\mathcal{A})$  the tree associated with the ancestral process  $\mathcal{A}$ .

*Proof.* As the completion of a real tree is still a real tree, it is enough to prove that  $(\mathfrak{T}, d)$  is a real tree, with  $\mathfrak{T}$  defined by (8) and  $d$  defined by (9) and (10).

*First case:  $\mathcal{I}$  finite.*

We can suppose that  $\mathcal{I} = \{1, \dots, n\}$  with  $x_1 < x_2 < \dots < x_n$  (with  $x_i \neq 0$  for  $i \in \mathcal{I}$ ). We consider the continuous, piece-wise affine function  $g$  on  $\mathbb{R}$  such that

- For  $1 \leq i \leq n$ ,  $g(x_i) = -\zeta_i$ ,
- For  $1 \leq i \leq n-1$ ,  $g\left(\frac{x_i+x_{i+1}}{2}\right) = 0$ ,
- $g(x_1-1) = g(x_n+1) = 0$ ,
- $g'(x) = -1$  for  $x < x_1-1$  and  $g'(x) = 1$  for  $x > x_n+1$ .

Then, it is easy to see that  $(\mathfrak{T}, d)$  is the tree  $T_g$  coded by  $g$  (see Section 2.4) and hence is a real tree. Notice that the number of leaves of  $\mathfrak{T}$  is  $\text{Card}(\mathcal{I}_0) = n+1$ .

*Second case:  $\mathcal{I}$  infinite,  $\sup_{i \in \mathcal{I}} x_i < +\infty$  and  $\inf_{i \in \mathcal{I}} x_i > -\infty$ .*

In that case, by Condition (ii) in Definition 3.1, we can order the set  $\mathcal{I}$  via a sequence  $(i_1, i_2, \dots)$  such that the sequence  $(\zeta_{i_k}, k \geq 1)$  is non-increasing. For every  $n \geq 1$ , we denote by  $(\mathfrak{T}_n, d_n)$  the tree associated with the ancestral process  $\sum_{k=1}^n \delta_{(x_{i_k}, \zeta_{i_k})}$  (which is indeed a tree according to the first case). Remark first that  $\mathfrak{T}_n \subset \mathfrak{T}_{n+1}$ . Moreover, as  $\zeta_{i_{n+1}} \leq \zeta_{i_k}$  for every  $1 \leq k \leq n$ , we deduce from (9) that  $d_n$  is equal to the restriction of  $d_{n+1}$  to  $\mathfrak{T}_n$ . Therefore, we have  $\mathfrak{T} = \bigcup_{n \geq 1} \mathfrak{T}_n$  and  $d$  is the distance induced by the distances  $d_n$ . We deduce that  $(\mathfrak{T}, d)$  is a real tree as limit of increasing real trees. Indeed, clearly  $\mathfrak{T}$  is connected (as the union of an increasing sequence of connected sets) and  $d$  satisfies the so-called "4-points condition" (see Lemma 3.12 in [14]). To conclude, use that those two conditions characterize real trees (see Theorem 3.40 in [14]). We deduce that  $(\mathfrak{T}, d)$  is a real tree.

*Third case:  $\mathcal{I}$  infinite and  $\sup_{i \in \mathcal{I}} x_i = +\infty$  or  $\inf_{i \in \mathcal{I}} x_i = -\infty$ .*

We consider in that case, for every integer  $m \geq 1$  the tree  $(\mathfrak{T}_m, d_m)$  induced by the ancestral process  $\mathcal{A}$  restricted to  $[-m, m] \times [0, +\infty)$  (which is indeed a tree by the second case). We still have  $\mathfrak{T} = \bigcup_{m \geq 1} \mathfrak{T}_m$  and the compatibility condition for the distances. We then conclude as for the second case that  $(\mathfrak{T}, d)$  is a real tree.

By construction of  $\mathfrak{T}$ , it is easy to check that  $\mathfrak{T}(\mathcal{A})$  has a unique semi-infinite branch.

Let us now prove that  $\mathfrak{T}(\mathcal{A})$  is locally compact. Let  $(y_n, n \in \mathbb{N})$  be a bounded sequence of  $\mathfrak{T}$ .

On one hand, let us assume that there exists  $i \in \mathcal{I}_0$  and a sub-sequence  $(y_{n_k}, k \in \mathbb{N})$  such that  $y_{n_k}$  belongs to  $S_i = \{x_i\} \times (-\zeta_i, 0]$ . Since, for  $i \in \mathcal{I}$ , there exists a unique  $j \in \mathcal{I}_0$  such that  $S_i \cup \{(x_j, -\zeta_j)\}$  is compact in  $(\mathfrak{T}, d)$ , see Remark 3.2, and for  $i = 0$ ,  $S_0 = \{0\} \times (-\infty, 0]$ , we deduce that the bounded sequence  $(y_{n_k}, k \in \mathbb{N})$  has an accumulation point in  $S_i \cup \{(x_j, -\zeta_j)\}$  if  $i \in \mathcal{I}$  or in  $\{0\} \times (-\infty, 0]$  if  $i = 0$ .

On the other hand, let us assume that for all  $i \in \mathcal{I}_0$  the sets  $\{n, y_n \in S_i\}$  are finite. For  $n \in \mathbb{N}$ , let  $i_n$  uniquely defined by  $y_n \in S_{i_n}$ . Since  $(y_n, n \in \mathbb{N})$  is bounded, we deduce from Conditions (ii-iii) in Definition 3.1, that the sequence  $(x_{i_n}, n \in \mathbb{N})$  is bounded in  $\mathbb{R}$ . In particular, there is a sub-sequence such that  $(x_{i_{n_k}}, k \in \mathbb{N})$  converges to a limit say  $a$ . Without loss of generality, we can assume that the sub-sequence is non-decreasing. We deduce from Condition (ii) in Definition (3.1) that  $\lim_{\varepsilon \downarrow 0} \max\{\zeta_i, a - \varepsilon < x_i < a\} = 0$ . This implies thanks to Definition (9) that

$(\{x_{i_{n_k}}\} \times \{0\}, k \in \mathbb{N})$  is Cauchy in  $\mathfrak{T}$  and using (ii) again that  $\lim_{k \rightarrow +\infty} \zeta_{i_{n_k}} = 0$ . Then use that

$$d(y_{n_k}, y_{n_{k'}}) \leq \zeta_{i_{n_k}} + \zeta_{i_{n_{k'}}} + d((x_{n_k}, 0), (x_{n_{k'}}, 0))$$

to conclude that the  $(y_{n_k}, k \in \mathbb{N})$  is Cauchy in  $\mathfrak{T}$ .

We deduce that all bounded sequence in  $\mathfrak{T}$  has a Cauchy sub-sequence. This proves that  $\mathfrak{T}(\mathcal{A})$ , the completion of  $\mathfrak{T}$  is locally compact.  $\square$

*Remark 3.4.* In the proof of Proposition 3.3, Conditions (i) and (ii) in Definition 3.1 insure that  $\mathfrak{T}(\mathcal{A})$  is a tree and Conditions (ii) and (iii) that  $\mathfrak{T}(\mathcal{A})$  is locally compact.

**3.2. The ancestral process of the Brownian forest.** Let  $\theta \geq 0$ . Let  $\mathcal{N}(dh, d\varepsilon, de) = \sum_{i \in I} \delta_{(h_i, \varepsilon_i, e_i)}(dh, d\varepsilon, de)$  be, under  $\mathbb{P}^{(\theta)}$ , a Poisson point measure on  $\mathbb{R} \times \{-1, 1\} \times \mathcal{E}$  with intensity  $\beta dh (\delta_{-1}(d\varepsilon) + \delta_1(d\varepsilon)) n^{(\theta)}(de)$ , and let  $\mathcal{F} = ((h_i, \tau_i), i \in I)$  be the associated Brownian forest where  $\tau_i = T_{e_i}$  is the tree associated with the excursion  $e_i$ , see Section 2.3. As explained in Section 3.2.4, this Brownian forest models the evolution of a stationary population directed by the branching mechanism  $\psi_\theta$  defined in (4).

We want to describe the genealogical tree of the extant population at some fixed time, say 0. The looked after genealogical tree is then  $\mathcal{G}_0(\mathcal{F})$  defined by (7). To describe the distribution of this tree, we use an ancestral process as described in the previous subsection. We first construct a contour process  $(B_t, t \in \mathbb{R})$  (obtained by the concatenation of two Brownian motions with drift) which codes for the tree  $\mathcal{F}_{(-\infty, 0]}$  (see Section 2.4 for the notations). The supplementary variables  $\varepsilon_i$  are needed at this point to decide if the tree  $\mathbf{t}_i$  is located on the left or on the right of the infinite spine. The atoms of the ancestral process are then the pairs formed by the points of growth of the local time at 0 of  $B$  and the depth of the associated excursion of  $B$  below 0.

**3.2.1. Construction of the contour process.** Set  $\mathcal{I} = \{i \in I; h_i < 0\}$ .

For every  $i \in \mathcal{I}$ , we set:

$$a_i = \sum_{j \in \mathcal{I}} \mathbf{1}_{\{\varepsilon_j = \varepsilon_i\}} \mathbf{1}_{\{h_j < h_i\}} \sigma(e_j) \quad \text{and} \quad b_i = a_i + \sigma(e_i),$$

where we recall that  $\sigma(e_i)$  is the length of excursion  $e_i$ . For every  $t \geq 0$ , we set  $i_t^d$  (resp.  $i_t^g$ ) the only index  $i \in \mathcal{I}$  such that  $\varepsilon_i = 1$  (resp.  $\varepsilon_i = -1$ ) and  $a_i \leq t < b_i$ . Notice that  $i_t^d$  and  $i_t^g$  are a.s. well defined but on a Lebesgue-null set of values of  $t$ . We set  $B^d = (B_t^d, t \geq 0)$  and  $B^g = (B_t^g, t \geq 0)$  where for  $t \geq 0$ :

$$B_t^d = h_{i_t^d} + e_{i_t^d}(t - a_{i_t^d}) \quad \text{and} \quad B_t^g = h_{i_t^g} + e_{i_t^g}(\sigma(e_{i_t^g}) - (t - a_{i_t^g})).$$

By standard excursion theory (decomposition of  $B^{(\theta)}$  above its minimum), we have the following result.

**Proposition 3.5.** *Let  $\theta \geq 0$ . The processes  $B^d$  and  $B^g$  are two independent Brownian motions distributed as  $B^{(\theta)}$ .*

We define the process  $B = (B_t, t \in \mathbb{R})$  by  $B_t = B_t^d \mathbf{1}_{\{t > 0\}} + B_{-t}^g \mathbf{1}_{\{t < 0\}}$ . By construction, the process  $B$  indeed codes for the tree  $\mathcal{F}_{(-\infty, 0]}$ .

**3.2.2. The ancestral process.** Let  $(L_t^\ell, t \geq 0)$  be the local time at 0 of the process  $B^\ell$ , where  $\ell \in \{g, d\}$ . We denote by  $((\alpha_i, \beta_i), i \in \mathcal{I}^\ell)$  the excursion intervals of  $B^\ell$  below 0, omitting the last infinite excursion if any, and, for every  $i \in \mathcal{I}^\ell$ , we set  $\zeta_i = -\min\{B_s^\ell, s \in (\alpha_i, \beta_i)\}$ .

We consider the point measure on  $\mathbb{R} \times \mathbb{R}_+$  defined by:

$$\mathcal{A}^{\mathcal{N}}(du, d\zeta) = \sum_{i \in \mathcal{I}^d} \delta_{(L_{\alpha_i}^d, \zeta_i)}(du, d\zeta) + \sum_{i \in \mathcal{I}^g} \delta_{(-L_{\alpha_i}^g, \zeta_i)}(du, d\zeta).$$

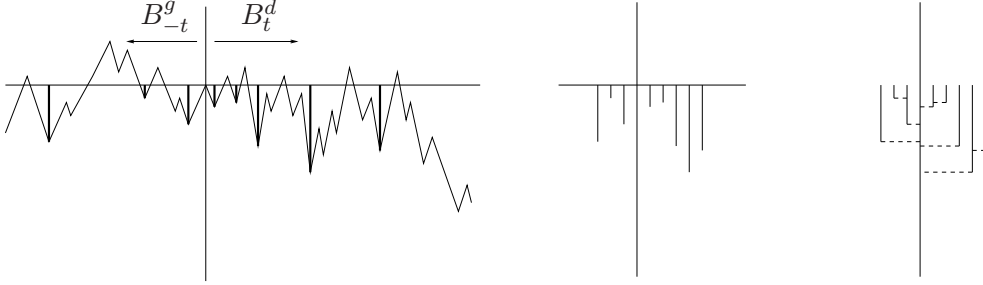


FIGURE 4. The Brownian motions with drift, the ancestral process and the associated genealogical tree

See Figure 4 for a representation of the contour process  $B$ , the ancestral process  $\mathcal{A}^{\mathcal{N}}$  and the genealogical tree  $\mathcal{G}_0(\mathcal{F})$ . In this figure, the horizontal axis represents the time for Brownian motion on the left-hand figure whereas it is in the scale of local time for the ancestral process on the two right-hand figures. This will always be the case in the rest of the paper dealing with ancestral processes.

Let  $[-E_g, E_d]$  be the closed support of the measure  $\mathcal{A}^{\mathcal{N}}(du, \mathbb{R}_+)$ :

$$E_d = \inf\{u \geq 0, \mathcal{A}([u, +\infty) \times \mathbb{R}_+) = 0\} \quad \text{and} \quad E_g = \inf\{u \geq 0, \mathcal{A}((-\infty, -u] \times \mathbb{R}_+) = 0\},$$

with the convention that  $\inf \emptyset = +\infty$ . Notice that, for  $\ell \in \{g, d\}$ , we also have  $E_\ell = L_\infty^\ell$ . We now give the distribution of the ancestral process  $\mathcal{A}^{\mathcal{N}}$ . Recall  $c_\theta$  defined by (5).

**Proposition 3.6.** *Let  $\theta \geq 0$ . Under  $\mathbb{P}^{(\theta)}$ , the random variables  $E_g, E_d$  are independent and exponentially distributed with parameter  $2\theta$  (and mean  $1/2\theta$ ) with the convention that  $E_d = E_g = +\infty$  if  $\theta = 0$ . Under  $\mathbb{P}^{(\theta)}$  and conditionally given  $(E_g, E_d)$ , the ancestral process  $\mathcal{A}^{\mathcal{N}}(du, d\zeta)$  is a Poisson point measure with intensity:*

$$\mathbf{1}_{(-E_g, E_d)}(u) du |c'_\theta(\zeta)| d\zeta.$$

Notice that the random measure  $\mathcal{A}^{\mathcal{N}}$  satisfies Conditions (i)-(iii) from Definition 3.1 and is thus indeed an ancestral process.

This result is very similar to Corollary 2 in [6]. The main additional ingredient here is the order (given by the  $u$  variable) which will be very useful in the simulation.

*Proof.* Since  $B^d$  and  $B^g$  are independent with the same distribution, we deduce that  $E_g$  and  $E_d$  are independent with the same distribution. Let  $\theta > 0$ . Since  $B^d$  is a Brownian motion with drift  $-2\theta$ , we deduce from [7], page 90, that  $E_d$  is exponential with mean  $1/2\theta$ . The case  $\theta = 0$  is immediate.

The excursions below zero of  $B^d$  conditionally given  $E_d$  are excursions of a Brownian motion  $B^{(-\theta)}$  with drift  $2\theta$  (after symmetry with respect to 0) conditioned on being finite, that is excursions of a Brownian motion  $B^{(\theta)}$  with drift  $-2\theta$ . Moreover, by (5),  $c_\theta$  is exactly the tail distribution of the maximum of an excursion under  $n^{(\theta)}$ . Standard theory of Brownian excursions gives then the result.  $\square$

**3.2.3. Identification of the trees.** Let  $\mathfrak{T}^{\mathcal{N}} = \mathfrak{T}(\mathcal{A}^{\mathcal{N}})$  denote the locally compact tree associated with the ancestral process  $\mathcal{A}^{\mathcal{N}}$ , see Proposition 3.3. According to the following proposition, we shall say that the ancestral process  $\mathcal{A}^{\mathcal{N}}$  codes for the genealogical tree of the extant population at time 0 for the forest  $\mathcal{F}$ .

**Proposition 3.7.** *Let  $\theta \geq 0$ . The trees  $\mathcal{G}_0(\mathcal{F})$  under  $\mathbb{P}^{(\theta)}$  and  $\mathfrak{T}^{\mathcal{N}}$  belong to the same equivalence class in  $\mathbb{T}_1$ .*

*Proof.* Let us first remark that the genealogical tree  $\mathcal{G}_0(\mathcal{F})$  can be directly constructed using the process  $B$  as described on Figure 5.

More precisely, recall that  $B$  is the contour function of the tree  $\mathcal{F}_{(-\infty,0]}$ . Let us denote by  $p_B$  the canonical projection from  $\mathbb{R}$  to  $\mathcal{F}_{(-\infty,0]}$  as defined in Section 2.4. Recall  $((\alpha_i, \beta_i), i \in \mathcal{I}^\ell)$ , with  $\ell \in \{g, d\}$ , are the excursion intervals of  $B^\ell$  below 0. Then  $\mathcal{G}_0(\mathcal{F})$  is the smallest complete sub-tree of  $\mathcal{F}_{(-\infty,0]}$  that contains the points  $(p_B(\alpha_i), i \in \mathcal{I}_g \cup \mathcal{I}_d)$  and the semi-infinite branch of  $\mathcal{F}_{(-\infty,0]}$ .

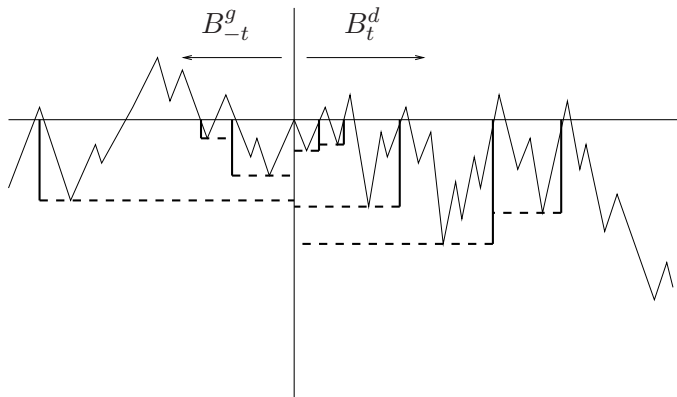


FIGURE 5. The genealogical tree inside the Brownian motions

Let  $i, j \in \mathcal{I}$  with  $0 < \alpha_i < \alpha_j$  for instance. By definition of the tree coded by a function, the distance between  $p_B(\alpha_i)$  and  $p_B(\alpha_j)$  in  $\mathcal{G}_0(\mathcal{F})$  is given by:

$$d(p_B(\alpha_i), p_B(\alpha_j)) = -2 \min_{u \in [\alpha_i, \alpha_j]} B_u.$$

But, by definition of  $\mathcal{A}^{\mathcal{N}}$ , we have:

$$\begin{aligned} - \min_{u \in [\alpha_i, \alpha_j]} B_u &= \max_{k \in \mathcal{I} \mid \alpha_i \leq \alpha_k < \alpha_j} \left( - \min_{u \in [\alpha_k, \beta_k]} B_u \right) \\ &= \max_{k \in \mathcal{I} \mid \alpha_i \leq \alpha_k < \alpha_j} \zeta_k. \end{aligned}$$

The other cases  $\alpha_j < \alpha_i < 0$  and  $\alpha_i < 0 < \alpha_j$  can be handled similarly. We deduce that the distances on a dense subset of leaves of  $\mathcal{G}_0(\mathcal{F})$  and  $\mathfrak{T}^{\mathcal{N}}$  coincide, which implies the result by completeness of the trees.  $\square$

**3.2.4. Local times.** The Brownian forest  $\mathcal{F}$  can be viewed as the genealogical tree of a stationary continuous-state branching process (associated with the branching mechanism  $\psi_\theta$  defined in (4)), see [9]. To be more precise, for every  $i \in I$  let  $(\ell_a^{(i)}, a \geq 0)$  be the local time measures of the tree  $\tau_i$ . For every  $t \in \mathbb{R}$ , we consider the measure  $\mathbf{Z}_t$  on  $\mathcal{Z}_t(\mathcal{F})$  defined by:

$$(11) \quad \mathbf{Z}_t(dx) = \sum_{i \in I} \mathbf{1}_{\tau_i}(x) \ell_{t-h_i}^{(i)}(dx),$$

and write  $Z_t = \mathbf{Z}_t(1)$  for its total mass which also represents the population size at time  $t$ . For  $\theta = 0$ , we have  $Z_t = +\infty$  a.s. for every  $t \in \mathbb{R}$ . For  $\theta > 0$ , the process  $(Z_t, t \geq 0)$  is a stationary Feller diffusion, solution of the SDE

$$dZ_t = \sqrt{2\beta Z_t} dB_t + 2\beta(1 - \theta Z_t)dt.$$

#### 4. SIMULATION OF THE GENEALOGICAL TREE ( $\theta > 0$ )

We use the representation of trees using ancestral process, see Section 3, which is an atomic measure on  $\mathbb{R}^* \times (0, +\infty)$  satisfying conditions of Definition 3.1.

Under  $\mathbb{P}^{(\theta)}$ , let  $\sum_{i \in I} \delta_{(h_i, \varepsilon_i, e_i)}$  be a Poisson point measure on  $\mathbb{R} \times \{-1, 1\} \times \mathcal{E}$  with intensity  $\beta dh (\delta_{-1}(d\varepsilon) + \delta_1(d\varepsilon)) n^{(\theta)}(de)$ , and let  $\mathcal{F} = ((h_i, \tau_i), i \in I)$  be the associated Brownian forest. We denote by  $\ell_a^{(i)}$  the local time measure of the tree  $\tau_i$  at level  $a$  (recall that this local time is zero for  $a \notin [0, H(\tau_i)]$ ) and we denote by  $\partial_i$  the root of  $\tau_i$ . Recall that the extant population at time  $h \in \mathbb{R}$  is given by  $\mathcal{Z}_h(\mathcal{F})$  defined in Section 2.4.4 and the measure  $\mathbf{Z}_h$  on  $\mathcal{Z}_h(\mathcal{F})$  is defined by (11).

Let  $(\mathfrak{X}_k, k \in \mathbb{N}^*)$  be, conditionally given  $\mathcal{F}$ , independent random variables distributed according to the probability measure  $\mathbf{Z}_0/Z_0$ . Remark that the normalization by  $Z_0$ , which is motivated by the sampling approach, is not usual in the branching setting, see for instance Theorem 4.7 in [9], where the sampling is according to  $\mathbf{Z}_0$  instead leading to the bias factor  $Z_0^n$ .

For every  $k \in \mathbb{N}^*$ , we set  $i_k$  the index in  $I$  such that  $\mathfrak{X}_k \in \tau_{i_k}$ . For every  $n \in \mathbb{N}^*$ , we set  $I_n = \{i_k, 1 \leq k \leq n\}$  and for every  $i \in I_n$ , we denote by  $\tau_i^{(n)}$  the sub-tree of  $\tau_i$  generated by the  $\mathfrak{X}_k$  such that  $i_k = i$  and  $1 \leq k \leq n$ , i.e.

$$\tau_i^{(n)} = \bigcup_{1 \leq k \leq n, i_k = i} [\partial_i, \mathfrak{X}_k].$$

We define the genealogical tree  $T_n$  of  $n$  individuals sampled at random among the population at time 0 by:

$$T_n = (-\infty, 0] \otimes_{i \in I_n} (\tau_i^{(n)}, h_i).$$

Notice that  $T_n \subset T_{n+1}$ . Since the support of  $\mathbf{Z}_h$  is  $\mathcal{Z}_h(\mathcal{F})$  a.s., we get that a.s.  $\text{cl}(\bigcup_{n \in \mathbb{N}^*} T_n) = \mathcal{G}_0(\mathcal{F})$ , where  $\mathcal{G}_0(\mathcal{F})$ , see Definition (7), is the genealogical tree of the forest  $\mathcal{F}$  at time 0.

Recall  $c_\theta$  defined by (5). For  $\delta > 0$ , we will consider in the next sections a positive random variable  $\zeta_\delta^*$  whose distribution is given by, for  $h > 0$ :

$$(12) \quad \mathbb{P}(\zeta_\delta^* < h) = e^{-\delta c_\theta(h)}.$$

This random variable is easy to simulate as, if  $U$  is uniformly distributed on  $[0, 1]$ , then  $\zeta_\delta^*$  has the same distribution as:

$$\frac{1}{2\theta\beta} \log \left( 1 - \frac{2\theta\delta}{\log(U)} \right).$$

This random variable appears naturally in the simulation of the ancestral process of  $\mathcal{F}$  as, if  $\sum_{i \in I} \delta_{(z_i, \zeta_i)}$  is a Poisson point measure on  $\mathbb{R} \times \mathbb{R}_+$  with intensity  $\mathbf{1}_{[0, \delta]}(z) dz |c'_\theta(\zeta)| d\zeta$  (see Proposition 3.6 for the interpretation), then  $\zeta_\delta^*$  is distributed as  $\max_{i \in I} \zeta_i$ .

We now present many ways to simulate  $T_n$ . This will be done by simulating ancestral processes, see Section 3, which code for trees distributed as  $T_n$ .

Recall that for an interval  $I$ , we write  $|I|$  for its length.

**4.1. Static simulation.** In what follows, S stands for static. Assume  $n \in \mathbb{N}^*$  is fixed. We present a way to simulate  $T_n$  under  $\mathbb{P}^{(\theta)}$  with  $\theta > 0$ . See Figures 6 and 7 for an illustration for  $n = 5$ .

- (i) The size of the population on the left (resp. right) of the origin is  $E_g$  (resp.  $E_d$ ), where  $E_g, E_d$  are independent exponential random variables with mean  $1/2\theta$ . Set  $Z_0 = E_g + E_d$  for the total size of the population at time 0. Let  $(U_k, k \in \mathbb{N}^*)$  be independent random variables uniformly distributed on  $[0, 1]$  and independent of  $(E_g, E_d)$ . Set  $X_0 = 0$ , and, for  $k \in \mathbb{N}^*$ ,  $X_k = Z_0 U_k - E_g$  as well as  $\mathcal{X}_k = \{-E_g, E_d, X_0, \dots, X_k\}$ .
- (ii) For  $1 \leq k \leq n$ , set  $X_{k,n}^g = \max\{x \in \mathcal{X}_n, x < X_k\}$  and  $X_{k,n}^d = \min\{x \in \mathcal{X}_n, x > X_k\}$ . We also set  $I_k^S = [X_{k,n}^g, X_k]$  if  $X_k > 0$  and  $I_k^S = [X_k, X_{k,n}^d]$  if  $X_k < 0$ .
- (iii) Conditionally on  $(E_g, E_d, X_1, \dots, X_n)$ , let  $(\zeta_k^S, 1 \leq k \leq n)$  be independent random variables such that for  $1 \leq k \leq n$ ,  $\zeta_k^S$  is distributed as  $\zeta_\delta^*$ , see (12), with  $\delta = |I_k^S|$ . Consider the tree  $\mathfrak{T}_n^S$  corresponding to the ancestral process  $\mathcal{A}_n^S = \sum_{k=1}^n \delta_{(X_k, \zeta_k^S)}$ .

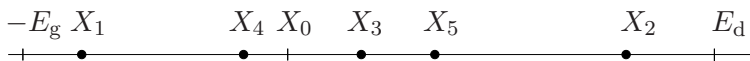


FIGURE 6. One realization of  $E_g, E_d, X_1, \dots, X_5$ .

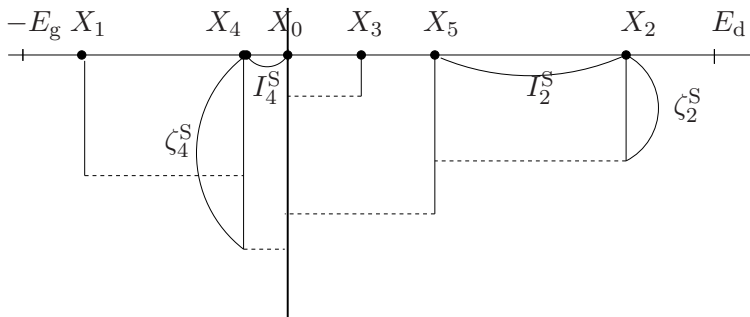


FIGURE 7. One realization of the tree  $\mathfrak{T}_5^S$ .

This gives an exact simulation of the tree  $T_n$  according to the following result.

**Lemma 4.1.** *Let  $\theta > 0$  and  $n \in \mathbb{N}^*$ . The tree  $\mathfrak{T}_n^S$  is distributed as  $T_n$  under  $\mathbb{P}^{(\theta)}$ .*

*Proof.* Let  $B = (B_t, t \in \mathbb{R})$  be the Brownian motion with drift defined in Section 3.2.1 and let  $(L_t, t \in \mathbb{R})$  be its local time at 0 i.e.

$$L_t = L_t^d \mathbf{1}_{t>0} + L_{-t}^g \mathbf{1}_{t<0}.$$

We set  $L_\infty = L_\infty^d + L_\infty^g$  and we consider i.i.d. variables  $(S_1, \dots, S_n)$  distributed according to  $dL_s/L_\infty$ . We denote by  $(S_{(1)}, \dots, S_{(n)})$  the order statistics of  $(S_1, \dots, S_n)$  and, for every  $i \leq n$ , we set

$$M_i = \begin{cases} -\min_{u \in [S_{(i)}, S_{(i+1)} \wedge 0]} B_u & \text{if } S_{(i)} < 0, \\ -\min_{u \in [S_{(i-1)} \vee 0, S_{(i)}]} B_u & \text{if } S_{(i)} > 0. \end{cases}$$



We set  $\mathcal{A}_n = \sum_{1 \leq i \leq n} \delta_{(L_{S(i)}, \zeta_i)}$  which is (see Definition 3.1) an ancestral process and let  $\mathfrak{T}(\mathcal{A}_n)$  be the associated tree. As  $B$  is the contour process of the tree  $\mathcal{F}_{(-\infty, 0]}$ , we get that  $T_n$  and  $\mathfrak{T}(\mathcal{A}_n)$  are equally distributed.

Moreover, by Proposition 3.6, Proposition 3.7 and standard results on Poisson point processes, we get that  $\mathfrak{T}(\mathcal{A}_n)$  and  $\mathfrak{T}_n^S$  are also equally distributed.  $\square$

**4.2. Dynamic simulation (I).** We can modify the static simulation of the previous section to provide a natural dynamic construction of the genealogical tree. In what follows, D stands for static. Let  $\theta > 0$ . We build recursively a family of ancestral processes  $(\mathcal{A}_n, n \in \mathbb{N})$ , with  $\mathcal{A}_0^D = 0$  and  $\mathcal{A}_n^D = \sum_{k=1}^n \delta_{(V_k, \zeta_k^D)}$  for  $n \in \mathbb{N}^*$ .

- (i) Let  $E_g, E_d, (X_n, n \in \mathbb{N})$  and  $(\mathcal{X}_n, n \in \mathbb{N}^*)$  be defined as in (i) of Section 4.1. For  $n \in \mathbb{N}^*$ , set  $X_n^g = \max\{x \in \mathcal{X}_n, x < X_n\}$  and  $X_n^d = \min\{x \in \mathcal{X}_n, x > X_n\}$ .

For  $n \in \mathbb{N}^*$  and  $\ell \in \{g, d\}$ , define the interval  $I_n^\ell = [X_n \wedge X_n^\ell, X_n \vee X_n^\ell]$  and its length  $|I_n^\ell| = |X_n - X_n^\ell|$ .

We shall consider and check by the induction the following hypothesis: for  $n \geq 2$  the random variables  $V_1, \dots, V_{n-1}$  are such that

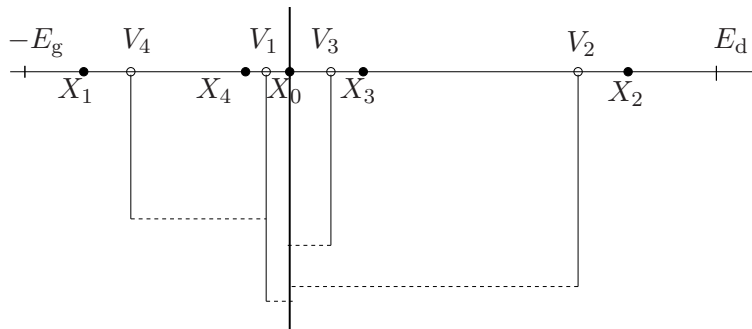
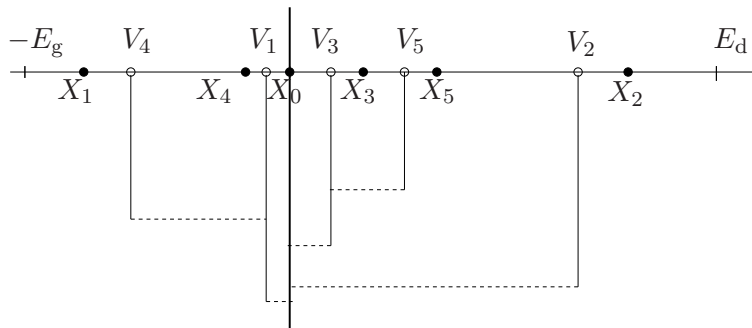
$$(13) \quad X_{(0, n-1)} < V_{(1, n-1)} < X_{(1, n-1)} < \dots < V_{(n-1, n-1)} < X_{(n-1, n-1)},$$

where  $(V_{(1, n-1)}, \dots, V_{(n, n)})$  and  $(X_{(0, n-1)}, \dots, X_{(n-1, n-1)})$  respectively are the order statistics of  $(V_1, \dots, V_{n-1})$  and of  $(X_0, \dots, X_{n-1})$  respectively. Notice that (13) holds trivially for  $n = 1$ .

We set  $\mathcal{W}_n^D = (E_g, E_d, X_1, \dots, X_n, V_1, \dots, V_{n-1}, \zeta_1^D, \dots, \zeta_{n-1}^D)$ .

- (ii) Assume  $n \geq 1$ . We work conditionally on  $\mathcal{W}_n^D$ . On the event  $\{X_n^d = E_d\}$  set  $I_n = I_n^g$  and on the event  $\{X_n^g = -E_g\}$  set  $I_n = I_n^d$ . On the event  $\{X_n^d = E_d\} \cup \{X_n^g = -E_g\}$ , let  $V_n$  be uniform on  $I_n$  and  $\zeta_n^D$  be distributed as  $\zeta_\delta^*$ , see (12), with  $\delta = |I_n|$ .
- (iii) Assume  $n \geq 2$  and that (13) holds. We work conditionally on  $\mathcal{W}_n^D$ . On the event  $\{-E_g < X_n^g, X_n^d < E_d\}$ , there exists a unique integer  $\kappa_n \in \{1, \dots, n-1\}$  such that  $V_{\kappa_n} \in [X_n^g, X_n^d]$ . If  $X_n \in [X_n^g, V_{\kappa_n})$ , set  $I_n = I_n^g$ ; and if  $X_n \in [V_{\kappa_n}, X_n^d]$ , set  $I_n = I_n^d$ . Then, let  $V_n$  be uniform on  $I_n$  and  $\zeta_n^D$  be distributed as  $\zeta_\delta^*$ , with  $\delta = |I_n|$ , conditionally on being less than  $\zeta_{\kappa_n}^D$ .
- (iv) Thanks to (ii) and (iii), notice that (13) holds with  $n-1$  replaced by  $n$ , so that the induction is valid. Set  $\mathcal{A}_n^D = \mathcal{A}_{n-1}^D + \delta_{(V_n, \zeta_n^D)}$  and consider the tree  $\mathfrak{T}_n^D$  corresponding to the ancestral process  $\mathcal{A}_n^D$ .

See Figures 8 and 9 for an instance of  $\mathfrak{T}_4^D$  and  $\mathfrak{T}_5^D$ .

FIGURE 8. An instance of the tree  $\mathfrak{T}_4^D$ .FIGURE 9. An instance of the tree  $\mathfrak{T}_5^D$ . The length of the new branch attached to  $V_5$  is conditioned to be less than the previous branch that was in the considered interval attached to  $V_2$ 

Then we have the following result.

**Lemma 4.2.** *Let  $\theta > 0$ . The sequences of trees  $(\mathfrak{T}_n^D, n \in \mathbb{N}^*)$  and  $(T_n, n \in \mathbb{N}^*)$  under  $\mathbb{P}^{(\theta)}$  have the same distribution.*

*Proof.* We consider  $\sum_{i \in \mathcal{I}} \delta_{(u_i, \zeta_i)}$  the ancestral process associated to the Poisson point measure  $\sum_{i \in \mathcal{I}} \delta_{(h_i, \varepsilon_i, e_i)}$  defined in Section 3.2.2. Let  $(X_k'', k \in \mathbb{N}^*)$  be independent uniform random variables on  $[-E_g, E_d]$ . Set  $X_0'' = 0$ . For  $n \geq 1$ , let us denote by  $(X_{(k,n)}'', 0 \leq k \leq n)$  the order statistic of  $(X_0'', \dots, X_n'')$ .

For every  $n \geq 1$  and every  $1 \leq k \leq n$ , we set  $i_{k,n}$  the index in  $\mathcal{I}$  such that

$$\zeta_{i_{k,n}} = \max_{X_{(k-1,n)}'' \leq u_i < X_{(k,n)}''} \zeta_i.$$

Remark that this index exists since, for every  $\varepsilon > 0$ , the set  $\{i \in \mathcal{I}, \zeta_i > \varepsilon\}$  is a.s. finite. We set  $V_{(k,n)}'' = u_{i_{k,n}}$  and notice that, by standard Poisson point measure properties,  $V_{(k,n)}''$  is, conditionally given  $(X_0'', \dots, X_n'')$ , uniformly distributed on  $[X_{(k-1,n)}'', X_{(k,n)}'']$ . We define

$$\mathcal{A}_n'' = \sum_{k=1}^n \delta_{(V_{(k,n)}'', \zeta_{i_{k,n}})}.$$

By construction, it is easy to check that the order statistics

$$X''_{(0,n)} < V''_{(1,n)} < X''_{(1,n)} < \cdots < V''_{(n,n)} < X''_{(n,n)}$$

is distributed as

$$X_{(0,n)} < V_{(1,n)} < X_{(1,n)} < \cdots < V_{(n,n)} < X_{(n,n)}.$$

For  $1 \leq k \leq n$ , let  $j_{k,n} \in \{1, \dots, n\}$  be the index such that  $V_{(k,n)} = V_{j_{k,n}}$ . By construction, we then deduce that  $((V_{(k,n)}, \zeta_{j_{k,n}}^D), 1 \leq k \leq n), n \in \mathbb{N}^*$  is distributed as  $((V''_{(k,n)}, \zeta_{i_{k,n}}), 1 \leq k \leq n), n \in \mathbb{N}^*$ . This implies that the sequence of ancestral processes  $(\mathcal{A}''_n, n \in \mathbb{N}^*)$  and  $(\mathcal{A}_n, n \in \mathbb{N}^*)$  have the same distribution. Then use Proposition 3.7 to get that the sequence of trees  $(T''_n, n \in \mathbb{N}^*)$ , with  $T''_n$  associated to  $\mathcal{A}''_n$ , is distributed as  $(T_n, n \in \mathbb{N}^*)$ .  $\square$

**4.3. Dynamic simulation (II).** In a sense, we had to introduce another random information corresponding to the position  $V_n$  of the largest spine of the sub-tree containing  $X_n$ . The construction in this sub-section provides a way to remove this additional information (which is now hidden) but at the expense to possibly exchange the new inserted branch with one of its neighbor. In what follows, H stands for hidden. An instance is provided for  $\mathfrak{T}_4^H$  and  $\mathfrak{T}_5^H$  in Figures 10, 11 and 12.

Let  $\theta > 0$ . We build recursively a family of ancestral processes  $(\mathcal{A}_n^H, n \in \mathbb{N})$ , with  $\mathcal{A}_0^H = 0$  and  $\mathcal{A}_n^H = \sum_{k=1}^n \delta_{(X_k, \zeta_{k,n}^H)}$  for  $n \in \mathbb{N}^*$ .

- (i) Let  $E_g, E_d, (X_n, n \in \mathbb{N})$  and  $(\mathcal{X}_n, n \in \mathbb{N}^*)$  be defined as in (i) of Section 4.1. For  $n \in \mathbb{N}^*$ , set  $X_n^g = \max\{x \in \mathcal{X}_n, x < X_n\}$  and  $X_n^d = \min\{x \in \mathcal{X}_n, x > X_n\}$ . For  $n \in \mathbb{N}^*$  and  $\ell \in \{g, d\}$ , define the interval  $I_n^\ell = [X_n \wedge X_n^\ell, X_n \vee X_n^\ell]$  and its length  $|I_n^\ell| = |X_n - X_n^\ell|$ .

We set  $\mathcal{W}_n^H = (E_g, E_d, X_1, \dots, X_n, \zeta_{1,n-1}^H, \dots, \zeta_{n-1,n-1}^H)$ .

- (ii) Assume  $n \geq 1$ . On the event  $\{X_n^d = E_d\}$  set  $I_n = I_n^g$  and on the event  $\{X_n^g = -E_g\}$  set  $I_n = I_n^d$ . Conditionally on  $\mathcal{W}_n^H$ , let  $\zeta_{n,n}^H$  be distributed as  $\zeta_\delta^*$ , see (12), with  $\delta = |I_n|$ ; and for  $1 \leq k \leq n-1$ , set  $\zeta_{k,n}^H = \zeta_{k,n-1}^H$ .
- (iii) Assume  $n \geq 2$ . We work conditionally on  $\mathcal{W}_n^H$ . We define:

$$p_d = \frac{|I_n^d|}{|I_n^d| + |I_n^g|} \quad \text{and} \quad p_g = 1 - p_d = \frac{|I_n^g|}{|I_n^d| + |I_n^g|}.$$

- (a) On the event  $\{0 \leq X_n^g, X_n^d < E_d\}$ , there exists a unique integer  $\kappa_n^d \in \{1, \dots, n-1\}$  such that  $X_{\kappa_n^d}^d = X_n^d$ . For  $1 \leq k \leq n-1$  and  $k \neq \kappa_n^d$ , set  $\zeta_{n,k}^H = \zeta_{n-1,k}^H$ . Write  $\zeta_n^H = \zeta_{n-1, \kappa_n^d}^H$ .

With probability  $p_d$ , set  $\zeta_{n, \kappa_n^d}^H = \zeta_n^H$  and let  $\zeta_{n,n}^H$  be distributed as  $\zeta_\delta^*$ , with  $\delta = |I_n^g|$ , conditionally on being less than  $\zeta_n^H$ .

With probability  $p_g$ , set  $\zeta_n^H = \zeta_{n, \kappa_n^d}^H$  and let  $\zeta_{n, \kappa_n^d}^H$  be distributed as  $\zeta_\delta^*$ , with  $\delta = |I_n^d|$ , conditionally on being less than  $\zeta_n^H$ .

- (b) On the event  $\{-E_g < X_n^g, X_n^d \leq 0\}$ , there exists a unique integer  $\kappa_n^g \in \{1, \dots, n-1\}$  such that  $X_{\kappa_n^g}^g = X_n^g$ . For  $1 \leq k \leq n-1$  and  $k \neq \kappa_n^g$ , set  $\zeta_{n,k}^H = \zeta_{n-1,k}^H$ . Write  $\zeta_n^H = \zeta_{n-1, \kappa_n^g}^H$ .

With probability  $p_g$ , set  $\zeta_{n, \kappa_n^g}^H = \zeta_n^H$  and let  $\zeta_{n,n}^H$  be distributed as  $\zeta_\delta^*$ , with  $\delta = |I_n^d|$ , conditionally on being less than  $\zeta_n^H$ .

With probability  $p_d$ , set  $\zeta_n^H = \zeta_{n, \kappa_n^g}^H$  and let  $\zeta_{n, \kappa_n^g}^H$  be distributed as  $\zeta_\delta^*$ , with  $\delta = |I_n^g|$ , conditionally on being less than  $\zeta_n^H$ .

- (iv) Let  $\mathfrak{T}_n^H$  be the tree corresponding to the ancestral process  $\mathcal{A}_n^H = \sum_{k=1}^n \delta_{(X_k, \zeta_{k,n}^H)}$ .

We have the next result.

**Lemma 4.3.** *Let  $\theta > 0$ . The sequences of trees  $(\mathfrak{T}_n^H, n \in \mathbb{N}^*)$  and  $(T_n, n \in \mathbb{N}^*)$  under  $\mathbb{P}^{(\theta)}$  have the same distribution.*

*Proof.* The proof is left to the reader. It is in the same spirit as the proof of Lemma 4.2, but here we consider the random variables  $((V''_{(k,n)}, 1 \leq k \leq n), n \in \mathbb{N}^*)$  as unobserved.  $\square$

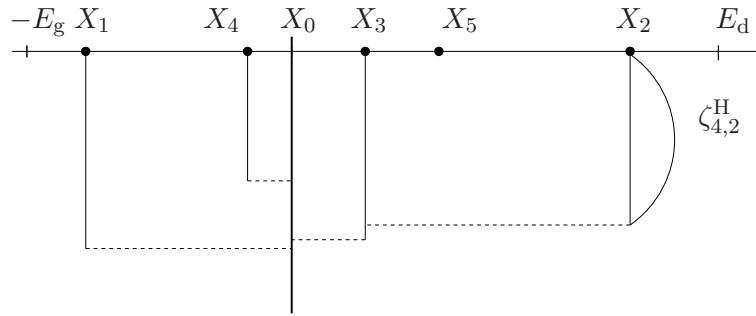


FIGURE 10. An instance of the tree  $\mathfrak{T}_4^H$  with the new individual  $X_5$ .

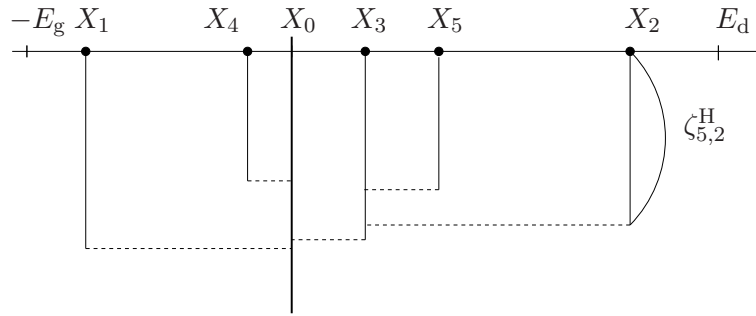


FIGURE 11. An instance of the tree  $\mathfrak{T}_5^H$  with  $\mathfrak{T}_4^H$  given in Figure 10 and the event associated with  $p_d$  (a new segment is attached to  $X_5$ ).

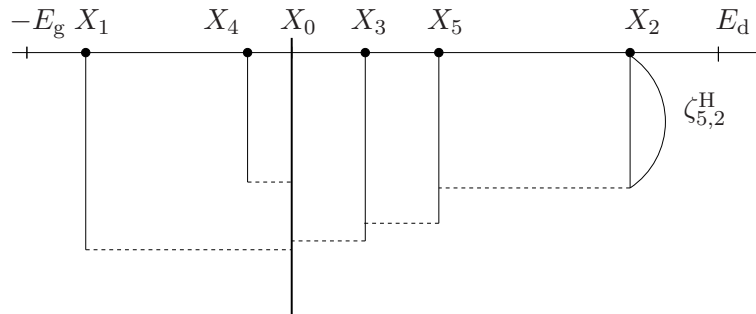


FIGURE 12. An instance of the tree  $\mathfrak{T}_5^H$  with  $\mathfrak{T}_4^H$  given in Figure 10 and the event associated with  $p_g$  (the segment previously attached to  $X_2$  is now attached to  $X_5$  and a new segment is attached to  $X_2$ ).

**4.4. Simulation of genealogical tree conditionally on its maximal height.** Let  $\mathcal{F} = ((\tau_i, h_i), i \in I)$  be a Brownian forest under  $\mathbb{P}^{(\theta)}$ . Recall the definition of  $A_0$  the time to the MRCA of the population living at time 0 given in (16). The goal of this section is to simulate the genealogical tree  $T_n$  of  $n$  individuals uniformly sampled in the population living at time 0, conditionally given the time to the MRCA of the whole population is  $h$ , that is given  $A_0 = h$ .

Let  $\mathcal{A}(du, d\zeta) = \sum_{j \in \mathcal{I}} \delta_{(u_j, \zeta_j)}(du, d\zeta)$  be the ancestral process of Definition 3.1. Recall the notations  $E_g, E_d$  from Section 3.2.2. Let  $\zeta_{\max} = \sup\{\zeta_j, j \in \mathcal{I}\}$  and define the random index  $J_0 \in \mathcal{I}$  such that  $\zeta_{\max} = \zeta_{J_0}$ . Note that  $J_0$  is well defined since for every  $\varepsilon > 0$ , the set  $\{j \in \mathcal{I}, \zeta_j > \varepsilon\}$  is finite. We set  $X = u_{J_0} \in (-E_g, E_d)$ . Remark that  $\zeta_{\max}$  is distributed as  $A_0$ .

For  $r \in \mathbb{R}$ , let  $r_+ = \max(0, r)$  and  $r_- = \max(0, -r)$  be respectively the positive and negative part of  $r$ . The proof of the next lemma is postponed to the end of this section.

**Lemma 4.4.** *Let  $\theta > 0$ . Under  $\mathbb{P}^{(\theta)}$ , conditionally given  $\zeta_{\max} = h$ , the random variables  $E_g + X_-$ ,  $|X|$ ,  $E_d - X_+$  and  $\mathbf{1}_{\{X \geq 0\}}$  are independent;  $E_g + X_-$ ,  $|X|$  and  $E_d - X_+$  are exponentially distributed with parameter  $2\theta + c_\theta(h)$  and  $\mathbf{1}_{\{X \geq 0\}}$  is Bernoulli  $1/2$ .*

Let  $h > 0$  be fixed. For  $\delta > 0$ , let  $\zeta_\delta^{*,h}$  be a positive random variable distributed as  $\zeta_\delta^*$  conditionally on  $\{\zeta_\delta^* \leq h\}$ , i.e., for  $0 \leq u \leq h$ :

$$\mathbb{P}(\zeta_\delta^{*,h} \leq u) = \mathbb{P}(\zeta_\delta^* \leq u \mid \zeta_\delta^* \leq h) = e^{-\delta(c_\theta(u) - c_\theta(h))}.$$

Then the static simulation runs as follows.

- (i) Simulate three independent random variables  $E_1, E_2, E_3$  exponentially distributed with parameter  $2\theta + c_\theta(h)$ , and another independent Bernoulli variable  $\xi$  with parameter  $1/2$ . If  $\xi = 0$ , set  $E_g = E_1$ ,  $X = E_2$ ,  $E_d = E_2 + E_3$ , and if  $\xi = 1$ , set  $E_g = E_1 + E_2$ ,  $X = -E_2$ ,  $E_d = E_3$ . Let  $X_k$  and  $\mathcal{X}_k$  be defined as in (i) of Section 4.1 for  $1 \leq k \leq n$ .
- (ii) Let the intervals  $I_k^S$  be defined as in (ii) of Section 4.1 for  $1 \leq k \leq n$ .
- (iii) Conditionally on  $(E_g, E_d, X, X_1, \dots, X_n)$ , let  $(\zeta_k^h, 1 \leq k \leq n)$  be independent random variables such that, for  $1 \leq k \leq n$ ,  $\zeta_k^h$  is distributed as  $\zeta_\delta^{*,h}$  with  $\delta = |I_k^S|$  if  $X \notin I_k^S$ ; and  $\zeta_k^h = h$  if  $X \in I_k^S$ . Consider the tree  $\mathfrak{T}_n^h$  corresponding to the ancestral process  $\mathcal{A}_n^h = \sum_{k=1}^n \delta_{(X_k, \zeta_k^h)}$ .

The proof of the following result which relies on Lemma 4.4 is similar to the one of Lemma 4.1, and is not reproduced here.

**Lemma 4.5.** *Let  $\theta > 0$ ,  $h > 0$  and  $n \in \mathbb{N}^*$ . The tree  $\mathfrak{T}_n^h$  is distributed as  $T_n$  under  $\mathbb{P}^{(\theta)}$  conditionally given  $A_0 = h$ .*

Notice that the height of  $\mathfrak{T}_n^h$  is less than or equal to  $h$ . When strictly less than  $h$ , it means that no individual of the oldest family has been sampled.

*Proof of Lemma 4.4.* By Proposition 3.6, the pair  $E = (E_g, E_d)$  under  $\mathbb{P}^{(\theta)}$  has density:

$$f_E(e_g, e_d) = (2\theta)^2 e^{-2\theta(e_g + e_d)} \mathbf{1}_{\{e_g \geq 0, e_d \geq 0\}}.$$

Moreover, by standard results on Poisson point measures, the conditional density of the pair  $(X, \zeta_{\max})$  given  $(E_d, E_g) = (e_d, e_g)$  exists and is:

$$\begin{aligned} f_{X, \zeta_{\max}}^{E=(e_g, e_d)}(x, h) &= \frac{1}{e_g + e_d} \mathbf{1}_{[-e_g, e_d]}(x) (e_g + e_d) |c'_\theta(h)| e^{-c_\theta(h)(e_g + e_d)} \mathbf{1}_{\{h \geq 0\}} \\ &= \mathbf{1}_{[-e_g, e_d]}(x) |c'_\theta(h)| e^{-c_\theta(h)(e_g + e_d)} \mathbf{1}_{\{h \geq 0\}}. \end{aligned}$$

We deduce that the vector  $(E_g, E_d, X, \zeta_{\max})$  has density:

$$f(e_g, e_d, x, h) = (2\theta)^2 |c'_\theta(h)| e^{-(2\theta+c_\theta(h))(e_g+e_d)} \mathbf{1}_{\{e_g \geq 0, e_d \geq 0, -e_g \leq x \leq e_d, h \geq 0\}}$$

and that the random variable  $\zeta_{\max}$  has density:

$$\begin{aligned} f_{\zeta_{\max}}(h) &= \int (2\theta)^2 |c'_\theta(h)| e^{-(2\theta+c_\theta(h))(e_g+e_d)} \mathbf{1}_{\{e_g \geq 0, e_d \geq 0, -e_g \leq x \leq e_d, h \geq 0\}} de_g de_d dx \\ &= (2\theta)^2 |c'_\theta(h)| \frac{2}{(2\theta + c_\theta(h))^3} \mathbf{1}_{\{h \geq 0\}}. \end{aligned}$$

Therefore, the conditional density of the vector  $(E_g, E_d, X)$  given  $\zeta_{\max} = h$  is:

$$f_{E, X}^{\zeta_{\max}=h}(e_g, e_d, x) = \frac{1}{2} (2\theta + c_\theta(h))^3 e^{-(2\theta+c_\theta(h))(e_g+e_d)} \mathbf{1}_{\{e_g \geq 0, e_d \geq 0, -e_g \leq x \leq e_d\}}.$$

For any nonnegative measurable function  $\varphi$ , we have:

$$\begin{aligned} \mathbb{E}^{(\theta)}[\varphi(E_g + X_-, |X|, E_d - X_+) \mathbf{1}_{\{X \geq 0\}} \mid \zeta_{\max} = h] \\ &= \mathbb{E}^{(\theta)}[\varphi(E_g, X, E_d - X) \mathbf{1}_{\{X \geq 0\}} \mid \zeta_{\max} = h] \\ &= \int \varphi(e_g, x, e_d - x) \frac{1}{2} (2\theta + c_\theta(h))^3 e^{-(2\theta+c_\theta(h))(e_g+e_d)} \mathbf{1}_{\{e_g \geq 0, e_d \geq x \geq 0\}} de_g de_d dx \\ &= \int \varphi(e_1, e_2, e_3) \frac{1}{2} (2\theta + c_\theta(h))^3 e^{-(2\theta+c_\theta(h))(e_1+e_2+e_3)} \mathbf{1}_{\{e_1 \geq 0, e_2 \geq 0, e_3 \geq 0\}} de_1 de_2 de_3, \end{aligned}$$

using an obvious change of variables. Similarly, we get:

$$\begin{aligned} \mathbb{E}^{(\theta)}[\varphi(E_g + X_-, |X|, E_d - X_+) \mathbf{1}_{\{X < 0\}} \mid \zeta_{\max} = h] \\ &= \mathbb{E}^{(\theta)}[\varphi(E_g + X_-, |X|, E_d - X_+) \mathbf{1}_{\{X \geq 0\}} \mid \zeta_{\max} = h]. \end{aligned}$$

This proves the lemma.  $\square$

## 5. RENORMALIZED TOTAL LENGTH OF THE GENEALOGICAL TREE

Let  $\mathcal{F} = ((h_i, \tau_i), i \in I)$  be a Brownian forest under  $\mathbb{P}^{(\theta)}$  with  $\theta > 0$ . Recall that the tree  $\mathcal{F}_{(-\infty, 0]}$  belongs to  $\mathbb{T}_1$ . For a forest  $\mathbf{f} \in \mathbb{T}_1$ , recall that  $\mathcal{Z}_h(\mathbf{f})$  denotes the set of vertices of  $\mathcal{F}_{(-\infty, 0]}$  at level  $h$ . We shall also consider  $\mathcal{Z}_h^*(\mathbf{f}) = \mathcal{Z}_h(\mathbf{f}) \cap \mathcal{S}(\mathbf{f}_{(-\infty, h]})^c$  the extant population at time  $h$  except the point on the semi-infinite branch  $(-\infty, h]$ . For  $r \leq h$ , we define the set of ancestors at time  $r$  in the past of the extant population at time  $h$  forgetting the individual in the infinite spine:

$$(14) \quad \mathcal{M}_r^h(\mathbf{f}) = \mathcal{G}_h(\mathbf{f}) \cap \mathcal{Z}_r^*(\mathbf{f})$$

and its cardinality

$$(15) \quad M_r^h(\mathbf{f}) = \text{Card}(\mathcal{M}_r^h(\mathbf{f})).$$

We also define the time to the MRCA of  $\mathcal{Z}_t(\mathcal{F})$  as

$$(16) \quad A_t = t - \sup \{r \leq t; M_r^t = 0\}.$$

We want to define the length of the genealogical tree  $\mathcal{G}_t(\mathcal{F})$  of all extant individuals at time  $t$  (which is a.s. infinite) by approximating this genealogical tree by trees with finite length and take compensated limits. Without loss of generality we can take  $t = 0$  (since the distribution of the Brownian forest is invariant by time translation).

Two approximations may be considered here. The first one is to consider for  $\varepsilon > 0$  the genealogical tree of individuals at time  $t - \varepsilon$ , with descendants at time  $t$ , and let  $\varepsilon$  goes down to 0. We define the total length of the genealogical tree of the current population up to  $\varepsilon > 0$  in the past as:

$$(17) \quad L_\varepsilon = \int_\varepsilon^\infty M_{-s}^0 ds.$$

Set  $L = (L_\varepsilon, \varepsilon > 0)$ . According to [6], we have  $\mathbb{E}[L_\varepsilon|Z_0] = -Z_0 \log(2\beta\theta\varepsilon)/\beta + O(\varepsilon)$  (see also (21) as  $\tilde{L}_\varepsilon$  is distributed as  $L_\varepsilon$ ), and that the sequence  $(L_\varepsilon - \mathbb{E}[L_\varepsilon|Z_0], \varepsilon > 0)$  converges a.s. as  $\varepsilon$  goes down to zero towards a limit say  $\mathcal{L}$ . Furthermore, for all  $\lambda > 0$ ,

$$\mathbb{E} \left[ e^{-\lambda\mathcal{L}} | Z_0 \right] = e^{\theta Z_0 \varphi(\lambda/(2\beta\theta))}, \quad \text{with} \quad \varphi(\lambda) = \lambda \int_0^1 \frac{1-v^\lambda}{1-v} dv.$$

The second approximation consists in looking at the genealogical tree associated with  $n$  individuals picked at random in the population at time 0. Recall Definition (11) of  $\mathbf{Z}_h$ . Let  $(X_k, k \in \mathbb{N}^*)$  be, conditionally on  $\mathcal{F}$ , independent random variables with distribution  $\mathbf{Z}_0(dx)/Z_0$ . This models individuals uniformly chosen among the population living at time 0. Define the ancestors of  $X_1, \dots, X_n$  at time  $s < 0$  as:

$$\mathcal{M}_s^{(n)}(\mathcal{F}) = \{x \in \mathcal{M}_s^0(\mathcal{F}); x \prec X_i \text{ for some } 1 \leq i \leq n\},$$

and  $M_s^{(n)} = \text{Card}(\mathcal{M}_s^{(n)}(\mathcal{F}))$  its cardinality. We define the total length of the genealogical tree of  $n$  individuals uniformly chosen in the current population as:

$$(18) \quad \Lambda_n = \int_0^\infty M_{-s}^{(n)} ds.$$

Set  $\Lambda = (\Lambda_n, n \in \mathbb{N}^*)$ . The next theorem states that the two approximations give the same limit a.s.

**Theorem 5.1.** *The sequence  $(\Lambda_n - \mathbb{E}[\Lambda_n|Z_0], n \in \mathbb{N}^*)$  converges a.s. and in  $L^2$  towards  $\mathcal{L}$  as  $n$  tends to  $+\infty$ . And we also have  $\mathbb{E}[\Lambda_n|Z_0] = \frac{Z_0}{\beta} \log\left(\frac{n}{2\theta Z_0}\right) + O(n^{-1} \log(n))$ .*

The rest of the section is devoted to the proof of this theorem.

**5.1. Preliminary results.** Let  $E_g$  and  $E_d$  be two independent exponential random variable with parameter  $2\theta$ . Let  $\mathcal{N} = \sum_{i \in I} \delta_{z_i, \tau_i}$  be, conditionally given  $(E_g, E_d)$ , distributed as a Poisson point measure with intensity  $\mathbf{1}_{[-E_g, E_d]}(z) dz \mathbb{N}^{(\theta)}[d\tau]$ . We define  $\tilde{L} = (\tilde{L}_\varepsilon, \varepsilon > 0)$  with:

$$\tilde{L}_\varepsilon = \sum_{i \in I} (\zeta_i - \varepsilon)_+,$$

where  $\zeta_i = H(\tau_i)$  is the height of  $\tau_i$ . Let  $(U_k, k \in \mathbb{N}^*)$  be independent random variables uniformly distributed on  $[0, 1]$  and independent of  $(\mathcal{N}, E_g, E_d)$ . We set  $X_0 = 0$ , and  $X_k = (E_g + E_d)U_k - E_g$  for  $k \in \mathbb{N}^*$ . Fix  $n \in \mathbb{N}^*$ . Let  $X_{(0,n)} \leq \dots \leq X_{(n,n)}$  be the corresponding order statistic of  $(X_0, \dots, X_n)$ . We set  $X_{(-1,n)} = -E_g$  and  $X_{(n+1,n)} = E_d$ . We define the interval  $I_{k,n} = (X_{(k-1,n)}, X_{(k,n)})$  and its length  $\Delta_{k,n} = X_{(k,n)} - X_{(k-1,n)}$  for  $0 \leq k \leq n+1$ . We set  $\tilde{\Lambda}_n = (\tilde{\Delta}_{k,n}, 0 \leq k \leq n+1)$ . For  $1 \leq k \leq n$ , we define  $\tilde{\Lambda} = (\tilde{\Lambda}_n, n \in \mathbb{N}^*)$  by:

$$\tilde{\Lambda}_n = \sum_{k=1}^n \zeta_{k,n}^* \quad \text{with} \quad \zeta_{k,n}^* = \max\{\zeta_i; z_i \in I_{k,n}\}.$$

Recall the definitions of  $Z_0$  in (11),  $L = (L_\varepsilon, \varepsilon > 0)$  in (17) and  $\Lambda = (\Lambda_n, n \in \mathbb{N}^*)$  in (18). Thanks to Proposition 3.6, we deduce that  $(Z_0, L, \Lambda)$  is distributed as  $(E_g + E_d, \tilde{L}, \tilde{\Lambda})$ . So to

prove Theorem 5.1, it is enough to prove the statement with  $\tilde{\Lambda}$  instead of  $\Lambda$ .

For convenience, we set  $Z_0 = E_g + E_d$ . Elementary computations give the following lemma. Recall that  $z_+ = \max(z, 0)$ .

**Lemma 5.2.** *Let  $\theta > 0$  and  $\varepsilon > 0$ . We have:*

$$(19) \quad \mathbb{N}^{(\theta)}[(\zeta - \varepsilon)_+] = \int_{\varepsilon}^{\infty} c_{\theta}(h) dh = -\frac{1}{\beta} \log(2\beta\theta\varepsilon) + O(\varepsilon),$$

$$(20) \quad \mathbb{N}^{(\theta)}[(\zeta - \varepsilon)_+^2] = 2 \int_{\varepsilon}^{\infty} hc_{\theta}(h) dh - 2\varepsilon \int_{\varepsilon}^{\infty} c_{\theta}(h) dh = 2 \int_0^{\infty} hc_{\theta}(h) dh + O(\varepsilon \log(\varepsilon)).$$

We deduce that:

$$(21) \quad \mathbb{E}[\tilde{L}_{\varepsilon}|Z_0] = -\frac{Z_0}{\beta} \log(2\beta\theta\varepsilon) + O(\varepsilon),$$

$$(22) \quad \mathbb{E}[\tilde{L}_{\varepsilon}^2|Z_0] = 2Z_0 \int_0^{\infty} hc_{\theta}(h) dh + \mathbb{E}[\tilde{L}_{\varepsilon}|Z_0]^2 + O(\varepsilon \log(\varepsilon)),$$

where we used that if  $\sum_{i \in I} \delta_{x_i}$  is a Poisson point measure with intensity  $\mu(dx)$ , then:

$$(23) \quad \mathbb{E} \left[ \left( \sum_{i \in I} f(x_i) \right)^2 \right] = \mu(f^2) + \mu(f)^2.$$

Eventually, let us notice that with the change of variable  $u = c_{\theta}(h)$  (so that  $dh = du/\beta u(u+2\theta\delta)$ ), we have:

$$(24) \quad 2 \int_0^{\infty} hc_{\theta}(h) dh = \frac{1}{\beta^2\theta} \int_0^{\infty} \frac{\log(v+1)}{v(v+1)} dv.$$

Recall the definition of  $\zeta_{\delta}^*$  for  $\delta > 0$ , see (12). Let  $\gamma$  be the Euler constant, and thus:

$$\gamma = - \int_0^{+\infty} \log(u) e^{-u} du.$$

We have the following lemma.

**Lemma 5.3.** *Let  $\delta > 0$ . We have:*

$$(25) \quad \mathbb{E}[\zeta_{\delta}^*] = -\frac{\delta}{\beta} \log(2\theta\delta) + \frac{\delta}{\beta}(1 - \gamma) + \frac{\delta}{\beta} g_1(2\theta\delta),$$

with  $|g_1(x)| \leq x(|\log(x)| + 2)$  for  $x > 0$  and

$$(26) \quad \mathbb{E}[(\zeta_{\delta}^*)^2] = 2\delta \int_0^{\infty} hc_{\theta}(h) dh + \frac{\delta}{\beta^2\theta} g_2(2\theta\delta),$$

with  $|g_2(x)| \leq x(|\log(x)| + 2)$  for  $x > 0$ . We also have:

$$(27) \quad \mathbb{E} \left[ \zeta_{\delta}^* \sum_{i \in I} (\zeta_i - \varepsilon)_+ \right] = 2\delta \int_0^{\infty} hc_{\theta}(h) dh + g_3(\delta)$$

and there exists a finite constant  $c$  such that for all  $x > 0$  and  $\varepsilon \in (0, 1]$ , we have  $|g_3(x)| \leq cx^2(1+x)(|\log(x)| + 1)(|\log(\varepsilon)| + 1) + c\varepsilon x(|\log(x)| + 1)(1+x) + \varepsilon^2$ .

The end of this section is devoted to the proof of Lemma 5.3.



5.1.1. *Proof of (25).* Using (12), we get:

$$(28) \quad \mathbb{E}[\zeta_\delta^*] = \int_0^\infty \mathbb{P}(\zeta_\delta^* > h) dh = \int_0^\infty (1 - e^{-\delta c_\theta(h)}) dh = \frac{\delta}{\beta} \int_0^\infty (1 - e^{-u}) \frac{du}{u(u + 2\theta\delta)},$$

where we used the change of variable  $u = \delta c_\theta(h)$ . It is easy to check that for  $a > 0$ ,

$$(29) \quad \log(1 + a) \leq |\log(a)| + \log(2).$$

Let  $a > 0$ . We have:

$$\begin{aligned} \int_0^1 (1 - e^{-u}) \frac{du}{u(u + a)} &= \int_0^1 (1 - u - e^{-u}) \frac{du}{u(u + a)} + \log(1 + a) - \log(a) \\ &= \int_0^1 (1 - u - e^{-u}) \frac{du}{u^2} + \log(1 + a) - \log(a) + ag_{1,0}(a), \end{aligned}$$

with

$$g_{1,0}(a) = - \int_0^1 (1 - u - e^{-u}) \frac{du}{u^2(u + a)} \leq \int_0^1 \frac{du}{2(u + a)} = \frac{1}{2}(\log(1 + a) - \log(a)) \leq |\log(a)| + \frac{1}{2}$$

and  $g_{1,0}(a) \geq 0$ , where we used that  $0 \leq -(1 - u - e^{-u}) \leq u^2/2$  for  $u \geq 0$ . We also have:

$$\int_1^\infty (1 - e^{-u}) \frac{du}{u(u + a)} = \int_1^\infty (1 - e^{-u}) \frac{du}{u^2} - ag_{1,1}(a),$$

with

$$g_{1,1}(a) = \int_1^\infty (1 - e^{-u}) \frac{du}{u^2(u + a)} \leq \int_1^\infty \frac{du}{u^3} \leq \frac{1}{2}.$$

Notice that, by integration by parts, we have:

$$\int_0^1 (1 - u - e^{-u}) \frac{du}{u^2} + \int_1^\infty (1 - e^{-u}) \frac{du}{u^2} = e^{-1} + \int_0^1 \log(u) e^{-u} du + 1 - e^{-1} + \int_1^\infty \log(u) e^{-u} du = 1 - \gamma.$$

We deduce that:

$$\int_0^\infty (1 - e^{-u}) \frac{du}{u(u + a)} = 1 - \gamma - \log(a) + g_1(a)$$

with  $g_1(a) = \log(1 + a) + ag_{1,0}(a) - ag_{1,1}(a)$  and

$$|g_1(a)| = |\log(1 + a) + ag_{1,0}(a) - ag_{1,1}(a)| \leq a(|\log(a)| + 2).$$

Then, use (28) to get (25).

5.1.2. *Proof of (26).* Using (12), we get:

$$(30) \quad \mathbb{E}[(\zeta_\delta^*)^2] = 2 \int_0^\infty h(1 - e^{-\delta c_\theta(h)}) dh = 2 \frac{\delta}{\beta} \int_0^\infty \frac{1}{2\beta\theta} \log\left(\frac{u + 2\theta\delta}{u}\right) (1 - e^{-u}) \frac{du}{u(u + 2\theta\delta)},$$

where we used the change of variable  $u = \delta c_\theta(h)$ . Let  $a > 0$ . We set:

$$g_{2,1}(a) = \int_1^\infty \log\left(\frac{u + a}{u}\right) (1 - e^{-u}) \frac{du}{u(u + a)}.$$

We have using that  $0 \leq \log(1 + x) \leq x$  for  $x > 0$ :

$$|g_{2,1}(a)| \leq a \int_1^\infty \frac{du}{u^3} \leq \frac{a}{2}.$$

We also have:

$$\begin{aligned} \int_0^1 \log\left(\frac{u+a}{u}\right) (1-e^{-u}) \frac{du}{u(u+a)} &= \int_0^1 \log\left(\frac{u+a}{u}\right) \frac{du}{u+a} + g_{2,2}(u) \\ &= \int_0^\infty \frac{\log(v+1)}{v(v+1)} dv - g_{2,3}(a) + g_{2,2}(a), \end{aligned}$$

with the change of variable  $v = a/u$  as well as:

$$g_{2,2}(a) = \int_0^1 \log\left(\frac{u+a}{u}\right) (1-u-e^{-u}) \frac{du}{u(u+a)} \quad \text{and} \quad g_{2,3}(a) = \int_0^a \frac{\log(v+1)}{v(v+1)} dv.$$

We have, using  $\log(1+v) \leq v$  for  $v > 0$  (twice), that:

$$0 \leq g_{2,3}(a) \leq \int_0^a \frac{dv}{v+1} \leq a.$$

We have, using  $|1-u-e^{-u}| \leq u^2/2$  if  $u > 0$  for the first inequality and (29) for the last, that:

$$|g_{2,2}(a)| \leq \frac{1}{2} \int_0^1 \log\left(1 + \frac{a}{u}\right) \frac{u du}{(u+a)} \leq \frac{a}{2} \int_0^1 \frac{du}{(u+a)} \leq a(|\log(a)| + \frac{1}{2}).$$

We deduce that:

$$\int_0^\infty \log\left(\frac{u+a}{u}\right) (1-e^{-u}) \frac{du}{u(u+a)} = \int_0^\infty \frac{\log(v+1)}{v(v+1)} dv + g_2(a)$$

and

$$|g_2(a)| = |g_{2,1}(a) - g_{2,3}(a) + g_{2,2}(a)| \leq a(|\log(a)| + 2).$$

Then, use (30) as well as the identity (24) to get (26).

5.1.3. *Proof of (27).* Using properties of Poisson point measures, we get that if  $\sum_{j \in J} \delta_{\zeta_j}$  is a Poisson point measure with intensity  $\delta \mathbb{N}[d\zeta]$  and  $\zeta_\delta^* = \max_{j \in J} \zeta_j$ , then for any measurable non-negative functions  $f$  and  $g$ , we have:

$$\mathbb{E} \left[ f(\zeta_\delta^*) e^{-\sum_{j \in J} g(\zeta_j)} \right] = \mathbb{E} \left[ f(\zeta_\delta^*) e^{-g(\zeta_\delta^*) - G(\zeta_\delta^*)} \right] \quad \text{with} \quad G(r) = \delta \mathbb{N} \left[ (1 - e^{-g(\zeta)}) \mathbf{1}_{\{\zeta < r\}} \right].$$

We deduce that:

$$\mathbb{E} \left[ \zeta_\delta^* \sum_{i \in I} (\zeta_i - \varepsilon)_+ \right] = \mathbb{E}[\zeta_\delta^* (\zeta_\delta^* - \varepsilon)_+] + \delta g_{3,1}(\delta),$$

with  $g_{3,1}(\delta) = \mathbb{E} \left[ \zeta_\delta^* \mathbb{N} \left[ (\zeta - \varepsilon)_+ \mathbf{1}_{\{\zeta < h\}} \right] \Big|_{h=\zeta_\delta^*} \right]$ . According to (25), there exists a finite constant  $c > 0$  such that for all  $\delta > 0$ , we have  $\mathbb{E}[\zeta_\delta^*] \leq c\delta(|\log(\delta)| + 1)(1 + \delta)$ . We deduce from (19) that there exists a finite constant  $c$  independent of  $\delta > 0$  and  $\varepsilon \in (0, 1]$  such that:

$$g_{3,1}(\delta) \leq \mathbb{E}[\zeta_\delta^*] \mathbb{N}[(\zeta - \varepsilon)_+] \leq c\delta(|\log(\delta)| + 1)(1 + \delta)(|\log(\varepsilon)| + 1).$$

We also have:

$$\mathbb{E}[\zeta_\delta^* (\zeta_\delta^* - \varepsilon)_+] = \mathbb{E}[(\zeta_\delta^*)^2] - \mathbb{E}[(\zeta_\delta^*)^2 \mathbf{1}_{\{\zeta_\delta^* < \varepsilon\}}] - \varepsilon \mathbb{E}[\zeta_\delta^* \mathbf{1}_{\{\zeta_\delta^* > \varepsilon\}}] = 2\delta \int_0^\infty hc_\theta(h) dh + g_{3,2}(\varepsilon, \delta),$$

with, thanks to (25) and (26),  $|g_{3,2}(\varepsilon, \delta)| \leq c\delta^2(|\log(\delta)| + 1) + \varepsilon^2 + c\varepsilon\delta(|\log(\delta)| + 1)(1 + \delta)$ , for some finite constant  $c$  independent of  $\delta > 0$  and  $\varepsilon > 0$ . We deduce that:

$$\mathbb{E} \left[ \zeta_\delta^* \sum_{i \in I} (\zeta_i - \varepsilon)_+ \right] = 2\delta \int_0^\infty hc_\theta(h) dh + g_3(\delta)$$

and for some finite constant  $c$  independent of  $\delta > 0$  and  $\varepsilon \in (0, 1]$ .

$$|g_3(\delta)| \leq c\delta^2(1+\delta)(|\log(\delta)|+1)(|\log(\varepsilon)|+1) + c\varepsilon\delta(|\log(\delta)|+1)(1+\delta) + \varepsilon^2.$$

**5.2. A technical lemma.** An elementary induction gives for  $n \in \mathbb{N}$  that:

$$\int_0^1 (1-x)^n |\log(x)| dx = \frac{H_{n+1}}{n+1} \quad \text{and} \quad \int_0^1 (1-x)^n \log^2(x) dx = \frac{2}{n+1} \sum_{k=1}^{n+1} \frac{H_k}{k},$$

where  $H_n = \sum_{k=1}^n k^{-1}$  is the harmonic sum. Recall that  $H_n = \log(n) + \gamma + (2n)^{-1} + O(n^{-2})$ . So we deduce that:

$$(31) \quad (n+1) \int_0^1 (1-x)^n |\log(x)| dx = \log(n) + \gamma + \frac{3}{2n} + O(n^{-2}).$$

It is also easy to deduce that for  $a, b \in \{1, 2\}$ :

$$(32) \quad \int_0^1 x^a (1-x)^n |\log(x)|^b dx = O\left(\frac{\log^b(n)}{n^{a+1}}\right).$$

Recall  $\tilde{\Lambda}_n$  and  $\Delta_n$  defined in Section 5.1. We give a technical lemma.

**Lemma 5.4.** *We have:*

$$(33) \quad \mathbb{E}[\tilde{\Lambda}_n | \Delta_n] = \frac{Z_0}{\beta}(1-\gamma) - \sum_{k=1}^n \frac{\Delta_{k,n}}{\beta} \log(2\theta\Delta_{k,n}) + W_n,$$

with  $\mathbb{E}[|W_n| | Z_0] = O(n^{-1} \log(n))$  and

$$(34) \quad \mathbb{E}[\tilde{\Lambda}_n | Z_0] = \frac{Z_0}{\beta} \log\left(\frac{n}{2\theta Z_0}\right) + O(n^{-1} \log(n)).$$

We have also:

$$(35) \quad \mathbb{E}[\tilde{\Lambda}_n^2 | Z_0] = 2Z_0 \int_0^\infty hc_\theta(h) dh + \mathbb{E}[\tilde{\Lambda}_n | Z_0]^2 + O(n^{-1} \log^2(n)).$$

*Proof.* We first prove (33). We have  $\mathbb{E}[\tilde{\Lambda}_n | \Delta_n] = \sum_{k=1}^n \mathbb{E}[\zeta_\delta^*]_{\delta=\Delta_{k,n}}$ . We deduce from (25) that (33) holds with:

$$W_n = \frac{\Delta_{0,n} + \Delta_{n+1,n}}{\beta}(\gamma - 1) + \frac{1}{\beta} \sum_{k=1}^n \Delta_{k,n} g_1(2\theta\Delta_{k,n}).$$

Since, conditionally on  $Z_0$ , the random variables  $\Delta_{k,n}$  are all distributed as  $Z_0 \tilde{U}_n$ , where  $\tilde{U}_n$  is independent of  $Z_0$  and has distribution  $\beta(1, n+1)$ , we deduce using (32) that:

$$\mathbb{E}[|W_n| | Z_0] \leq 2 \frac{(1-\gamma)Z_0}{\beta} \mathbb{E}[\tilde{U}_n] + n \frac{2\theta Z_0^2}{\beta} \mathbb{E}[\tilde{U}_n^2 (|\log(2\theta Z_0 \tilde{U}_n)| + 2) | Z_0] = O(n^{-1} \log(n)).$$

We then prove (34). Taking the expectation in (33) conditionally on  $Z_0$ , we get:

$$\mathbb{E}[\tilde{\Lambda}_n | Z_0] = \frac{Z_0}{\beta}(1-\gamma) - n \frac{Z_0}{\beta} \mathcal{H}(2\theta Z_0) + \mathbb{E}[W_n | Z_0],$$

where

$$(36) \quad \mathcal{H}(a) = \mathbb{E}[\tilde{U}_n \log(a\tilde{U}_n)].$$

We deduce from (31) that:

$$(37) \quad n\mathcal{H}(a) = \log(a) - \log(n) + 1 - \gamma + O(n^{-1} \log(n)).$$

This gives:

$$\mathbb{E}[\tilde{\Lambda}_n | Z_0] = \frac{Z_0}{\beta} \log \left( \frac{n}{2\theta Z_0} \right) + O(n^{-1} \log(n)).$$

We finally prove (35). We have:

$$(38) \quad \mathbb{E} \left[ \tilde{\Lambda}_n^2 | \Delta_n \right] = \sum_{k=1}^n \mathbb{E} \left[ (\zeta_\delta^*)^2 \right]_{|\delta=\Delta_{k,n}} - \sum_{k=1}^n \mathbb{E} \left[ \zeta_\delta^* \right]_{|\delta=\Delta_{k,n}}^2 + \mathbb{E} \left[ \tilde{\Lambda}_n | \Delta_n \right]^2.$$

We have thanks to (26):

$$\sum_{k=1}^n \mathbb{E} \left[ (\zeta_\delta^*)^2 \right]_{|\delta=\Delta_{k,n}} = 2Z_0 \int_0^\infty hc_\theta(h) dh + W_{1,n},$$

with

$$W_{1,n} = -2(\Delta_{0,n} + \Delta_{n+1,n}) \int_0^\infty hc_\theta(h) dh + \sum_{k=1}^n \frac{\Delta_{k,n}}{\beta^2 \theta} g_2(2\theta \Delta_{k,n}).$$

Using similar computations as the ones used to bound  $\mathbb{E}[|W_n| | Z_0]$ , we get  $\mathbb{E}[|W_{1,n}| | Z_0] = O(n^{-1} \log(n))$  so that

$$\mathbb{E} \left[ \sum_{k=1}^n \mathbb{E} \left[ (\zeta_\delta^*)^2 \right]_{|\delta=\Delta_{k,n}} \mid Z_0 \right] = 2Z_0 \int_0^\infty hc_\theta(h) dh + O(n^{-1} \log(n)).$$

Thanks to (25), we have  $\mathbb{E}[\zeta_\delta^*]^2 \leq c\delta^2(|\log(\delta)| + 1)^2(1 + \delta)^2$  for some finite constant  $c$  which does not depend on  $\delta$ . We set  $\mathcal{H}_2(a) = \mathbb{E} \left[ \tilde{U}_n^2 \log^2(a\tilde{U}_n)(1 + \tilde{U}_n)^2 \right]$ , and using (32), we get:

$$(39) \quad \mathcal{H}_2(a) = O(n^{-3} \log^2(n)) = O(n^{-2} \log^2(n)).$$

We deduce that:

$$\mathbb{E} \left[ \sum_{k=1}^n \mathbb{E} \left[ \zeta_\delta^* \right]_{|\delta=\Delta_{k,n}}^2 \mid Z_0 \right] = O(n^{-1} \log^2(n)).$$

Then using (34), elementary computations give:

$$\mathbb{E} \left[ \mathbb{E} \left[ \tilde{\Lambda}_n | \Delta_n \right]^2 \mid Z_0 \right] = 2 \frac{Z_0}{\beta} (1 - \gamma) \mathbb{E}[\tilde{\Lambda}_n | Z_0] - \frac{Z_0^2}{\beta^2} (1 - \gamma)^2 + \frac{1}{\beta^2} J_{1,n} + J_{2,n} - \frac{2}{\beta} J_{3,n},$$

with  $J_{2,n} = \mathbb{E}[W_n^2 | Z_0]$ ,

$$J_{1,n} = \mathbb{E} \left[ \left( \sum_{k=1}^n \Delta_{k,n} \log(2\theta \Delta_{k,n}) \right)^2 \mid Z_0 \right] \quad \text{and} \quad J_{3,n} = \mathbb{E} \left[ W_n \left( \sum_{k=1}^n \Delta_{k,n} \log(2\theta \Delta_{k,n}) \right) \mid Z_0 \right].$$

By Cauchy-Schwartz, we have  $|J_{3,n}| \leq \sqrt{J_{1,n} J_{2,n}}$ . Using  $(\sum_{k=1}^n a_k)^2 \leq n \sum_{k=1}^n a_k^2$ , we also get:

$$J_{2,n} \leq \frac{8}{\beta^2} (\gamma - 1)^2 Z_0^2 \mathbb{E}[\tilde{U}_n^2] + \frac{2n}{\beta^2} Z_0^2 \mathbb{E} \left[ \tilde{U}_n^2 g_1^2(2\theta Z_0 \tilde{U}_n) \right] = O(n^{-2}).$$

By independence, we obtain:

$$J_{1,n} = n(n-1) \mathbb{E}[\Delta_{1,n} \log(2\theta \Delta_{1,n}) | Z_0]^2 + n \mathbb{E}[\Delta_{1,n}^2 \log^2(2\theta \Delta_{1,n}) | Z_0].$$

Recall the function  $\mathcal{H}$  defined in (36) and its asymptotic expansion (37). We have, using (39), that:

$$J_{1,n} = n(n-1) Z_0^2 \mathcal{H}(2\theta Z_0)^2 + n Z_0^2 \mathcal{H}_2(2Z_0) = Z_0^2 \left( -\log \left( \frac{n}{2\theta Z_0} \right) + 1 - \gamma \right)^2 + O(n^{-1} \log^2(n)).$$

So we deduce that:

$$\begin{aligned} \frac{1}{\beta^2} J_{1,n} + J_{2,n} - \frac{2}{\beta} J_{3,n} &= \left( -\frac{Z_0}{\beta} \log \left( \frac{n}{2\theta Z_0} \right) + \frac{Z_0}{\beta} (1 - \gamma) \right)^2 + O(n^{-1} \log^2(n)) \\ &= \left( -\mathbb{E}[\tilde{\Lambda}_n | Z_0] + \frac{Z_0}{\beta} (1 - \gamma) \right)^2 + O(n^{-1} \log^2(n)). \end{aligned}$$

We deduce that:

$$\mathbb{E} \left[ \mathbb{E} \left[ \tilde{\Lambda}_n | \Delta_n \right]^2 \mid Z_0 \right] = \mathbb{E}[\tilde{\Lambda}_n | Z_0]^2 + O(n^{-1} \log^2(n)).$$

So in the end, using (38), we get:

$$\mathbb{E} \left[ \tilde{\Lambda}_n^2 \mid Z_0 \right] = 2Z_0 \int_0^\infty hc_\theta(h) dh + \mathbb{E}[\tilde{\Lambda}_n | Z_0]^2 + O(n^{-1} \log^2(n)).$$

□

**5.3. Proof of Theorem 5.1.** We shall keep notations from Section 5.1. We set  $J_n(\varepsilon) = \mathbb{E} \left[ \left( \tilde{\Lambda}_n - \tilde{L}_\varepsilon \right)^2 \mid Z_0 \right]$ . We have:

$$J_n(\varepsilon) = \mathbb{E}[\tilde{\Lambda}_n^2 | Z_0] + \mathbb{E}[\tilde{L}_\varepsilon^2 | Z_0] - 2\mathbb{E}[\tilde{\Lambda}_n \tilde{L}_\varepsilon | Z_0].$$

By conditioning with respect to  $\Delta_n$ , and using the independence, we get:

$$\mathbb{E}[\tilde{\Lambda}_n \tilde{L}_\varepsilon | Z_0] = \mathbb{E} \left[ \mathbb{E}[\tilde{\Lambda}_n \tilde{L}_\varepsilon | \Delta_n] \mid Z_0 \right] = \Sigma_n + \mathbb{E} \left[ \mathbb{E}[\tilde{\Lambda}_n | \Delta_n] \mathbb{E}[\tilde{L}_\varepsilon | \Delta_n] \mid Z_0 \right] = \Sigma_n + \mathbb{E}[\tilde{\Lambda}_n | Z_0] \mathbb{E}[\tilde{L}_\varepsilon | Z_0],$$

where we used that  $\mathbb{E}[\tilde{L}_\varepsilon | \Delta_n] = \mathbb{E}[\tilde{L}_\varepsilon | Z_0]$  for the last equality, and:

$$\Sigma_n = \mathbb{E} \left[ \sum_{k=1}^n \mathbb{E} \left[ \zeta_{k,n}^* \sum_{z_i \in I_{k,n}} (\zeta_i - \varepsilon)_+ \mid \Delta_n \right] - \sum_{k=1}^n \mathbb{E}[\zeta_{k,n}^* | \Delta_n] \mathbb{E} \left[ \sum_{z_i \in I_{k,n}} (\zeta_i - \varepsilon)_+ \mid \Delta_n \right] \mid Z_0 \right].$$

So using (22) and (35), we get:

$$J_n(\varepsilon) = 4Z_0 \int_0^\infty hc_\theta(h) dh - 2\Sigma_n + \left( \mathbb{E}[\tilde{\Lambda}_n | Z_0] - \mathbb{E}[\tilde{L}_\varepsilon | Z_0] \right)^2 + O(\varepsilon \log(\varepsilon)) + O(n^{-1} \log^2(n)).$$

Then taking  $\varepsilon \asymp n^{-1}$ , we get, using (21), (34) and Lemma 5.5 below:

$$J_n(\varepsilon) = \frac{Z_0^2}{\beta^2} \log^2 \left( n\varepsilon \frac{\beta}{Z_0} \right) + O(n^{-1} \log^2(n)).$$

We deduce that  $\tilde{\Lambda}_n - \tilde{L}_{Z_0/(n\beta)}$  converges in probability to 0 and, by Borel-Cantelli lemma almost surely along the sub-sequence  $n^3$ . Recall that the sequence  $(\tilde{L}_\varepsilon - \mathbb{E}[\tilde{L}_\varepsilon | Z_0], \varepsilon > 0)$  converges a.s., as  $\varepsilon$  goes down to 0, towards a limit say  $\tilde{\mathcal{L}}$ . Notice that  $\mathbb{E}[\tilde{L}_{Z_0/n\beta} | Z_0] = \mathbb{E}[\tilde{\Lambda}_n | Z_0] + O(n^{-1} \log(n))$  and thus, we deduce that  $(\tilde{\Lambda}_{n^3} - \mathbb{E}[\tilde{\Lambda}_{n^3} | Z_0], n \in \mathbb{N}^*)$  converges also a.s. towards  $\tilde{\mathcal{L}}$ . Then use (33) to get that for  $k \in [n^3, (n+1)^3)$ :

$$\tilde{\Lambda}_{n^3} - \mathbb{E}[\tilde{\Lambda}_{n^3} | Z_0] + O(n^{-1} \log(n)) \leq \tilde{\Lambda}_k - \mathbb{E}[\tilde{\Lambda}_k | Z_0] \leq \tilde{\Lambda}_{(n+1)^3} - \mathbb{E}[\tilde{\Lambda}_{(n+1)^3} | Z_0] + O(n^{-1} \log(n)).$$

Then conclude that  $(\tilde{\Lambda}_n - \mathbb{E}[\tilde{\Lambda}_n | Z_0], n \in \mathbb{N}^*)$  converges also a.s. towards  $\tilde{\mathcal{L}}$ .

**Lemma 5.5.** *Let  $\varepsilon \asymp n^{-1}$ . We have:*

$$\Sigma_n = 2Z_0 \int_0^\infty hc_\theta(h) dh + O(n^{-1} \log^2(n)).$$

*Proof.* We have  $\mathbb{E} \left[ \sum_{z_i \in I_{k,n}} (\zeta_i - \varepsilon)_+ \mid \Delta_n \right] = \Delta_{k,n} \mathbb{N}[(\zeta - \varepsilon)_+]$ . Thanks to (25), (31) and (32), we get:

$$\begin{aligned} \mathbb{E} \left[ \sum_{k=1}^n \Delta_{k,n} \mathbb{E}[\zeta_{k,n}^* \mid \Delta_n] \mid Z_0 \right] &= \frac{nZ_0^2}{\beta} \mathbb{E} \left[ \tilde{U}_n^2 \left( \log(2\theta Z_0 \tilde{U}_n) + (1 - \gamma) + g_1(2\theta Z_0 \tilde{U}_n) \right) \mid Z_0 \right] \\ &= O(n^{-2} \log(n)). \end{aligned}$$

We deduce from (19) with  $\varepsilon \asymp n^{-1}$  that:

$$\mathbb{E} \left[ \sum_{k=1}^n \mathbb{E}[\zeta_{k,n}^* \mid \Delta_n] \mathbb{E} \left[ \sum_{z_i \in I_{k,n}} (\zeta_i - \varepsilon)_+ \mid \Delta_n \right] \mid Z_0 \right] = O(n^{-1} \log^2(n)).$$

According to (27), we have:

$$\sum_{k=1}^n \mathbb{E} \left[ \zeta_{k,n}^* \sum_{z_i \in I_{k,n}} (\zeta_i - \varepsilon)_+ \mid \Delta_n \right] = 2Z_0 \int_0^\infty hc_\theta(h) dh + W_n''',$$

with

$$W_n''' = -2(\Delta_{0,n} + \Delta_{n+1,n}) \int_0^\infty hc_\theta(h) dh + \sum_{k=1}^n g_3(\Delta_{k,n}).$$

Since  $\varepsilon \asymp n^{-1}$ , we deduce that

$$\mathbb{E}[|W_n'''| \mid Z_0] \leq \frac{2Z_0}{n+1} \int_0^\infty hc_\theta(h) dh + O(n^{-1} \log^2(n)).$$

This gives the result.  $\square$

## REFERENCES

- [1] R. Abraham and J.-F. Delmas. Williams' decomposition of the Lévy continuum random tree and simultaneous extinction probability for populations with neutral mutations. *Stochastic Process. Appl.*, 119(4):1124–1143, 2009.
- [2] R. Abraham and J.-F. Delmas. Reversal property of the Brownian tree. arXiv:1710.03460, 2017.
- [3] R. Abraham, J.-F. Delmas, and P. Hoscheit. A note on the Gromov-Hausdorff-Prokhorov distance between (locally) compact metric measure spaces. *Electron. J. Probab.*, 18:no. 14, 21, 2013.
- [4] D. Aldous. The continuum random tree. I. *Ann. Probab.*, 19(1):1–28, 1991.
- [5] D. Aldous and L. Popovic. A critical branching process model for biodiversity. *Adv. in Appl. Probab.*, 37(4):1094–1115, 12 2005.
- [6] H. Bi and J.-F. Delmas. Total length of the genealogical tree for quadratic stationary continuous-state branching processes. *Ann. Inst. Henri Poincaré Probab. Stat.*, 52(3):1321–1350, 2016.
- [7] A. N. Borodin and P. Salminen. *Handbook of Brownian motion—facts and formulae*. Probability and its Applications. Birkhäuser Verlag, Basel, second edition, 2002.
- [8] D. Burago, Y. D. Burago, and S. Ivanov. *A course in metric geometry*. American Mathematical Society, Providence, 2001.
- [9] Y.-T. Chen and J.-F. Delmas. Smaller population size at the MRCA time for stationary branching processes. *The Annals of Probability*, 40(5):2034–2068, 2012.
- [10] A. Dress, V. Moulton, and W. Terhalle. *T-theory: an overview*. *European J. Combin.*, 17(2-3):161–175, 1996. Discrete metric spaces (Bielefeld, 1994).
- [11] T. Duquesne. The coding of compact real trees by real valued functions. arXiv:math/0604106, 2006.
- [12] T. Duquesne and J.-F. Le Gall. Random trees, Lévy processes and spatial branching processes. *Astérisque*, (281):vi+147, 2002.
- [13] T. Duquesne and J.-F. Le Gall. Probabilistic and fractal aspects of Lévy trees. *Probab. Theory Related Fields*, 131(4):553–603, 2005.

- [14] S. N. Evans. *Probability and real trees*, volume 1920 of *Lecture Notes in Mathematics*. Springer, Berlin, 2008. Lectures from the 35th Summer School on Probability Theory held in Saint-Flour, July 6–23, 2005.
- [15] S. N. Evans, J. Pitman, and A. Winter. Rayleigh processes, real trees, and root growth with re-grafting. *Probab. Theory Related Fields*, 134(1):81–126, 2006.
- [16] M. Gromov. *Metric structures for Riemannian and non-Riemannian spaces*, volume 152 of *Progress in Mathematics*. Birkhäuser Boston, Inc., Boston, MA, 1999. Based on the 1981 French original [ MR0682063 (85e:53051)], With appendices by M. Katz, P. Pansu and S. Semmes, Translated from the French by Sean Michael Bates.
- [17] A. Lambert. Coalescence times for the branching process. *Advances in Applied Probability*, 35(4):1071–1089, 2003.
- [18] A. Lambert. Quasi-stationary distributions and the continuous-state branching process conditioned to be never extinct. *Electron. J. Probab.*, 12:no. 14, 420–446, 2007.
- [19] J.-F. Le Gall and Y. Le Jan. Branching processes in Lévy processes: the exploration process. *Ann. Probab.*, 26(1):213–252, 1998.
- [20] P. Pfaffelhuber, A. Wakolbinger, and H. Weisshaupt. The tree length of an evolving coalescent. *Probab. Theory Related Fields*, 151(3-4):529–557, 2011.
- [21] L. Popovic. Asymptotic genealogy of a critical branching process. *Ann. Appl. Probab.*, 14(4):2120–2148, 2004.
- [22] S. Roelly-Coppoletta and A. Rouault. Processus de Dawson-Watanabe conditionné par le futur lointain. *C. R. Acad. Sci. Paris Sér. I Math.*, 309(14):867–872, 1989.

ROMAIN ABRAHAM, INSTITUT DENIS POISSON, UNIVERSITÉ D'ORLÉANS, UNIVERSITÉ DE TOURS, CNRS, FRANCE

*E-mail address:* `romain.abraham@univ-orleans.fr`

JEAN-FRANÇOIS DELMAS, UNIVERSITÉ PARIS-EST, CERMICS (ENPC), FRANCE

*E-mail address:* `delmas@cermics.enpc.fr`