



HAL
open science

The split-and-drift random graph, a null model for speciation

François Bienvenu, Florence Débarre, Amaury Lambert

► **To cite this version:**

François Bienvenu, Florence Débarre, Amaury Lambert. The split-and-drift random graph, a null model for speciation. 2017. hal-01533404v1

HAL Id: hal-01533404

<https://hal.science/hal-01533404v1>

Preprint submitted on 6 Jun 2017 (v1), last revised 16 Feb 2024 (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.

The split-and-drift random graph, a null model for speciation

François Bienvenu^{*1, 2}, Florence Débarre¹, and Amaury Lambert^{1, 2}

¹*Center for Interdisciplinary Research in Biology (CIRB), Collège de France, CNRS
UMR 7241, Paris, France*

²*Laboratoire de Probabilités et Modèles Aléatoires (LPMA), CNRS UMR 7599, UPMC
Univ Paris 06, Paris, France*

June 6, 2017

Abstract

We introduce a new random graph model motivated by biological questions relating to speciation. This random graph is defined as the stationary distribution of a Markov chain on the space of graphs on $\{1, \dots, n\}$. The dynamics of this Markov chain is governed by two types of events: vertex duplication, where at constant rate a pair of vertices is sampled uniformly, one of these vertices loses its incident edges and is rewired to the other vertex and its neighbors; and edge removal, where each edge disappears at constant rate. Besides the number of vertices n , the model has a single parameter r_n .

Using a coalescent approach, we obtain explicit formulas for the first moments of several graph invariants such as the number of edges or the number of complete subgraphs of order k . These are then used to identify five non-trivial regimes depending on the asymptotics of the parameter r_n . We derive an explicit expression for the degree distribution, and show that under appropriate rescaling it converges to classical distributions when the number of vertices goes to infinity. Finally, we give asymptotic bounds for the number of connected components, and show that in the sparse regime the number of edges is Poissonian.

*Corresponding author: francois.bienvenu@normalesup.org

Contents

1	Introduction	3
1.1	Biological context	3
1.2	Formal description of the model	4
1.3	Notation	5
1.4	Statement of results	6
2	Coalescent constructions of G_{n,r_n}	8
2.1	The standard Moran process	8
2.2	Backward construction	9
2.3	Forward construction	12
3	First moments methods	13
3.1	First moments of graph invariants	13
3.1.1	Degree and number of edges	13
3.1.2	Complete subgraphs	15
3.2	Identification of different regimes	18
4	The degree distribution	20
4.1	Ideas of the proof of Theorem 4.1	20
4.2	Formal proof of Theorem 4.1	22
4.2.1	The vertex-marking process	22
4.2.2	The branching process	24
4.2.3	Relabeling and end of proof	25
5	Connected components in the intermediate regime	25
5.1	Lower bound on the number of connected components	26
5.2	Upper bound on the number of connected components	27
6	Number of edges in the sparse regime	28
6.1	Proof of the positive relation between the edges	29
6.1.1	Preliminary lemmas	29
6.1.2	Stein–Chen coupling	30
6.2	Proof of Theorem 6.1	33
A	Proofs of Proposition 3.5 and Corollary 3.6	35
A.1	Proof of Proposition 3.5	35
A.2	Proof of Corollary 3.6	37
B	Proof of Theorem 4.2	38
B.1	Outline of the proof	38
B.2	Step 1	39
B.3	Step 2	39
B.3.1	Proof of (i)	40
B.3.2	Proof of (ii)	41

1 Introduction

In this paper, we introduce a random graph derived from a minimalistic model of speciation. This random graph bears superficial resemblance to classic models of protein interaction networks [3, 6, 13, 16] in that the events shaping the graph are the duplication of vertices and the loss of edges. However, our model is obtained as the steady state of a Markov process (rather than by repeatedly adding vertices), and has the crucial feature that the duplication of vertices is independent from the loss of edges. These differences result in a very different behavior of the model.

Before describing the model formally in Section 1.2, let us briefly explain the motivation behind its introduction.

1.1 Biological context

Although it is often presented as central to biology, there is no consensus about how the concept of species should be defined. A widely held view is that it should be based on the capacity of individuals to interbreed. This is the so-called “biological species concept”, wherein a species is defined as a group of potentially interbreeding populations that cannot interbreed with populations outside the group.

This view, whose origins can be traced back to the beginning of the 20th century [11], was most famously promoted by Ernst Mayr [9] and has been most influential in biology [4]. However, it remains quite imprecise: indeed, groups of populations such that (1) all pairs of populations can interbreed and (2) no population can interbreed with a population outside the group are probably not common in nature – and, at any rate, do not correspond to what is considered a species in practice. Therefore, some leniency is required when applying conditions (1) and (2). But once we allow for this, there are several ways to formalize the biological species concept, as illustrated in Figure 1. Thus, it seems arbitrary to favor one over the others in the absence of a mechanism to explain why some kind of groups should be more relevant (e.g., arise more frequently) than others.

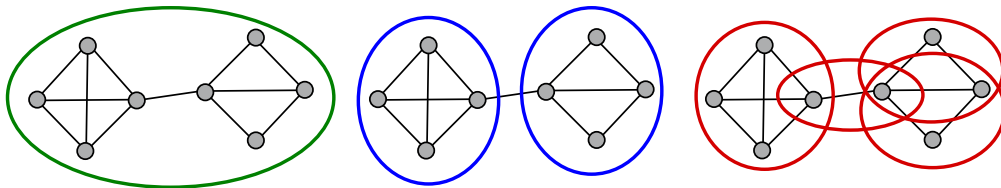


Figure 1: The vertices of the graph represent populations and its edges denote interbreeding potential (that is, individuals from two linked population could interbreed, if given the chance). Even with such perfect information, it is not obvious how to delineate “groups of potentially interbreeding populations that cannot interbreed with populations outside the group”: should these correspond to connected components (on the left, in green), maximal complete subgraphs (on the right, in red), or be based on some other clustering method (middle, in blue)?

Our aim is to build a minimal model of speciation that would make predictions about the structure and dynamics of the interbreeding network and allow one to recover species as an emergent property. To do so, we model speciation at the level of populations. Thus, we consider a set of n populations and we track the interbreeding ability for every pair of populations. All this information is encoded in a graph whose vertices correspond to populations and whose edges indicate potential interbreeding, i.e., two vertices are linked if and only if the corresponding populations can interbreed.

Speciation will result from the interplay between two mechanisms. First, populations can sometimes “split” into two initially identical populations which then behave as independent entities; this could happen as a result of the fragmentation of the habitat or of the colonization of a new patch. Second, because they behave as independent units, two initially identical populations will diverge (e.g., as a result of genetic drift) until they can no longer interbreed.

1.2 Formal description of the model

Start from any graph with vertex set $V = \{1, \dots, n\}$, and let it evolve according to the following rules

1. **Vertex duplication:** each vertex “duplicates” at rate 1; when a vertex duplicates, it chooses another vertex uniformly at random among the other vertices and replaces it with a copy of itself. The replacement of j by a copy of i means that j loses its incident edges and is then linked to i and to all of its neighbors, as depicted in Figure 2.

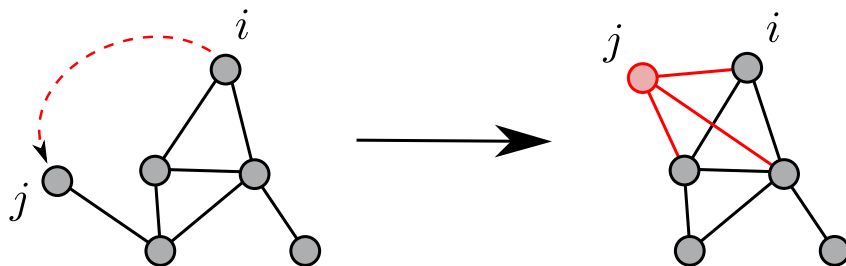


Figure 2: An illustration of vertex duplication. Here, i duplicates and replaces j . After the duplication, j is linked to i and to each of its neighbors.

2. **Edge removal:** each edge disappears at constant rate ρ .

This procedure defines a continuous-time Markov chain $(G_n(t))_{t \geq 0}$ on the finite state space of all graphs whose vertices are the integers $1, \dots, n$. It is easy to see that this Markov chain is irreducible. Indeed, to go from any graph G_1 to any graph G_2 , one can consider the following sequence of events: first, a vertex is duplicated repeatedly in order to obtain the complete graph of order n (e.g., $\forall k \in \{2, \dots, n\}$, vertex k is replaced by a copy of vertex 1); then, all the edges that are not in G_2 are removed.

Because the Markov chain $(G_n(t))_{t \geq 0}$ is irreducible, it has a unique stationary probability distribution $\mu_{n,\rho}$. This probability distribution on the set of graphs of order n defines a random graph that is the object of study of this paper.

1.3 Notation

To study the asymptotic behavior of our model as $n \rightarrow +\infty$, we can let ρ , the ratio of the edge removal rate to the vertex duplication rate, be a function of n . As will become evident, it is more convenient to parametrize the model by

$$r_n := \frac{n-1}{2} \rho_n.$$

Thus, we write G_{n,r_n} to refer to a random graph whose law is $\mu_{n, \frac{2r_n}{n-1}}$.

Although some of our results hold for any (n, r) , in many cases we will be interested in asymptotic properties that are going to depend on the asymptotics of r_n . To quantify these, we will use the Bachmann–Landau notation, which for positive sequences r_n and $f(n)$ can be summarized as:

- $r_n \sim f(n)$ when $r_n/f(n) \rightarrow 1$.
- $r_n = o(f(n))$ when $r_n/f(n) \rightarrow 0$.
- $r_n = \Theta(f(n))$ when there exists positive constants α and β such that, asymptotically, $\beta f(n) \leq r_n \leq \beta f(n)$.
- $r_n = \omega(f(n))$ when $r_n/f(n) \rightarrow +\infty$.

These notations also have stochastic counterparts, whose meaning will be recalled when we use them.

Finally, we occasionally use the expression *asymptotically almost surely* (abbreviated as a.a.s.) to that a property holds with probability that goes to 1 as n tends to infinity:

$$\mathcal{Q}_n \text{ a.a.s.} \iff \mathbb{P}(\mathcal{Q}_n) \xrightarrow[n \rightarrow +\infty]{} 1.$$

1.4 Statement of results

Table 1 lists the first moments of several graph invariants obtained in Section 3.1. These are then used to identify different regimes, depending on the asymptotic behavior of the parameter r_n , as stated in Theorem 3.10.

Variable	Expectation	Variance	Covariance
$\mathbb{1}_{\{i \leftrightarrow j\}}$	$\frac{1}{1+r}$	$\frac{r}{(1+r)^2}$	$\frac{r}{(1+r)^2(3+2r)}$ if vertex in common, $\frac{2r}{(1+r)^2(3+r)(3+2r)}$ otherwise.
$D_n^{(i)}$	$\frac{n-1}{1+r}$	$\frac{r(n-1)(1+2r+n)}{(1+r)^2(3+2r)}$	$\frac{r}{(1+r)^2} \left(1 + \frac{3(n-2)}{3+2r} + \frac{2(n-2)(n-3)}{(3+r)(3+2r)} \right)$
$ E_n $	$\frac{n(n-1)}{2(1+r)}$	$\frac{rn(n-1)(n^2+2r^2+2nr+n+5r+3)}{2(1+r)^2(3+r)(3+2r)}$	—
$X_{n,k}$	$\binom{n}{k} \left(\frac{1}{1+r} \right)^{k-1}$?	—

Table 1: First and second moments of several graph invariants of $G_{n,r}$: $\mathbb{1}_{\{i \leftrightarrow j\}}$ is the variable indicating that $\{ij\}$ is an edge, $D_n^{(i)}$ the degree of vertex i , $|E_n|$ the number of edges and $X_{n,k}$ the number of complete subgraphs of order k . These expressions hold for every n and r . The variance of the number of complete subgraphs of order k is not known. Finally, the covariance of the indicator variables of two edges depends on whether these edges share a common end, hence the two expressions.

Theorem 3.10. *Let D_n be the degree of a fixed vertex of G_{n,r_n} . In the limit as $n \rightarrow +\infty$, depending on the asymptotics of r_n we have the following behaviors for G_{n,r_n}*

- (i) Transition for the complete graph: when $r_n = o(1/n)$, $\mathbb{P}(G_{n,r_n} \text{ is complete})$ goes to 1, while when $r_n = \omega(1/n)$ it goes to 0; when $r_n = \Theta(1/n)$, this probability is bounded away from 0 and from 1.
- (ii) Dense regime: when $r_n = o(1)$, $\mathbb{P}(D_n = n - 1) \rightarrow 1$.
- (iii) Sparse regime: when $r_n = \omega(n)$, $\mathbb{P}(D_n = 0) \rightarrow 1$.
- (iv) Transition for the empty graph: when $r_n = o(n^2)$, $\mathbb{P}(G_{n,r_n} \text{ is empty})$ goes to 0 while when $r_n = \omega(n^2)$ it goes to 1; when $r_n = \Theta(n^2)$, this probability is bounded away from 0 and from 1.

In Section 4, we derive an explicit expression for the degree distribution, which holds for every value of n and r_n . We then show that, under appropriate rescaling, this degree converges to classical distributions.

Theorem 4.1.

Let D_n be the degree of a fixed vertex of G_{n,r_n} . Then, for each $k \in \{0, \dots, n - 1\}$,

$$\mathbb{P}(D_n = k) = \frac{2r_n(2r_n + 1)}{(n + 2r_n)(n - 1 + 2r_n)} (k + 1) \prod_{i=1}^k \frac{n - i}{n - i + 2r_n - 1}, \quad (1)$$

where the empty product is 1.

Theorem 4.2.

- (i) If $r_n \rightarrow r > 0$, then $\frac{D_n}{n}$ converges in distribution to a Beta(2, 2r) random variable.
- (ii) If r_n is both $\omega(1)$ and $o(n)$, then $\frac{D_n}{n/r_n}$ converges in distribution to a size-biased exponential variable with parameter 2.
- (iii) If $2r_n/n \rightarrow \rho > 0$, then $D_n + 1$ converges in distribution to a size-biased geometric variable with parameter $\rho/(1 + \rho)$.

Asymptotic bounds for the number of connected components are obtained in Section 5, where the following theorem is proved.

Theorem 5.1. Let $\#\text{CC}_n$ be the number of connected components of G_{n,r_n} . If r_n is both $\omega(1)$ and $o(n)$, then

$$\frac{r_n}{2} + o_p(r_n) \leq \#\text{CC}_n \leq 2r_n \log n + o_p(r_n \log n)$$

where, for a positive sequence (u_n) , $o_p(u_n)$ stands for some sequence of random variables (X_n) such that $X_n/u_n \rightarrow 0$ in probability.

Because the method used to obtain the upper bound in Theorem 5.1 is rather crude, we formulate the following conjecture, which is well supported by simulations.

Conjecture 5.4.

$$\exists \alpha, \beta > 0 \text{ s.t. } \mathbb{P}(\alpha r_n \leq \#\text{CC}_n \leq \beta r_n) \xrightarrow[n \rightarrow \infty]{} 1.$$

Finally, in Section 6 we use the Stein–Chen method to show that the number of edges is Poissonian in the sparse regime, as shown by Theorem 6.1.

Theorem 6.1. Let $|E_n|$ be the number of edges of G_{n,r_n} . If $r_n = \omega(n)$ then

$$d_{\text{TV}}(|E_n|, \text{Poisson}(\lambda_n)) \xrightarrow[n \rightarrow +\infty]{} 0,$$

where d_{TV} stands for the total variation distance and $\lambda_n = \mathbb{E}(|E_n|) \sim \frac{n^2}{2r_n}$. If in addition $r_n = o(n^2)$, then $\lambda_n \rightarrow +\infty$ and as a result

$$\frac{|E_n| - \lambda_n}{\sqrt{\lambda_n}} \xrightarrow[n \rightarrow +\infty]{\mathcal{D}} \mathcal{N}(0, 1),$$

where $\mathcal{N}(0, 1)$ denotes the standard normal distribution.

These results are summarized in Figure 3.

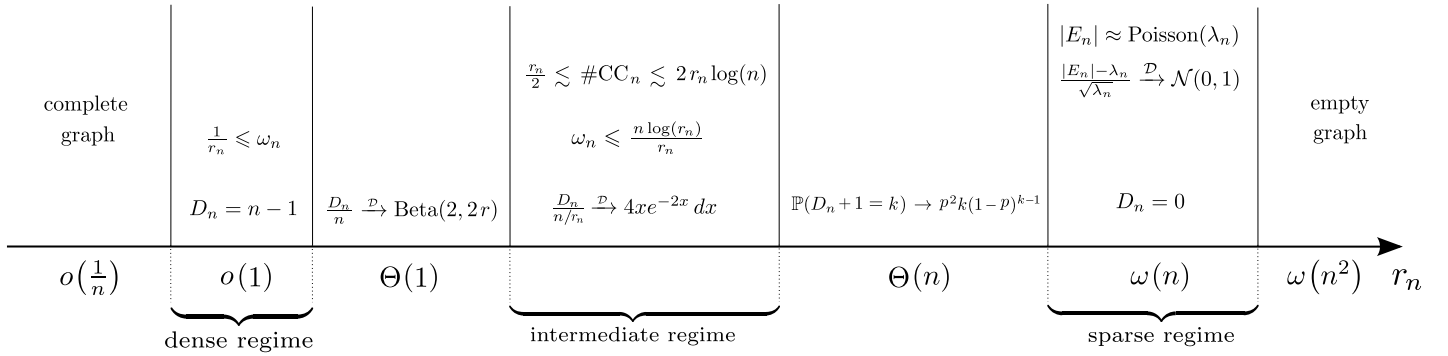


Figure 3: A graphical summary of the main results established in the paper; D_n is the degree of a fixed vertex, $|E_n|$ the number of edges, $\#CC_n$ the number of connected components, and ω_n the clique number. All equalities and inequalities are to be understood “asymptotically almost surely” (i.e. hold with probability that goes to 1 as n tends to infinity).

2 Coalescent constructions of G_{n,r_n}

In this section, we detail coalescent constructions of G_{n,r_n} that will be used throughout the rest of the paper. Let us start by recalling some results about the Moran model.

2.1 The standard Moran process

The Moran model [10] is a classic model of population genetics. It consists in a set of n particles governed by the following dynamics: after an exponential waiting time with parameter $\binom{n}{2}$, a pair of particles is sampled uniformly at random. One of these particles is then removed (death) and replaced by a copy of the other (birth), and we iterate the procedure.

In this document, we will use the Poissonian representation of the Moran process detailed in the next definition.

Definition 2.1. The *driving measure of a standard Moran process on V* is a collection $\mathcal{M} = (M_{(ij)})_{(ij) \in V^2}$ of i.i.d. Poisson point processes with rate $1/2$ on \mathbb{R} .

We think of the elements of V as sites, each occupied by a single particle. In forward time, each atom $t \in M_{(ij)}$ indicates the replacement, at time t , of the particle in i by a copy of the particle in j .

For any given time $\alpha \in \mathbb{R}$, \mathcal{M} defines a *genealogy of V on $]-\infty, \alpha]$* . Taking $\alpha = 0$ and working in backward time, i.e. writing $t \geq 0$ to refer to the absolute time $-t$, this genealogy is described by a collection of *ancestor functions* a_t , $t \in [0, +\infty[$, $a_t: V \rightarrow V$, defined as follows: $(a_t)_{t \geq 0}$ is the piecewise constant process such that

- (i) a_0 is the identity on V .
- (ii) If $t \in M_{(ij)}$ then $a_t(i) = a_{t-}(j)$ and, for all $k \neq i$, $a_t(k) = a_{t-}(k)$.
- (iii) If for all $(ij) \in V^2$, $M_{(ij)} \cap [s, t] = \emptyset$, then $a_t = a_s$.

We refer to $a_t(i)$ as the ancestor of i at time t before the present – or, more simply, as the *ancestor of i at time t* .

The standard Moran process is closely related to the Kingman coalescent [7]. Indeed, let \mathcal{R}_t denote the equivalence relation on V defined by

$$i \mathcal{R}_t j \iff a_t(i) = a_t(j),$$

and let $K_t = V/\mathcal{R}_t$ be the partition of V induced by \mathcal{R}_t . Then, $(K_t)_{t \geq 0}$ is a Kingman coalescent on V . In particular, we will frequently use the next lemma.

Lemma 2.2. *Let $(a_t)_{t \geq 0}$ be the ancestor functions of a standard Moran process on V . For any $i \neq j$, let*

$$T_{\{ij\}} = \inf\{t \geq 0 : a_t(i) = a_t(j)\}$$

be the coalescence time of i and j , and for any $S \subset V$ let

$$T_S = \inf\{T_{\{ij\}} : i, j \in S, i \neq j\}$$

Then, for all $t \geq 0$, conditional on $\{T_S > t\}$, $(T_S - t)$ is an exponential variable with parameter $\binom{|S|}{2}$.

2.2 Backward construction

We now turn to the description of the coalescent framework on which our study relies.

First, consider the two-sided extension of $(G_n(t))_{t \geq 0}$, i.e. the corresponding stationary process on \mathbb{R} (see, e.g., Section 7.1 of [5]), which by a slight abuse of notation we note $(G_n(t))_{t \in \mathbb{R}}$.

To study the stationary distribution of $(G_n(t))_{t \in \mathbb{R}}$, it is equivalent to consider the time-rescaled process $(G'(t))_{t \in \mathbb{R}}$ defined by

$$G'(t) = G_n(t(n-1)/2).$$

In this process, each vertex duplicates at rate $(n-1)/2$ and each edge disappears at rate $r_n = (n-1)\rho_n/2$. All these events being independent, we see that the vertex duplications correspond to the atoms of a standard Moran process on $V = \{1, \dots, n\}$, and the edge removals to the atoms of $\binom{n}{2}$ i.i.d. Poisson point processes with rate r_n on \mathbb{R} that are independent of the Moran process. More specifically, writing $\mathcal{M} = (M_{(ij)})_{(ij) \in V^2}$ for the vertex duplication process and $\mathcal{P} = (P_{\{ij\}})_{\{ij\} \in V^{(2)}}$, where $V^{(2)}$ denotes the set of unordered pairs of elements of V , for the edge removal process,

- If $t \in M_{(ij)}$, then j duplicates and replaces i at time t .
- If $t \in P_{\{ij\}}$, then if there is an edge between i and j at time t , it is removed.

Now let us re-express the event “ i is linked to j in $G'(0)$ ” in terms of \mathcal{M} and \mathcal{P} . Using the notation of Definition 2.1 for the ancestor functions (i.e. writing $a_t(i)$ for the ancestor of vertex i at time t before the present), set

$$\tau = \sup\{t \geq 0 : a_t(i) = i, a_t(j) = j\}.$$

If $\tau < T_{\{ij\}}$, then

$$\{i \leftrightarrow j \text{ in } G'(0)\} = \{P_{\{ij\}} \cap]-\tau, 0] = \emptyset\} \cap \{a_\tau(i) \leftrightarrow a_\tau(j) \text{ in } G'(-\tau)\},$$

while if $\tau = T_{\{ij\}}$ then

$$\{i \leftrightarrow j \text{ in } G'(0)\} = \{P_{\{ij\}} \cap]-\tau, 0] = \emptyset\},$$

because either i or j was duplicated from the other one at time $-\tau$, and a duplicated vertex and its copy are always linked right after the duplication.

Proceeding similarly, we let $\sigma = \sup\{t > \tau : a_t(i) = a_\tau(i), a_t(j) = a_\tau(j)\}$ in order to express $\{a_\tau(i) \leftrightarrow a_\tau(j) \text{ in } G'(-\tau)\}$ as

- $\{P_{\{a_\tau(i)a_\tau(j)\}} \cap]-\sigma, -\tau] = \emptyset\} \cap \{a_\sigma(i) \leftrightarrow a_\sigma(j) \text{ in } G'(-\sigma)\}$ if $\sigma < T_{\{ij\}}$,
- $\{P_{\{a_\tau(i)a_\tau(j)\}} \cap]-\sigma, -\tau] = \emptyset\}$ if $\sigma = T_{\{ij\}}$.

And, likewise, we reexpress $\{a_\sigma(i) \leftrightarrow a_\sigma(j) \text{ in } G'(-\sigma)\}$, etc. This iterative procedure will always stop before the finite time $T_{\text{MRCA}} = \sup\{T_{\{ij\}} : \{ij\} \in V^{(2)}\}$, i.e. after at most $n - 1$ steps.

Since $T_{\{ij\}}$ and the intermediate times τ, σ, \dots are determined by \mathcal{M} , this shows that $G'(0)$ is fully specified by \mathcal{M} and \mathcal{P} . Note that we do not need any initial condition, reflecting the fact that every edge of $G'(0)$ is descended from an edge that was introduced between a vertex and its copy.

We record these observations in the following proposition.

Proposition 2.3. *Let $V = \{1, \dots, n\}$ and let $V^{(2)}$ be the set of unordered pairs of elements of V . Let \mathcal{M} be the driving measure of a standard Moran process on V , and $(a_t)_{t \geq 0}$ the associated ancestor functions (that is, for each i in V , $a_t(i)$ is the ancestor of i at time t). Let $\mathcal{P} = (P_{\{ij\}})_{\{ij\} \in V^{(2)}}$ be a collection of i.i.d. Poisson point processes with rate r_n on $[0, +\infty[$ such that \mathcal{M} and \mathcal{P} are independent. For every pair $\{ij\} \in V^{(2)}$, define*

$$P_{\{ij\}}^* = \{t \geq 0 : t \in P_{\{a_t(i)a_t(j)\}}\},$$

with the convention that, $\forall k \in V$, $P_{\{k\}} = \emptyset$. Finally, let $G = (V, E)$ be the graph defined by

$$E = \{\{ij\} \in V^{(2)} : P_{\{ij\}}^* = \emptyset\}.$$

Then, $G \sim G_{n,r_n}$.

In the rest of this document, we write G_{n,r_n} for the graph obtained by the procedure of Proposition 2.3.

Finally, our study of G_{n,r_n} is based on the following result.

Lemma 2.4. *Let S be a subset of $V^{(2)}$. Conditional on the measure \mathcal{M} , for any interval I such that*

(i) *For all $\{ij\} \in S$, $\forall t \in I$, $a_t(i) \neq a_t(j)$.*

(ii) *For all $\{kl\} \neq \{ij\}$ in S , $\forall t \in I$, $\{a_t(i), a_t(j)\} \neq \{a_t(k), a_t(\ell)\}$,*

$P_{\{ij\}}^ \cap I$, $\{ij\} \in S$, are independent Poisson point processes with rate r_n on I .*

Moreover, for any disjoint intervals I and J , $(P_{\{ij\}}^ \cap I)_{\{ij\} \in S}$ is independent of $(P_{\{ij\}}^* \cap J)_{\{ij\} \in S}$.*

Proof of Lemma 2.4. For all $t \geq 0$, define S_t by

$$S_t = \{\{a_t(i), a_t(j)\} : \{ij\} \in S\}.$$

Set $t_0 = \inf I$ and let t_1, \dots, t_{m-1} be the jump times of $(S_t)_{t \geq 0}$ on I , i.e.

$$t_p = \inf\{t > t_{p-1} : S_t \neq S_{t_{p-1}}\}, \quad p = 1, \dots, m-1.$$

Finally, set $t_m = \sup I$ and, for $p = 0, \dots, m-1$, let $I_p = [t_p, t_{p+1}[$ and $\tilde{a}_p = a_{t_p}$. With this notation, for all $\{ij\} \in S$,

$$P_{\{ij\}}^* \cap I = \bigcup_{p=0}^{m-1} \left(P_{\{\tilde{a}_p(i), \tilde{a}_p(j)\}} \cap I_p \right),$$

where for $p \neq q$, $I_p \cap I_q = \emptyset$, and $P_{\{uv\}}$, $\{uv\} \in V^{(2)}$, are i.i.d. Poisson point processes on $[0, +\infty[$ with rate r_n . By assumption, for all $p = 0, \dots, m-1$, for all $\{ij\} \neq \{kl\}$ in S , $\tilde{a}_p(i) \neq \tilde{a}_p(j)$, $\tilde{a}_p(k) \neq \tilde{a}_p(\ell)$ and $\{\tilde{a}_p(i), \tilde{a}_p(j)\} \neq \{\tilde{a}_p(k), \tilde{a}_p(\ell)\}$. This shows that $(P_{\{\tilde{a}_p(i), \tilde{a}_p(j)\}} \cap I_p)$, $\{ij\} \in S$ and $p = 0, \dots, m-1$, are i.i.d. Poisson point processes with rate r_n on the corresponding intervals, proving the first part of the lemma.

The second assertion is proved similarly. Adapting the previous notation to work with two disjoint intervals I and J , i.e. letting $(\tilde{a}_p^I)_{p \in \{0, \dots, m_I\}}$ be the embedded chain of $(a_t)_{t \in I}$ and $(\tilde{a}_p^J)_{p \in \{0, \dots, m_J\}}$ that of $(a_t)_{t \in J}$, for all $\{ij\} \in S$ we write

$$P_{\{ij\}}^* \cap I = \bigcup_{p=0}^{m_I-1} \left(P_{\{\tilde{a}_p^I(i), \tilde{a}_p^I(j)\}} \cap I_p \right),$$

and

$$P_{\{ij\}}^* \cap J = \bigcup_{p=0}^{m_J-1} \left(P_{\{\tilde{a}_p^J(i), \tilde{a}_p^J(j)\}} \cap J_p \right).$$

We conclude the proof by noting that the families

$$\left(P_{\{\tilde{a}_p^I(i), \tilde{a}_p^I(j)\}} \cap I_p \right)_{\{ij\} \in S, p \in \{0, \dots, m_I\}}$$

and

$$\left(P_{\{\tilde{a}_p^J(i), \tilde{a}_p^J(j)\}} \cap J_p \right)_{\{ij\} \in S, p \in \{0, \dots, m_J\}}$$

are independent, because the elements of these families are either deterministic (if, e.g. $\tilde{a}_p^I(i) = \tilde{a}_p^I(j)$, in which case $P_{\{\tilde{a}_p^I(i), \tilde{a}_p^I(j)\}} = \emptyset$) or Poisson point processes on intervals that are disjoint from each of the intervals involved in the definition of the other family. \square

If we think of $\{a_t(i), a_t(j)\}$ as being the ancestor of $\{ij\}$ at time t , then the genealogy of vertices induces a genealogy of pairs of vertices, as illustrated by Figure 4. In words, Proposition 2.3 and Lemma 2.4 say that edge-removal events occur at constant rate r_n along the branches of this genealogy, and that the events affecting disjoint sections of branches are independent. Thus, we can think of $P_{\{ij\}}^*$, $\{ij\} \in V^{(2)}$, as a single Poisson point process P^* on the lineages of pairs of vertices; a pair of vertices corresponds to an edge in G_{n,r_n} if and only if there is no atom of P^* on its lineage.

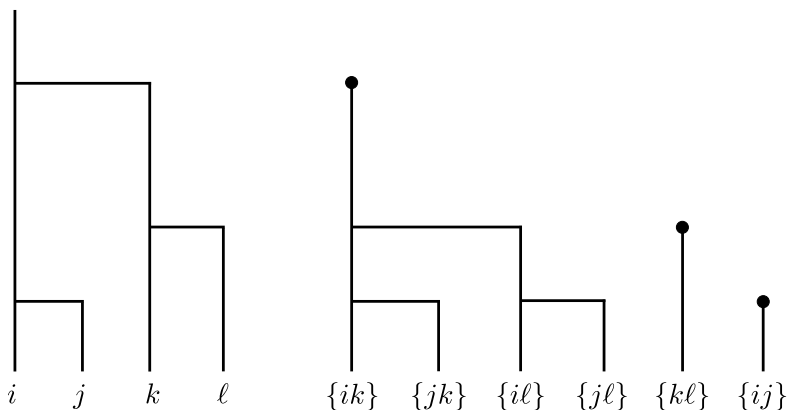


Figure 4: On the left, a genealogy on $\{i, j, k, \ell\}$ and on the right the corresponding genealogy of the pairs. Edge removal events occur at constant rate along the lineages of pairs of vertices, and a pair of vertices is an edge of G_{n,r_n} if and only if there is no atom on its lineage.

2.3 Forward construction

We now give a forward version of the coalescent construction presented in the previous section.

Consider the Markov process $(G_r^\dagger(t))_{t \geq 0}$ defined by

- (i) $G_r^\dagger(0) = (\{1, 2\}, \{\{1, 2\}\})$ is the complete graph of order 2.
- (ii) Conditional on $|V_t| = k$, where V_t is the set of vertices of $G_r^\dagger(t)$: after an exponential time with parameter $\binom{k}{2}$, a vertex is sampled uniformly in V_t and duplicated without replacement (i.e., we copy the vertex and all incident edges, resulting in a graph with $k + 1$ vertices).
- (iii) During the whole process, each edge disappears at constant rate r .

Next, for $k \in \mathbb{N}$, $k \geq 2$, let $G_r^*(k) = G_r^\dagger(t_k^-)$, where

$$t_k = \sup\{t \geq 0 : G_r^*(t) \text{ has } k \text{ vertices}\}.$$

Then, $(G_r^*(k))_{k \geq 2}$ is a sequence of graphs of increasing order with the property that

$$\forall r > 0, \forall n \geq 2, \quad G_r^*(n) \sim G_{n,r}.$$

In particular, for any given sequence (r_n) , for any $n \geq 2$ $G_{r_n}^*(n) \sim G_{n,r_n}$. Note, however, that unless (r_n) is constant, this is not a compatible construction of a sequence $(G_{n,r_n})_{n \geq 2}$, as all elements of a sequence $(G_r^*(n))_{n \geq 2}$ that results from the procedure described above are associated to the same value of r , while each term of a sequence $(G_{n,r_n})_{n \geq 2}$ corresponds to a different value of r_n .

3 First moments methods

In this section, we apply Proposition 3.7 and Lemma 2.4 to obtain the expressions presented in Table 1. These are then used to identify different regimes for G_{n,r_n} , depending on the asymptotic behavior of the parameter r_n .

In order to be able to use Lemma 2.4, we will always reason conditionally on the genealogy of the vertices (i.e. on the vertex duplication process \mathcal{M}) and then integrate against its law.

3.1 First moments of graph invariants

3.1.1 Degree and number of edges

Proposition 3.1. *For any fixed vertices i and j , $i \neq j$, the probability that i and j are linked in G_{n,r_n} is*

$$\mathbb{P}(i \leftrightarrow j) = \frac{1}{1 + r_n}.$$

Corollary 3.2. *Let D_n be the degree of a fixed vertex of G_{n,r_n} , and $|E_n|$ be the number of edges of G_{n,r_n} . Then,*

$$\mathbb{E}(D_n) = \frac{n-1}{1+r_n} \quad \text{and} \quad \mathbb{E}(|E_n|) = \binom{n}{2} \frac{1}{1+r_n}.$$

Proof. By Proposition 2.3,

$$\{i \leftrightarrow j\} \iff P_{\{ij\}}^* \cap [0, T_{\{ij\}}[= \emptyset.$$

Reasoning conditionally on $T_{\{ij\}}$ and applying Lemma 2.4 to $S = \{\{ij\}\}$ and $I = [0, T_{\{ij\}}[$, we see that $P_{\{ij\}}^*$ is a Poisson point process with rate r_n on I . Since $T_{\{ij\}} \sim \text{Exp}(1)$,

$$\mathbb{P}(i \leftrightarrow j) = \mathbb{P}(e_1 > T_{\{ij\}}),$$

where $e_1 = \inf P_{\{ij\}}^*$ is an exponential variable with rate r_n that is independent of $T_{\{ij\}}$.

The corollary follows directly from the fact that the degree of a vertex v can be written as

$$D_n^{(v)} = \sum_{i \neq v} \mathbb{1}_{\{i \leftrightarrow v\}}$$

and that the number of edges of G_{n,r_n} is

$$|E_n| = \sum_{\{ij\} \in V^{(2)}} \mathbf{1}_{\{i \leftrightarrow j\}}. \quad \square$$

Proposition 3.3. *Let i, j and k be three distinct vertices of G_{n,r_n} . We have*

$$\text{Cov}\left(\mathbf{1}_{\{i \leftrightarrow j\}}, \mathbf{1}_{\{i \leftrightarrow k\}}\right) = \frac{r_n}{(3 + 2r_n)(1 + r_n)^2}$$

Corollary 3.4. *Let D_n be the degree of a fixed vertex of G_{n,r_n} . We have*

$$\text{Var}(D_n) = \frac{r_n(n-1)(1+2r_n+n)}{(1+r_n)^2(3+2r_n)}$$

Proof. For all $t \geq 0$, let $S_t = \{a_t(i), a_t(j), a_t(k)\}$. Let $\tau_1 = \inf\{t \geq 0 : |S_t| = 2\}$ and $\tau_2 = \inf\{t \geq \tau_1 : |S_t| = 1\}$. Recall from Lemma 2.2 that τ_1 and $\tau_2 - \tau_1$ are independent exponential variables with parameter 3 and 1, respectively. Finally, let $\{u, v\} = S_{\tau_1}$.

By Proposition 2.3, $\{ij\}$ and $\{ik\}$ are edges of G_{n,r_n} if and only if $P_{\{ij\}}^* \cap [0, T_{\{ij}\}]$ and $P_{\{ik\}}^* \cap [0, T_{\{ik}\}]$ are empty, which can also be written

$$\left(P_{\{ij\}}^* \cap [0, \tau_1[\right) \cup \left(P_{\{ik\}}^* \cap [0, \tau_1[\right) \cup \left(P_{\{uv\}}^* \cap [\tau_1, \tau_2[\right) = \emptyset$$

Conditionally on τ_1 and τ_2 , by Lemma 2.4, $(P_{\{ij\}}^* \cap [0, \tau_1[\cup P_{\{ik\}}^* \cap [0, \tau_1[$ is independent of $P_{\{uv\}}^* \cap [\tau_1, \tau_2[$, $P_{\{ij\}}^*$ and $P_{\{ik\}}^*$ are independent Poisson point processes with rate r_n on $[0, \tau_1[$, and $P_{\{uv\}}^*$ is a Poisson point process with rate r_n on $[\tau_1, \tau_2[$. Therefore,

$$\mathbb{P}(i \leftrightarrow j, i \leftrightarrow k) = \mathbb{P}(e_1 > \tau_1) \mathbb{P}(e_2 > \tau_2 - \tau_1),$$

where $e_1 = \inf(P_{\{ij\}}^* \cup P_{\{ik\}}^*) \sim \text{Exp}(2r_n)$ is independent of τ_1 and $e_2 = \inf(P_{\{uv\}}^* \cap [\tau_1, +\infty[) \sim \text{Exp}(r_n)$ is independent of $\tau_2 - \tau_1$. As a result,

$$\mathbb{P}(i \leftrightarrow j, i \leftrightarrow k) = \frac{3}{3 + 2r_n} \times \frac{1}{1 + r_n}.$$

A short calculation shows that

$$\text{Cov}\left(\mathbf{1}_{\{i \leftrightarrow j\}}, \mathbf{1}_{\{i \leftrightarrow k\}}\right) = \frac{r_n}{(3 + 2r_n)(1 + r_n)^2},$$

proving the proposition.

As before, the corollary follows from writing the degree of v as $D_n^{(v)} = \sum_{i \neq v} \mathbf{1}_{\{i \leftrightarrow v\}}$, which gives

$$\text{Var}(D_n^{(v)}) = (n-1) \text{Var}\left(\mathbf{1}_{\{i \leftrightarrow v\}}\right) + (n-1)(n-2) \text{Cov}\left(\mathbf{1}_{\{i \leftrightarrow v\}}, \mathbf{1}_{\{j \leftrightarrow v\}}\right).$$

Substituting $\text{Var}\left(\mathbf{1}_{\{i \leftrightarrow v\}}\right) = r_n/(1+r_n)^2$ and $\text{Cov}\left(\mathbf{1}_{\{i \leftrightarrow v\}}, \mathbf{1}_{\{j \leftrightarrow v\}}\right)$ yields the desired expression. \square

Proposition 3.5. *Let i, j, k and ℓ be four distinct vertices of G_{n,r_n} . We have*

$$\text{Cov}\left(\mathbf{1}_{\{i \leftrightarrow j\}}, \mathbf{1}_{\{k \leftrightarrow \ell\}}\right) = \frac{2r_n}{(1+r_n)^2(3+r_n)(3+2r_n)}$$

Corollary 3.6. *Let $D_n^{(i)}$ and $D_n^{(j)}$ be the respective degrees of two fixed vertices i and j , and let $|E_n|$ be the number of edges of G_{n,r_n} . We have*

$$\text{Cov}\left(D_n^{(i)}, D_n^{(j)}\right) = \frac{r_n}{(1+r_n)^2} \left(1 + \frac{3(n-2)}{3+2r_n} + \frac{2(n-2)(n-3)}{(3+r_n)(3+2r_n)}\right)$$

and

$$\text{Var}(|E_n|) = \frac{r_n n (n-1)(n^2 + 2r_n^2 + 2nr_n + n + 5r_n + 3)}{2(1+r_n)^2(3+r_n)(3+2r_n)}$$

The proof of Proposition 3.5 and its corollary are conceptually identical to the proofs of Propositions 3.1 and 3.3 and their corollaries, but the calculations are more tedious and so they have been relegated to Section A of the Appendix.

3.1.2 Complete subgraphs

From a biological perspective, complete subgraphs are interesting because they are related to how fine the partition of the set of populations into species can be. Indeed, the vertices of a complete subgraph – and especially of a large one – should be considered as part of the same species. A complementary point of view will be brought by connected components in Section 5.

In this section we establish the following results.

Proposition 3.7. *The expected number of complete subgraphs of order k in G_{n,r_n} is*

$$\mathbb{E}(X_{n,k}) = \binom{n}{k} \left(\frac{1}{1+r_n}\right)^{k-1}.$$

Corollary 3.8. *Let ω_n be the clique number of G_{n,r_n} , i.e. the maximal number of vertices in a complete subgraph of G_{n,r_n} . If (k_n) is such that*

$$\binom{n}{k_n} \left(\frac{1}{1+r_n}\right)^{k_n-1} \xrightarrow{n \rightarrow \infty} 0,$$

then k_n is asymptotically almost surely an upper bound on ω_n , i.e. $\mathbb{P}(\omega_n \leq k_n) \rightarrow 1$ as $n \rightarrow +\infty$. In particular, when $r_n \rightarrow +\infty$,

- (i) If $r_n = o(n)$, then $\omega_n \leq \log(r_n)n/r_n$ a.a.s.
- (ii) If $r_n = O(n/\log(n))$, $\omega_n = O_p(n/r_n)$, i.e.

$$\forall \varepsilon > 0, \exists M > 0, \exists N \text{ s.t. } \forall n \geq N, \mathbb{P}(\omega_n > Mn/r_n) < \varepsilon.$$

Proof of Proposition 3.7. The number of complete subgraphs of order k of G_{n,r_n} is

$$X_{n,k} = \sum_{S \in V^{(k)}} \mathbb{1}_{\{G_{n,r_n}[S] \text{ is complete}\}}$$

where the elements of $V^{(k)}$ are the k -subsets of $V = \{1, \dots, n\}$, and $G_{n,r_n}[S]$ is the subgraph of G_{n,r_n} induced by S . By exchangeability,

$$\mathbb{E}(X_{n,k}) = \binom{n}{k} \mathbb{P}(G_{n,r_n}[S] \text{ is complete}),$$

where S is any fixed set of k vertices. Using the notation of Proposition 2.3,

$$G_{n,r_n}[S] \text{ is complete} \iff \forall \{ij\} \in S, P_{\{ij\}}^* = \emptyset.$$

For all $t \geq 0$, let $A_t = \{a_t(i) : i \in S\}$ be the set of ancestors of S at t . Let $\tau_0 = 0$ and for each $\ell = 1, \dots, k-1$ let τ_ℓ be the time of the ℓ -th coalescence between two lineages of S , i.e.

$$\tau_\ell = \inf\{t > \tau_{\ell-1} : |A_t| = |A_{\tau_{\ell-1}}| - 1\}$$

Finally, let $\tilde{A}_\ell = A_{\tau_\ell}$ and $I_\ell = [\tau_\ell, \tau_{\ell+1}[$. With this notation,

$$\{\forall \{ij\} \in S, P_{\{ij\}}^* = \emptyset\} = \bigcap_{\ell=0}^{k-2} B_\ell,$$

where

$$B_\ell = \bigcap_{\{ij\} \in \tilde{A}_\ell^{(2)}} \{P_{\{ij\}}^* \cap I_\ell = \emptyset\}$$

and $\tilde{A}_\ell^{(2)}$ denotes the (unordered) pairs of \tilde{A}_ℓ . Since for $\ell \neq m$, $I_\ell \cap I_m = \emptyset$, Lemma 2.4 shows that conditionally on I_0, \dots, I_{k-1} , the events B_0, \dots, B_{k-2} are independent. By construction, for all $\{ij\} \neq \{uv\}$ in $\tilde{A}_\ell^{(2)}$,

$$\forall t \in I_\ell, \{a_t(i), a_t(j)\} \neq \{a_t(u), a_t(v)\} \neq \emptyset$$

and so it follows from Lemma 2.4 that, conditional on I_ℓ , $(P_{\{ij\}}^* \cap I_\ell)$, $\{ij\} \in \tilde{A}_\ell^{(2)}$, are i.i.d. Poisson point processes with rate r_n on I_ℓ . Therefore,

$$\mathbb{P}(B_\ell) = \mathbb{P}\left(\min\{e_{\{ij\}}^{(\ell)} : \{ij\} \in \tilde{A}_\ell^{(2)}\} > |I_\ell|\right),$$

where $e_{\{ij\}}^{(\ell)}$, $\{ij\} \in \tilde{A}_\ell^{(2)}$, are $\binom{k-\ell}{2}$ i.i.d. exponential variables with parameter r_n that are also independent of $|I_\ell|$. Since $|I_\ell| \sim \text{Exp}\left(\binom{k-\ell}{2}\right)$,

$$\mathbb{P}(B_\ell) = \frac{1}{1 + r_n}$$

and Proposition 3.7 follows. \square

Proof of Corollary 3.8. The first part of the corollary is a direct consequence of Proposition 3.7. First, note that

$$X_{n,k_n} = 0 \iff \omega_n < k_n$$

that a complete subgraph of order k contains complete subgraphs of order ℓ for all $\ell < k$. As a result, any k_n such that $\mathbb{P}(X_{n,k_n} = 0) \rightarrow 1$ is asymptotically almost surely an upper bound on the clique number ω_n . Now, observe that since X_{n,k_n} is a non-negative integer, $X_{n,k_n} \geq \mathbb{1}_{\{X_{n,k_n} \neq 0\}}$ and therefore,

$$\mathbb{E}(X_{n,k_n}) \geq \mathbb{P}(X_{n,k_n} \neq 0).$$

Finally, because $X_{n,k}$ is integer-valued, $\mathbb{P}(X_{n,k_n} \neq 0) \rightarrow 0$ implies $\mathbb{P}(X_{n,k_n} = 0) \rightarrow 1$.

To prove the second part of the corollary, using Stirling's formula we find that whenever r_n and k_n are $o(n)$ and go to $+\infty$ as $n \rightarrow +\infty$,

$$\binom{n}{k_n} \left(\frac{1}{1+r_n}\right)^{k_n-1} \sim \frac{C}{\sqrt{k_n} k_n^{k_n} (n-k_n)^{n-k_n}} \left(\frac{1}{1+r_n}\right)^{k_n-1},$$

where $C = \sqrt{2\pi}$. The right-hand side goes to zero if and only if its logarithm goes to $-\infty$, i.e. if and only if

$$A_n := k_n \log\left(\frac{n-k_n}{k_n(1+r_n)}\right) - n \log\left(1 - \frac{k_n}{n}\right) + \log\left(\frac{1+r_n}{\sqrt{k_n}}\right)$$

goes to $-\infty$. Now let $k_n = ng_n/r_n$, where $g_n \rightarrow +\infty$ and is $o(r_n)$, so that $k_n = o(n)$. Then,

$$k_n \log\left(\frac{n-k_n}{k_n(1+r_n)}\right) \sim -k_n \log(g_n)$$

and

$$-n \log\left(1 - \frac{k_n}{n}\right) \sim k_n.$$

Moreover, as long as it does not go to zero,

$$\log\left(\frac{1+r_n}{\sqrt{k_n}}\right) \sim \frac{3}{2} \log(r_n) - \frac{1}{2} \log(ng_n).$$

Putting the pieces together, we find that A_n is asymptotically equivalent to

$$-\frac{ng_n}{r_n} \log(g_n) + \frac{3}{2} \log(r_n) - \frac{1}{2} \log(ng_n).$$

Taking $g_n = \log(r_n)$, this expression goes to $-\infty$ as $n \rightarrow +\infty$, yielding (i). If $r_n = O(n/\log(n))$, then it goes to $-\infty$ for any $g_n \rightarrow +\infty$, which proves (ii). Indeed, if there exists $\varepsilon > 0$ such that

$$\forall M > 0, \forall N, \exists n \geq N \text{ s.t. } \mathbb{P}(\omega_n > Mn/r_n) \geq \varepsilon,$$

then considering successively $M = 1, 2, \dots$, we can find $n_1 < n_2 < \dots$ such that

$$\forall k \in \mathbb{N}, \mathbb{P}(\omega_{n_k} > kn_k/r_{n_k}) \geq \varepsilon.$$

Defining (g_n) by

$$\forall n \in \{n_k, \dots, n_{k+1} - 1\}, g_n = k,$$

we obtain a sequence (g_n) that goes to infinity and yet is such that for all N there exists $n := \min\{n_k : n_k \geq N\}$ such that $\mathbb{P}(\omega_n > g_n n/r_n) \geq \varepsilon$. \square

A natural pendant to Proposition 3.7 and Corollary 3.8 would be to use the variance of $X_{n,k}$ to find a lower bound on the clique number. Indeed, it follows from Chebychev's inequality that

$$\mathbb{P}(X_{n,k} = 0) \leq \frac{\text{Var}(X_{n,k})}{\mathbb{E}(X_{n,k})^2}.$$

However, computing $\text{Var}(X_{n,k})$ requires being able to compute the probability that two subsets of k vertices S and S' both induce a complete subgraph, which we have not managed to do. Using the probability that $G_{n,r_n}[S]$ is complete as an upper bound for this quantity, we have the very crude inequality

$$\text{Var}(X_{n,k}) \leq \binom{n}{k}^2 p(1-p),$$

where $p = 1/(1+r_n)^{k-1}$. This shows that when $r_n \rightarrow 0$ and $k_n = o(1/r_n)$, $\mathbb{P}(X_{n,k_n} = 0)$ tends to zero, proving that ω_n is at least $\Theta(1/r_n)$.

Finally, because we expect our model to form dense connected components, whose number we conjecture to be on the order of r_n in the intermediate regime (see Theorem 5.1 and Conjecture 5.4), and since the degree of a typical vertex is approximately n/r_n in that regime, it seems reasonable to conjecture

Conjecture 3.9. *In the intermediate regime, i.e. when $r_n \rightarrow +\infty$ and $r_n = o(n)$,*

$$\exists \alpha, \beta > 0 \text{ s.t. } \mathbb{P}(\alpha n/r_n \leq \omega_n \leq \beta n/r_n) \xrightarrow{n \rightarrow +\infty} 1.$$

3.2 Identification of different regimes

We now use the results of the previous section to identify different regimes for the behavior of G_{n,r_n} . The proof of our next theorem relies in part on results proved later in the paper (namely, Theorems 4.1 and 6.1), but no subsequent result depends on it, avoiding cyclic dependencies. While this section could have been placed at the end of the paper, it makes more sense to present it here because it relies mostly on Section 3.1 and because it helps structure the rest of the paper.

Theorem 3.10. *Let D_n be the degree of a fixed vertex of G_{n,r_n} . In the limit as $n \rightarrow +\infty$, depending on the asymptotics of r_n we have the following behaviors for G_{n,r_n}*

- (i) Transition for the complete graph: *when $r_n = o(1/n)$, $\mathbb{P}(G_{n,r_n}$ is complete) goes to 1, while when $r_n = \omega(1/n)$ it goes to 0; when $r_n = \Theta(1/n)$, this probability is bounded away from 0 and from 1.*
- (ii) Dense regime: *when $r_n = o(1)$, $\mathbb{P}(D_n = n - 1) \rightarrow 1$.*
- (iii) Sparse regime: *when $r_n = \omega(n)$, $\mathbb{P}(D_n = 0) \rightarrow 1$.*
- (iv) Transition for the empty graph: *when $r_n = o(n^2)$, $\mathbb{P}(G_{n,r_n}$ is empty) goes to 0 while when $r_n = \omega(n^2)$ it goes to 1; when $r_n = \Theta(n^2)$, this probability is bounded away from 0 and from 1.*

Proof. (i) is a direct consequence of Proposition 3.7 which, applied to $k = n$, yields

$$\mathbb{P}(G_{n,r_n} \text{ is complete}) = \left(\frac{1}{1+r_n} \right)^{n-1}.$$

(ii) is intuitive since $\mathbb{E}(D_n) = (n-1)/(1+r_n)$; but because it takes $r_n = o(1/n^2)$ for $\text{Var}(D_n)$ to go to zero, a second moment method is not sufficient to prove it. However, using Theorem 4.1, we see that $\mathbb{P}(D_n = n-1)$ can be written as

$$\mathbb{P}(D_n = n-1) = \frac{\Gamma(2+2r_n)\Gamma(n+1)}{\Gamma(n+1+2r_n)},$$

where Γ is the gamma function. The results follows by letting r_n go to zero and using the continuity of Γ .

(iii) follows from the same argument as in the proof of Corollary 3.8, by which, D_n being a non-negative integer, $\mathbb{P}(D_n \neq 0) \leq \mathbb{E}(D_n) = \frac{n-1}{1+r_n}$.

In (iv), the fact that G_n is empty when $r_n = \omega(n^2)$ is yet another application of this argument, but this time using the expected number of edges, $\mathbb{E}(|E_n|) = \frac{n(n-1)}{2(1+r_n)}$, in conjunction with the fact that G_n is empty if and only if $|E_n| = 0$; to see why the graph cannot be empty when $r_n = o(n^2)$, consider the edge that was created between the duplicated vertex and its copy in the most recent duplication. Clearly, if this edge has not disappeared yet then G_n cannot be empty. But the probability that this edge has disappeared is just

$$\frac{r_n}{\binom{n}{2} + r_n},$$

which goes to zero when $r_n = o(n^2)$. Finally, the fact that $\mathbb{P}(G_n \text{ is empty})$ is bounded away from 0 and from 1 when $r_n = \Theta(n^2)$ is a consequence of Theorem 6.1, which shows that the number of edges is Poissonian when $r_n = \omega(n)$. As a result, $\mathbb{P}(|E_n| = 0) \sim e^{-\mathbb{E}(|E_n|)}$. \square

Remark 3.11. Note that when $r_n = o(1)$, $\text{Var}(D_n) \sim r_n n^2/3$ can go to infinity even though $D_n = n - 1$ with probability that goes to 1. Similarly, when $r_n = o(1/n)$, $\text{Var}(|E_n|) \sim r_n n^4/18$ and $|E_n| = \binom{n}{2}$ a.a.s. Notably, $\overline{D}_n = (n - 1) - D_n$ converges to 0 in probability while $\text{Var}(\overline{D}_n)$ goes to infinity.

4 The degree distribution

The degree distribution is one of the most widely studied graph invariants in network science. Our model makes it possible to obtain an exact expression for its probability distribution:

Theorem 4.1 (degree distribution).

Let D_n be the degree of a fixed vertex of G_{n,r_n} . Then, for each $k \in \{0, \dots, n - 1\}$,

$$\mathbb{P}(D_n = k) = \frac{2r_n(2r_n + 1)}{(n + 2r_n)(n - 1 + 2r_n)} (k + 1) \prod_{i=1}^k \frac{n - i}{n - i + 2r_n - 1}, \quad (2)$$

where the empty product is 1.

The expression above holds for any positive sequence (r_n) and any n ; but as $n \rightarrow +\infty$ it becomes much simpler and, under appropriate rescaling, the degree converges to classical distributions:

Theorem 4.2 (convergence of the rescaled degree).

- (i) If $r_n \rightarrow r > 0$, then $\frac{D_n}{n}$ converges in distribution to a $\text{Beta}(2, 2r)$ random variable.
- (ii) If r_n is both $\omega(1)$ and $o(n)$, then $\frac{D_n}{n/r_n}$ converges in distribution to a size-biased exponential variable with parameter 2.
- (iii) If $2r_n/n \rightarrow \rho > 0$, then $D_n + 1$ converges in distribution to a size-biased geometric variable with parameter $\rho/(1 + \rho)$.

In this section we prove Theorem 4.1 by coupling the degree to the number of individuals descended from a founder in a branching process with immigration. Theorem 4.2 is then easily deduced by a standard study that has been relegated to Section B of the Appendix.

4.1 Ideas of the proof of Theorem 4.1

Before jumping to the formal proof of Theorem 4.1, we give a verbal account of the main ideas of the proof.

In order to find the degree of a fixed vertex v , we have to consider all pairs $\{iv\}$ and look at their ancestry to assess the absence/presence of atoms in the corresponding Poisson point processes. To do so, we can restrict our attention to the genealogy

of the vertices, and consider that edge-removal events occur along the lineages of this genealogy: a point that falls on the lineage of vertex i at time t means that $t \in P_{\{iv\}}^*$. In this setting, edge-removal events occur at constant rate r_n on every lineage different from that of v .

Next, the closed neighborhood of v (i.e. the set of vertices that are linked to v , plus v itself) can be obtained through the following procedure: we trace the genealogy of vertices, backwards in time; if we encounter an edge-removal event on lineage i at time t , then we mark all vertices that descend from this lineage, i.e. all vertices whose ancestor at time t is i ; only the lineages of unmarked vertices are considered after t . We stop when there is only one lineage left in the genealogy. The unmarked vertices are then exactly the neighbors of v (plus v itself). The procedure is illustrated in Figure 5.

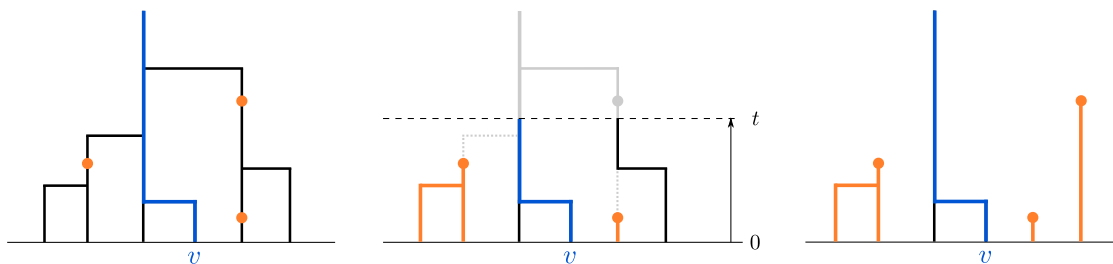


Figure 5: Illustration of the procedure used to find the neighborhood of v . On the left, the genealogy of the vertices, where the lineage of the focal vertex v has been highlighted in blue. A dot on lineage k corresponds to a point in $P_{\{ka_t(v)\}}$. In the middle, we uncover the genealogy and edge-removal events in backward time, as described in the main text. On the right, the forest that we get when the procedure is complete. The non-colored (black) branches are exactly the neighbors of v .

This vertex marking process is not convenient to describe in backward time because we typically mark several vertices simultaneously. By contrast, the forest that results from the completed process seems much easier to describe in forward time. Indeed, the arrival of a new lineage corresponds either to the addition of a new unmarked vertex or to the addition of a marked one, depending on whether the new lineage belongs to the same tree as v or not.

Moreover, in forward time, the process is reminiscent of a branching process with immigration: new lineages are either grafted to existing ones (branching) or sprout spontaneously (immigration). Let us try to find what the branching and immigration rates should be. In backward time, when there are $k + 1$ lineages then a coalescence occurs at rate $\binom{k+1}{2}$, while an edge-removal event occurs at rate $k r_n$. Reversing time, these events occur at the same rates. As a result, when going from k to $k + 1$ lineages, the probability that the next event is a branching is $(k + 1)/(k + 1 + 2 r_n)$.

Next, we have to find the probability that each lineage has to branch, given that the next event is a branching. Here, a classic argument suggests that every lineage branches at rate 1, except for the lineage of v , which branches at rate 2: in backward time, when going from $k + 1$ to k lineages, there are k pairs out of $\binom{k+1}{2}$ that involve the lineage of v . Therefore, the probability that the lineage of v is involved in the next coalescence is $2/(k + 1)$.

If this heuristic is correct, then in forward time it is easy to track the number of branches of the tree of v versus the number of branches of other trees: when there are p branches in the tree of v and q branches in the other trees, the probability that the next branch is added to the tree of v is just $(p + 1)/(p + 1 + q + 2r_n)$. Moreover, when the total number of branches reaches n , the number of branches in the tree of v is also the number of unmarked vertices at the end of the vertex marking procedure, which is itself $D_n^{(v)} + 1$, the degree of v plus one.

4.2 Formal proof of Theorem 4.1

The ideas and outline of the proof parallels the account given in the previous section: first, given a realization of the vertex-duplication process \mathcal{M} and of the edge-removal process \mathcal{P} , we describe a deterministic procedure that gives the closed neighborhood of any vertex v ,

$$N_G[v] = \{i \in V : \{iv\} \in E\} \cup \{v\},$$

where $G = (V, E)$ is the graph associated to \mathcal{M} and \mathcal{P} ; then, we identify the law of the process $(F_t)_{t \geq 0}$ corresponding to this procedure, and recognize it as the law of a branching process with immigration.

Definition 4.3. A *rooted forest with marked vertices* is a triple $F = (V^\circ, V^\bullet, \vec{E})$ such that

- (i) $V^\circ \cap V^\bullet = \emptyset$.
- (ii) Letting $V = V^\circ \cup V^\bullet$, (V, \vec{E}) is an acyclic digraph with the property that $\forall i \in V$, $\deg^+(i) \in \{0, 1\}$, where $\deg^+(i)$ is the out-degree of vertex i .

The marked vertices are the elements of V^\bullet ; the roots of F are the vertices with out-degree 0 (that is, edges are oriented towards the root), whose set we denote by $R(F)$; finally, the trees of F are its connected components (in the weak sense, i.e. considering the underlying undirected graph), and we write $T_F(i)$ for the tree containing i in F .

4.2.1 The vertex-marking process

We now define the backward-time process $(F_t)_{t \geq 0}$ that corresponds to the procedure described informally in Section 4.1. Recall the notation of Proposition 2.3. For a given realization of \mathcal{M} and \mathcal{P} , and for any fixed vertex v , let $(F_t)_{t \geq 0}$ be the piecewise constant process defined deterministically by

- $F_0 = (V, \emptyset, \emptyset)$.
- If $t \in M_{(ij)}$, then $\forall k, \ell \in R(F_{t-}) \cap V_{t-}^\circ$ such that $(a_{t-}(k), a_{t-}(\ell)) = (i, j)$,

$$\vec{E}_t = \vec{E}_{t-} \cup \{(k, \ell)\}.$$

- If $t \in P_{\{ia_t(v)\}}$, then letting $d_t(i) = \{j \in V : a_t(j) = i\}$ be the set of descendants of i born after time t ,

$$\begin{cases} V_t^\circ = V_{t-}^\circ \setminus d_t(i) \\ V_t^\bullet = V_{t-}^\bullet \cup d_t(i). \end{cases}$$

What makes $(F_t)_{t \geq 0}$ interesting is that

$$N_G[v] = V_\infty^\circ.$$

Indeed, by construction,

$$i \in V_t^\circ \iff \bigcup_{s \in [0, t]} \left(\bigcup_{j: i \in d_s(j)} P_{\{ja_s(v)\}} \right) = \emptyset,$$

and since for every s the unique j such that $i \in d_s(j)$ is $a_s(i)$, we have

$$V_t^\circ = \left\{ i \in V : P_{\{iv\}}^* \cap [0, t] = \emptyset \right\}.$$

The Poissonian construction given above shows that $(F_t, a_t)_{t \geq 0}$ is a Markov process. Now, observe that conditional on a_t

- (i) $M_{(ij)} \cap]t, +\infty[\sim M_{(a_t(i)a_t(j))} \cap]t, +\infty[$ and is independent of $(F_s, a_s)_{s \leq t}$
- (ii) $P_{\{ia_t(v)\}} \cap]t, +\infty[\sim P_{\{a_t(i)a_t(v)\}} \cap]t, +\infty[$ and is independent of $(F_s, a_s)_{s \leq t}$
- (iii) $j \in d_t(i) \iff i \in R(F_t)$ and $j \in T_{F_t}(i)$

As a consequence, $(F_t)_{t \geq 0}$ is also a Markov process, whose law is characterized by

- $F_0 = (V, \emptyset, \emptyset)$.
- F_t goes from $(V_t^\circ, V_t^\bullet, \vec{E}_t)$ to
 - $(V_t^\circ, V_t^\bullet, \vec{E}_t \cup \{(i, j)\})$ at rate $1/2$, for all i, j in $R(F_t)$
 - $(V_t^\circ \setminus T_{F_t}(i), V_t^\bullet \cup T_{F_t}(i), \vec{E}_t)$ at rate r_n , for all i in $R(F_t)$.

Let $(\tilde{F}_k)_{k \in \{1, \dots, n\}}$ be the chain embedded in $(F_t)_{t \geq 0}$, i.e. defined by

$$\tilde{F}_k = F_{t_k}, \text{ where } t_k = \inf\{t \geq 0 : |R(F_t)| = n - k + 1\}.$$

The rooted forests with marked vertices that correspond to realizations of \tilde{F}_n are exactly the $f_n = (V^\circ, V^\bullet, \vec{E})$ that have n vertices and are such that $V^\circ = T_{f_n}(v)$. Moreover, for each of these there exists a unique trajectory (f_1, \dots, f_n) of $(\tilde{F}_1, \dots, \tilde{F}_n)$ such that $\tilde{F}_n = f_n$ and it follows from the transition rates of $(F_t)_{t \geq 0}$ that

$$\begin{aligned} \mathbb{P}(\tilde{F}_n = f_n) &= \frac{(1/2)^{n-|R(f_n)|} r_n^{|R(f_n)|-1}}{\prod_{k=2}^n (k(k-1)/2 + (k-1)r_n)} \\ &= \frac{1}{(n-1)!} \times \frac{(2r_n)^{|R(f_n)|-1}}{\prod_{k=2}^n (k+2r_n)} \end{aligned} \quad (3)$$

Finally, note that $\tilde{V}_n^\circ = V_\infty^\circ$ is the closed neighborhood of v in our graph.

4.2.2 The branching process

The process with which we will couple the vertex-marking process described in the previous section is a simple random function of the trajectories of a branching process with immigration $(Z_t)_{t \geq 0}$. In this branching process, immigration occurs at rate $2r_n$ and each particle gives birth to a new particle at rate 1 – except for one particle, which carries a special item that enables it to give birth at rate 2; when this lineage reproduces, it keeps the item with probability 1/2, and passes it to its offspring with probability 1/2.

Formally, we consider the Markov process on the set of rooted forests with marked vertices (augmented with an indication of the carrier of the item), defined by $Z_0 = (\{1\}, \emptyset, \emptyset, 1)$ and by the following transition rates:

$(Z_t)_{t \geq 0}$ goes from $(W^\circ, W^\bullet, \vec{E}, c)$ to

- $(W^\circ \cup \{N\}, W^\bullet, \vec{E} \cup \{(N, i)\}, c)$ at rate 1, for all $i \in W^\circ$
- $(W^\circ, W^\bullet \cup \{N\}, \vec{E} \cup \{(N, i)\}, c)$ at rate 1, for all $i \in W^\bullet$
- $(W^\circ \cup \{N\}, W^\bullet, \vec{E} \cup \{(N, c)\}, N)$ at rate 1
- $(W^\circ, W^\bullet \cup \{N\}, \vec{E}, c)$ at rate $2r_n$

where $N = |W^\circ \cup W^\bullet| + 1$ is the label of the new particle. The fourth coordinate of $(Z_t)_{t \geq 0}$ tracks the carrier of the item.

As previously, the Markov chain $(\tilde{Z}_k)_{k \in \mathbb{N}^*}$ embedded in $(Z_t)_{t \geq 0}$ is defined by

$$\tilde{Z}_k = Z_{t_k}, \text{ where } t_k = \inf\{t \geq 0 : |W_t^\circ \cup W_t^\bullet| = k\}.$$

The realizations of \tilde{Z}_n are exactly the $(W_n^\circ, W_n^\bullet, \vec{E}_n, c_n)$ such that $f_n = (W_n^\circ, W_n^\bullet, \vec{E}_n)$ is a rooted forest with marked vertices on $\{1, \dots, n\}$ and $W_n^\circ = T_{f_n}(1) = T_{f_n}(c_n)$. For these, it follows from the transition rates of $(Z_t)_{t \geq 0}$ that

$$\mathbb{P}(\tilde{Z}_n = (W_n^\circ, W_n^\bullet, \vec{E}_n, c_n)) = \frac{(2r_n)^{|R(f_n)|-1}}{\prod_{k=1}^{n-1} (k+1+2r_n)}. \quad (4)$$

Finally, note that $(X_k)_{k \in \mathbb{N}^*} = (|\tilde{W}_k^\circ|, |\tilde{W}_k^\bullet|)_{k \in \mathbb{N}^*}$, which counts the number of descendants of the first particle and the number of descendants of immigrants, is a Markov chain whose law is characterized by $X_1 = (1, 0)$ and X_k goes from (p, q) to

- $(p+1, q)$ with probability $\frac{p+1}{p+1+q+2r_n}$
- $(p, q+1)$ with probability $\frac{q+2r_n}{p+1+q+2r_n}$.

4.2.3 Relabeling and end of proof

The last step before finishing the proof of Theorem 4.1 is to shuffle the vertices of the forest associated to \tilde{Z}_n appropriately. For any fixed n, v and c in $\{1, \dots, n\}$, let $\Phi_{(c,v)}$ be uniformly and independently of anything else picked among the permutations of $\{1, \dots, n\}$ that map c to v ; define $\Phi_v(\tilde{Z}_n)$ by

$$\Phi_v(\tilde{W}_n^\circ, \tilde{W}_n^\bullet, \tilde{E}_n, \tilde{c}_n) = \left(\Phi_{(\tilde{c}_n,v)}(\tilde{W}_n^\circ), \Phi_{(\tilde{c}_n,v)}(\tilde{W}_n^\bullet), \Phi_{(\tilde{c}_n,v)}(\tilde{E}_n) \right)$$

where $\Phi_{(\tilde{c}_n,v)}(\tilde{E}_n)$ is to be understood as $\left\{ \left(\Phi_{(\tilde{c}_n,v)}(i), \Phi_{(\tilde{c}_n,v)}(j) \right) : (i, j) \in \tilde{E}_n \right\}$.

With all these elements, the proof of Theorem 4.1 goes as follows. First, from equations (3) and (4) and the definition of Φ_v , we see that for all rooted forest with marked vertices f_n ,

$$\mathbb{P}(\tilde{F}_n = f_n) = \mathbb{P}(\Phi_v(\tilde{Z}_n) = f_n).$$

In particular, \tilde{V}_n° , the set of unmarked vertices in the vertex-marking process, and $\Phi_{(\tilde{c}_n,v)}(\tilde{W}_n^\circ)$, the relabeled set of descendants of the first particle in the branching process, have the same law. Now, on the one hand we have

$$|\tilde{V}_n^\circ| = |N_G[v]| = D_n^{(v)} + 1,$$

and on the other hand we have

$$|\Phi_{(\tilde{c}_n,v)}(\tilde{W}_n^\circ)| = |\tilde{W}_n^\circ|.$$

Since $|\tilde{W}_n^\circ|$ is the first coordinate of the Markov chain $(X_k)_{k \in \mathbb{N}^*}$ introduced in the previous section, it follows directly from the transition probabilities of $(X_k)_{k \in \mathbb{N}^*}$ that

$$\mathbb{P}(X_n = (k+1, n-k-1)) = \binom{n-1}{k} \frac{\prod_{p=1}^k (p+1) \prod_{q=0}^{n-k-2} (q+2r_n)}{\prod_{(p+q)=1}^{n-1} ((p+q)+1+2r_n)},$$

from which the expression of Theorem 4.1 can be deduced through elementary calculations.

5 Connected components in the intermediate regime

From a biological perspective, connected components are good candidates to define species, and have frequently been used to that end. Moreover, among the possible definitions of species, they play a special role because they indicate how coarse the partition of the set of populations into species can be; indeed, it would not make sense biologically for distinct connected components to be part of the same species.

As a result, connected components are in a sense the “loosest” possible definition of species. This complements the perspective brought by complete subgraphs, which inform us on how fine the species partition can be (see Section 3.1.2). For a discussion of the definition of species in a context where traits and ancestral relationships between individuals are known, see [8].

The aim of this section is to prove the following theorem.

Theorem 5.1. *Let $\#\text{CC}_n$ be the number of connected components of G_{n,r_n} . If r_n is both $\omega(1)$ and $o(n)$, then*

$$\frac{r_n}{2} + o_p(r_n) \leq \#\text{CC}_n \leq 2r_n \log n + o_p(r_n \log n)$$

where, for a positive sequence (u_n) , $o_p(u_n)$ stands for some sequence of random variables (X_n) such that $X_n/u_n \rightarrow 0$ in probability.

5.1 Lower bound on the number of connected components

The proof of the lower bound on the number of connected components uses the forward construction introduced in Section 2.2 and the associated notation. It relies on the simple observation that, letting $\#\text{CC}(G)$ denote the number of connected components of a graph G , $\#\text{CC}(G_{r_n}^*(k))$ is a nondecreasing function of k . Indeed, in the sequence of events defining $(G_{r_n}^*(k))_{k \geq 2}$, vertex duplications do not change the number of connected components – because a new vertex is always linked to an existing vertex (its ‘mother’) and her neighbors – and edge removals can only increase it. Thus, if $m_n \leq n$ and ℓ_n are such that $\mathbb{P}(\#\text{CC}(G_{r_n}^*(m_n)) \geq \ell_n) \xrightarrow{n \rightarrow \infty} 1$, then ℓ_n is asymptotically almost surely a lower bound on the number of connected components of $G_{r_n}^*(n)$ — and therefore of G_{n,r_n} .

To find such a pair (m_n, ℓ_n) , note that, for every graph G of order m ,

$$\#\text{CC}(G) \geq m - \#\text{edges}(G).$$

Moreover, since for any fixed n , $G_{r_n}^*(m_n)$ has the same law as G_{m_n, r_n} , the exact expressions for the expectation and the variance of $|E_{m_n}^*|$, the number of edges of $G_{r_n}^*(m_n)$, are given in Table 1. We see that, if r_n and m_n are both $\omega(1)$ and $o(n)$,

$$\mathbb{E}(|E_{m_n}^*|) \sim \frac{m_n^2}{2r_n} \quad \text{and} \quad \text{Var}(|E_{m_n}^*|) \sim \frac{m_n^2}{4r_n^3} (m_n^2 + 2r_n^2).$$

By Chebychev’s inequality,

$$\mathbb{P}\left(\left||E_{m_n}^*| - \mathbb{E}(|E_{m_n}^*|)\right| \geq m_n^{1-\varepsilon}\right) \leq \frac{\text{Var}(|E_{m_n}^*|)}{m_n^{2-2\varepsilon}}.$$

When $m_n = \Theta(r_n)$, since $r_n = \omega(1)$ the right-hand side of this inequality goes to 0 as $n \rightarrow +\infty$, for all $\varepsilon < 1/2$. It follows that

$$|E_{m_n}^*| = \mathbb{E}(|E_{m_n}^*|) + o_p(r_n).$$

Taking $m_n := \lfloor \alpha r_n \rfloor$, we find that

$$\#\text{CC}(G_{r_n}^*(m_n)) \geq m_n - |E_{m_n}^*| = \alpha \left(1 - \frac{\alpha}{2}\right) r_n + o_p(r_n).$$

The right-hand side is maximal for $\alpha = 1$ and is then $r_n/2 + o_p(r_n)$.

5.2 Upper bound on the number of connected components

Our strategy to get an upper bound on the number of connected components is to find a spanning subgraph whose number of connected components we can estimate. A natural idea is to look for a spanning forest, because forests have the property that their number of connected components is their number of vertices minus their number of edges.

Definition 5.2. A pair of vertices $\{ij\}$ is said to be a *founder* if it has no ancestor other than itself, i.e., letting $T_{\{ij\}} = \sup\{t \geq 0 : a_t(i) \neq a_t(j)\}$ be the coalescence time of i and j , $\{ij\}$ is a founder if and only if $\forall t < T_{\{ij\}}, \{a_t(i), a_t(j)\} = \{ij\}$.

Let \mathcal{F} be the set of founders of $G_{n,r_n} = (V, E)$, and let $T_n = (V, \mathcal{F})$. Note that $\#\mathcal{F} = n - 1$ and that T_n is a tree. Therefore, letting $F_n = (V, \mathcal{F} \cap E)$ be the spanning forest of G_{n,r_n} induced by T_n , we have

$$\#\text{CC}_n \leq n - \#\text{edges}(F_n).$$

Let us estimate the number of edges of F_n . Recall Proposition 2.3. By construction, $\forall \{ij\} \in \mathcal{F}$, $P_{\{ij\}}^* = P_{\{ij\}} \cap [0, T_{\{ij\}}]$. It follows that

$$\#\text{edges}(F_n) = \sum_{\{ij\} \in \mathcal{F}} \mathbb{1}_{\{P_{\{ij\}} \cap [0, T_{\{ij\}}] = \emptyset\}}$$

and, as a consequence,

$$\#\text{CC}_n \leq 1 + \sum_{\{ij\} \in \mathcal{F}} \mathbb{1}_{\{P_{\{ij\}} \cap [0, T_{\{ij\}}] \neq \emptyset\}}.$$

Now, $\mathbb{1}_{\{P_{\{ij\}} \cap [0, T_{\{ij\}}] \neq \emptyset\}} \leq \#(P_{\{ij\}} \cap [0, T_{\{ij\}}])$, and since $(P_{\{ij\}})_{\{ij\} \in \mathcal{F}}$ are i.i.d. Poisson point processes with intensity r_n that are also independent of $(T_{\{ij\}})_{\{ij\} \in \mathcal{F}}$,

$$\sum_{\{ij\} \in \mathcal{F}} \#(P_{\{ij\}} \cap [0, T_{\{ij\}}]) \leq \#(P \cap [0, L_n]),$$

where P is a Poisson point process on $]0, +\infty]$ with intensity r_n and $L_n = T_{\text{MRCA}} + \sum_{\{ij\} \in \mathcal{F}} T_{\{ij\}}$ is the total branch length of the genealogy of the vertices. Putting the pieces together,

$$\#\text{CC}_n \leq 1 + \#(P \cap [0, L_n]).$$

Conditional on L_n , $\#(P \cap [0, L_n])$ is a Poisson random variable with parameter $r_n L_n$. Moreover, it is known [15] that

$$\mathbb{E}(L_n) = 2 \sum_{i=1}^{n-1} \frac{1}{i} \quad \text{and} \quad \text{Var}(L_n) = 4 \sum_{i=1}^{n-1} \frac{1}{i^2}$$

As a result,

$$\mathbb{E}(\#(P \cap [0, L_n])) = r_n \mathbb{E}(L_n) \sim 2 r_n \log n$$

and

$$\text{Var}(\#(P \cap [0, L_n])) = r_n \mathbb{E}(L_n) + \text{Var}(r_n L_n) \sim 2 r_n \log n + \alpha r_n^2,$$

with $\alpha = 2\pi^2/3$. Using Chebychev's inequality, we find that for all $\varepsilon > 0$,

$$\mathbb{P}\left(|\#(P \cap [0, L_n]) - 2 r_n \log n| \geq \varepsilon r_n \log n\right) = O\left(\frac{2}{\varepsilon^2 r_n \log(n)} + \frac{\alpha}{\varepsilon^2 \log(n)^2}\right).$$

The right-hand side goes to 0 as $n \rightarrow +\infty$, which shows that $\#(P \cap [0, L_n]) - 2 r_n \log n = o_p(r_n \log n)$ and finishes the proof.

Remark 5.3. Using $\#(P \cap [0, L_n])$ as an upper bound for $\sum_{\{ij\} \in \mathcal{F}} \mathbb{1}_{\{P_{\{ij\}} \cap [0, T_{\{ij\}}] \neq \emptyset\}}$ turns out not to be a great source of imprecision, because most of the total branch length of a Kingman coalescent comes from very short branches. As a result, when $r_n = o(n)$, only a negligible proportion of the $P_{\{ij\}} \cap [0, T_{\{ij\}}]$'s, $\{ij\} \in \mathcal{F}$, have more than one point.

By contrast, using $n - \#\text{edges}(F_n)$ as an upper bound on $\#\text{CC}_n$ is very crude. This leads us to formulate the following conjecture:

Conjecture 5.4.

$$\exists \alpha, \beta > 0 \text{ s.t. } \mathbb{P}(\alpha r_n \leq \#\text{CC}_n \leq \beta r_n) \xrightarrow[n \rightarrow \infty]{} 1.$$

6 Number of edges in the sparse regime

From the expressions obtained in section 3.1.1 and recapitulated in Table 1, we see that when $r_n = \omega(n)$,

$$\mathbb{E}(|E_n|) \sim \text{Var}(|E_n|) \sim \frac{n^2}{2 r_n}.$$

This suggests that the number of edges is Poissonian in the sparse regime, and this is what the next theorem states.

Theorem 6.1. *Let $|E_n|$ be the number of edges of G_{n,r_n} . If $r_n = \omega(n)$ then*

$$d_{\text{TV}}(|E_n|, \text{Poisson}(\lambda_n)) \xrightarrow[n \rightarrow +\infty]{} 0,$$

where d_{TV} stands for the total variation distance and $\lambda_n = \mathbb{E}(|E_n|) \sim \frac{n^2}{2 r_n}$. If in addition $r_n = o(n^2)$, then $\lambda_n \rightarrow +\infty$ and as a result

$$\frac{|E_n| - \lambda_n}{\sqrt{\lambda_n}} \xrightarrow[n \rightarrow +\infty]{\mathcal{D}} \mathcal{N}(0, 1),$$

where $\mathcal{N}(0, 1)$ denotes the standard normal distribution.

The proof of Theorem 6.1 is a standard application of the Stein–Chen method [14, 2]. A reference on the topic is [1], and another excellent survey is given in [12]. Let us state briefly the results that we will need.

Definition A. The Bernoulli variables X_1, \dots, X_N are said to be *positively related* if for each $i = 1, \dots, N$ there exists $(X_1^{(i)}, \dots, X_N^{(i)})$, built on the same space as (X_1, \dots, X_N) , such that

- (i) $(X_1^{(i)}, \dots, X_N^{(i)}) \sim (X_1, \dots, X_N) \mid X_i = 1$.
- (ii) For all $j = 1, \dots, N$, $X_j^{(i)} \geq X_j$.

Note that other equivalent definitions of positive relation are known (see e.g. Lemma 4.27 in [12]).

We will also need the following classic theorem, which appears, e.g., as Theorem 4.20 in [12]:

Theorem A. *Let X_1, \dots, X_N be positively related Bernoulli variables with $\mathbb{P}(X_i = 1) = p_i$. Let $W = \sum_{i=1}^N X_i$ and $\lambda = \mathbb{E}(W)$. Then,*

$$d_{\text{TV}}(W, \text{Poisson}(\lambda)) \leq \min\{1, \lambda^{-1}\} \left(\text{Var}(W) - \lambda + 2 \sum_{i=1}^N p_i^2 \right).$$

It is intuitive that the variables indicating the presence of edges in our graph are positively related, because the only way through which these variables depend on each other is through the fact that the edges share ancestors. But we need to prove it before we can use Theorem A.

6.1 Proof of the positive relation between the edges

6.1.1 Preliminary lemmas

In this section we isolate the proof of two useful results that are not tied to the particular setting of our model.

Lemma 6.2. *Let $\mathbf{X} = (X_1, \dots, X_N)$ be a vector of Bernoulli variables. The distribution of \mathbf{X} is characterized by the quantities*

$$\mathbb{E} \left(\prod_{i \in I} X_i \right), \quad I \subset \{1, \dots, N\}, I \neq \emptyset$$

Proof. For all $I \subset \{1, \dots, N\}$, $I \neq \emptyset$, let

$$p_I = \mathbb{E} \left(\prod_{i \in I} X_i \right) \quad \text{and} \quad q_I = \mathbb{E} \left(\prod_{i \in I} X_i \prod_{j \in I^c} (1 - X_j) \right)$$

where the empty product is understood to be 1.

Clearly, the distribution of \mathbf{X} is fully specified by (q_I) . Now, observe that

$$\prod_{j \in I^c} (1 - X_j) = \sum_{K \subset I^c} (-1)^{|K|} \prod_{k \in K} X_k.$$

Therefore,

$$\begin{aligned} q_I &= \sum_{K \subset I^c} (-1)^{|K|} p_{I \cup K} \\ &= \sum_{J \supset I} (-1)^{|J| - |I|} p_J, \end{aligned}$$

which terminates the proof. \square

Lemma 6.3. *Let X_1, \dots, X_N be independent random nondecreasing functions from $[0, +\infty[$ to $\{0, 1\}$ such that, $\forall i \in \{1, \dots, N\}$, $\inf\{t \geq 0 : X_i(t) = 1\} < +\infty$ almost surely. Let T be a non-negative random variable that is independent of (X_1, \dots, X_N) . Then, $X_1(T), \dots, X_N(T)$ are positively related.*

Proof. Pick $i \in \{1, \dots, N\}$. Now, let $\tau_i = \inf\{t \geq 0 : X_i(t) = 1\}$. Assume without loss of generality that X_i is left-continuous, so that $\{X_i(T) = 1\} = \{T > \tau_i\}$. Next, note that,

$$\forall x, t \geq 0, \quad \mathbb{P}(T > x, T > t) \geq \mathbb{P}(T > x)\mathbb{P}(T > t).$$

Integrating in t against the law of τ_i , we find that

$$\forall x \geq 0, \quad \mathbb{P}(T > x \mid T > \tau_i) \geq \mathbb{P}(T > x).$$

This shows that T is stochastically dominated by $T^{(i)}$, where $T^{(i)}$ has the law of T conditioned on $\{T > \tau_i\}$. As a result, there exists S , built on the same space as X_1, \dots, X_N and independent of $(X_j)_{j \neq i}$, such that $S \sim T^{(i)}$ and $S \geq T$. For all $j \neq i$, let $X_j^{(i)} = X_j(S)$. Since X_j is nondecreasing, $X_j^{(i)} \geq X_j(T)$, and since $(X_j)_{j \neq i} \perp (T, \tau_i)$, $(X_j^{(i)})_{j \neq i} \sim ((X_j(T))_{j \neq i} \mid X_i(T) = 1)$. This shows that $X_1(T), \dots, X_N(T)$ are positively related. \square

Remark 6.4. Lemma 6.3 and its proof are easily adapted to the case where X_1, \dots, X_N are nonincreasing and such that $\inf\{t \geq 0 : X_i(t) = 0\} < +\infty$ a.s.

6.1.2 Stein–Chen coupling

To prove that the variables indicating the presence of the edges of G_{n, r_n} are positively related, we use the forward construction described in Section 2.3, and proceed by induction. Note that, in this section, we make no assumptions on r_n .

To keep the notation light, throughout the rest of this section we index the pairs of vertices of $G_r^*(n) = (\{1, \dots, n\}, E_n^*)$ by the integers from 1 to $N = \binom{n}{2}$ and, for $i \in \{1, \dots, N\}$, we let $X_i = \mathbb{1}_{\{i \in E_n^*\}}$. We also make consistent use of bold letters to denote vectors, i.e., given any family of random variables Z_1, \dots, Z_p , we write \mathbf{Z} for (Z_1, \dots, Z_p) .

For small values of n (e.g., $n = 3$), the law of $G_{n,r}$ — and hence of $G_r^*(n)$ — can be worked out explicitly. It is then straightforward to check that X_1, \dots, X_N are positively related for all values of r .

Now assume that X_1, \dots, X_N are positively related in $G_r^*(n)$, i.e.

$$\begin{aligned} \forall i \leq N, \exists \mathbf{Y}^{(i)} = (Y_1^{(i)}, \dots, Y_N^{(i)}) \text{ such that} \\ \text{(i) } \mathbf{Y}^{(i)} \sim (\mathbf{X} \mid X_i = 1) \\ \text{(ii) } \forall k \leq N, Y_k^{(i)} \geq X_k \end{aligned} \tag{5}$$

Remember that $G_r^*(n+1)$ is obtained by (1) adding a vertex to $G_r^*(n)$ (which, without loss of generality, we label $n+1$) and linking it to a uniformly chosen vertex u_n of $G_r^*(n)$ as well as to the neighbors of u_n ; and (2) waiting for an exponential time T with parameter $\binom{n}{2}$ while removing each edge at constant rate r .

Formally, $\forall k \leq N+n$, we define the “mother” of k , $M_k \in \{1, \dots, N\} \cup \{\emptyset\}$, by

- If $k \leq N$ (i.e., if k is the label of $\{u, v\}$, with $1 \leq u < v \leq n$), then $M_k = k$.
- If $k > N$ is the label of $\{v, n+1\}$, with $1 \leq v \leq n$, then $M_k = \ell$, where ℓ is the label of $\{u_n, v\}$.
- If $k > N$ is the label of $\{u_n, n+1\}$, then $M_k = \emptyset$.

Letting $X'_k = \mathbb{1}_{\{k \in E_{n+1}^*\}}$, we then have

$$X'_k = \begin{cases} A_k & \text{if } M_k = \emptyset \\ X_{M_k} A_k & \text{otherwise} \end{cases}$$

with $A_k = \mathbb{1}_{\{e_k > T\}}$, where we recall that $T \sim \text{Exp}(N)$ and, e_1, \dots, e_{N+n} are i.i.d. exponential variables with parameter r that are also independent of everything else.

Note that the random functions $\tilde{A}_k: t \mapsto \mathbb{1}_{\{e_k > t\}}$, $k \in \{1, \dots, N+n\}$ are non-increasing and such that $\inf\{t \geq 0 : \tilde{A}_k(t) = 0\} < +\infty$ almost surely. By Lemma 6.3 (see also Remark 6.4), it follows that A_1, \dots, A_{N+n} are positively related.

We now pick any $i \leq \binom{n+1}{2} = N+n$ and build a vector $\mathbf{Y}'^{(i)}$ that has the same law as $(\mathbf{X}' \mid X_i = 1)$ and satisfies $\mathbf{Y}'^{(i)} \geq \mathbf{X}'$.

Assume that $M_i \neq \emptyset$. In that case,

1. By the induction hypothesis, there exists $\mathbf{Y}^{(M_i)}$ that satisfies (5).
2. Since by A_1, \dots, A_{N+n} are positively related, $\exists \mathbf{B}^{(i)} \sim (\mathbf{A} \mid A_i = 1)$ such that $\mathbf{B}^{(i)} \geq \mathbf{A}$.

Note that \mathbf{A} , $\mathbf{B}^{(i)}$, \mathbf{X} and $\mathbf{Y}^{(M_i)}$ are all built on the same space. Therefore, omitting the (M_i) and (i) superscripts to keep the notation light, we can set $Y'_i = 1$ and, for $k \neq i$,

$$Y'_k = \begin{cases} B_k & \text{if } M_k = \emptyset \\ Y_{M_k} B_k & \text{otherwise.} \end{cases}$$

With this definition, $\forall J \subset \{1, \dots, N+n\}$,

$$\mathbb{E} \left(\prod_{j \in J} Y'_j \right) = \mathbb{E} \left(\prod_{j \in \tilde{J}} Y_j \right) \mathbb{E} \left(\prod_{j \in J} B_j \right),$$

where $\tilde{J} = \{M_j : j \in J, M_j \neq \emptyset\}$. By hypothesis,

$$\mathbb{E} \left(\prod_{j \in \tilde{J}} Y_j \right) = \mathbb{E} \left(\prod_{j \in \tilde{J}} X_j \mid X_{M_i} = 1 \right) = \mathbb{E} \left(X_{M_i} \prod_{j \in \tilde{J}} X_j \right) / \mathbb{E}(X_{M_i})$$

Similarly,

$$\mathbb{E} \left(\prod_{j \in J} B_j \right) = \mathbb{E} \left(A_i \prod_{j \in J} A_j \right) / \mathbb{E}(A_i).$$

As a result,

$$\begin{aligned} \mathbb{E} \left(\prod_{j \in J} Y'_j \right) &= \frac{\mathbb{E}(X_{M_i} \prod_{j \in \tilde{J}} X_j) \mathbb{E}(A_i \prod_{j \in J} A_j)}{\mathbb{E}(X_{M_i}) \mathbb{E}(A_i)} \\ &= \frac{\mathbb{E}(X_{M_i} A_i \prod_{j \in J} X'_j)}{\mathbb{E}(X_{M_i} A_i)} \\ &= \mathbb{E} \left(\prod_{j \in J} X'_j \mid X'_i = 1 \right) \end{aligned}$$

By Lemma 6.2, this shows that $\mathbf{Y}' \sim (\mathbf{X}' \mid X'_i = 1)$.

If $M_i = \emptyset$, we can no longer choose $\mathbf{Y}^{(M_i)}$. However, in that case, X'_i depends only on A_i . Therefore, we set $Y'_i = 1$ and, for $k \neq i$,

$$Y'_k = X_{M_k} B_k$$

Remembering that $X'_i = A_i$, we then check that

$$\mathbb{E} \left(\prod_{j \in J} Y'_j \right) = \frac{\mathbb{E}(\prod_{j \in \tilde{J}} X_j) \mathbb{E}(A_i \prod_{j \in J} A_j)}{\mathbb{E}(A_i)} = \mathbb{E} \left(\prod_{j \in J} X'_j \mid X'_i = 1 \right).$$

Finally, it is clear that, with both constructions of $\mathbf{Y}^{(i)}$, $\mathbf{Y}'_k^{(i)} \geq \mathbf{X}'_k$.

6.2 Proof of Theorem 6.1

Applying Theorem A to $|E_n| = \sum_{\{ij\}} \mathbb{1}_{\{i \leftrightarrow j\}}$ and using the expressions in table 1, we get

$$d_{\text{TV}}(|E_n|, \text{Poisson}(\lambda_n)) \leq \min\{1, \lambda_n^{-1}\} C_n,$$

with $\lambda_n = \frac{n(n-1)}{2(r_n+1)}$ and

$$C_n = \frac{n(n-1)(n^2 r_n + 2n r_n^2 + n r_n - 2r_n^2 + 3r_n + 9)}{2(2r_n + 3)(r_n + 3)(r_n + 1)^2}.$$

When $r_n = \omega(n)$,

$$C_n = \Theta\left(\frac{n^4}{r_n^3} + \frac{n^3}{r_n^2}\right).$$

Now, if $r_n > \frac{n(n-1)}{2} - 1$, so that $\min\{1, \lambda_n^{-1}\} = 1$, we see that $C_n = \Theta(n^3/r_n^2)$. If by contrast $r_n \leq \frac{n(n-1)}{2} - 1$ then $\lambda_n^{-1} C_n = \Theta(n/r_n)$. In both cases, $\min\{1, \lambda_n^{-1}\} C_n$ goes to zero as $n \rightarrow +\infty$, proving the first part of Theorem 6.1.

The convergence of $\frac{|E_n| - \lambda_n}{\sqrt{\lambda_n}}$ to the standard normal distribution is a classic consequence of the conjunction of $d_{\text{TV}}(|E_n|, \text{Poisson}(\lambda_n)) \rightarrow 0$ with $\lambda_n \rightarrow +\infty$. See, e.g., [1], page 17, where this is recovered as a consequence of inequality (1.39).

Acknowledgements

François Bienvenu thanks Jean-Jil Duchamps for helpful discussions. All authors thank the *Center for Interdisciplinary Research in Biology* (CIRB) for funding. Florence Débarre thanks the Agence Nationale de la Recherche for funding (grant ANR-14-ACHN- 0003-01).

References

- [1] A. D. Barbour, L. Holst, and S. Janson. *Poisson approximation*. Oxford Studies in Probability. Oxford University Press, 1992.
- [2] L. H. Y. Chen. Poisson approximation for dependent trials. *Annals of Probability*, 3(3):534–545, 1975.
- [3] F. Chung, L. Lu, T. G. Dewey, and D. J. Galas. Duplication models for biological networks. *Journal of Computational Biology*, 10(5):677–687, 2003.
- [4] J. A. Coyne and H. A. Orr. *Speciation*. Sinauer Associates, 2004.
- [5] R. Durrett. *Probability: Theory and Examples*. Cambridge University Press, 4th edition, 2010.
- [6] I. Ispolatov, P. Krapivsky, and A. Yuryev. Duplication-divergence model of protein interaction network. *Physical Review E*, 71(6):061911, 2005.
- [7] J. F. C. Kingman. The coalescent. *Stochastic Processes and their Applications*, 13(3):235–248, 1982.
- [8] M. Manceau and A. Lambert. The species problem from the modeler’s point of view. *bioRxiv*, 2016. DOI: [10.1101/075580](https://doi.org/10.1101/075580).
- [9] E. Mayr. *Systematics and the Origin of Species from the Viewpoint of a Zoologist*. 1942.
- [10] P. A. P. Moran. Random processes in genetics. *Mathematical Proceedings of the Cambridge Philosophical Society*, 54(1):60–71, 1958.
- [11] E. B. Poulton. What is a species? In *Proceedings of the Entomological Society of London 1903*, pages lxxvii–cxvi, 1904.
- [12] N. Ross. Fundamentals of Stein’s method. *Probability Surveys*, 8:201–293, 2011.
- [13] R. Solé, R. Pastor-Satorras, E. Smith, and T. B. Kepler. A model of large-scale proteome evolution. *Advances in Complex Systems*, 5(1):43–54, 2002.
- [14] C. M. Stein. A bound for the error in the normal approximation to the distribution of a sum of dependent random variables. In *Proceedings of the Sixth Berkeley Symposium on Mathematical Statistics and Probability*, volume 2, pages 583–602. University of California Press, 1972.
- [15] S. Tavaré. Line-of-descent and genealogical processes, and their applications in population genetics models. *Theoretical Population Biology*, 26(2):119–164, 1984.
- [16] A. Vázquez, A. Flammini, A. Maritan, and A. Vespignani. Modeling of protein interaction networks. *ComPlexUs*, 1:38–44, 2003.

A Proofs of Proposition 3.5 and Corollary 3.6

In this section we prove Proposition 3.5 and its corollary.

Proposition 3.5. *Let i, j, k and ℓ be four distinct vertices of G_{n,r_n} . We have*

$$\text{Cov}\left(\mathbb{1}_{\{i \leftrightarrow j\}}, \mathbb{1}_{\{k \leftrightarrow \ell\}}\right) = \frac{2r_n}{(1+r_n)^2(3+r_n)(3+2r_n)}$$

Corollary 3.6. *Let $D_n^{(i)}$ and $D_n^{(j)}$ be the respective degrees of two fixed vertices i and j , and let $|E_n|$ be the number of edges of G_{n,r_n} . We have*

$$\text{Cov}\left(D_n^{(i)}, D_n^{(j)}\right) = \frac{r_n}{(1+r_n)^2} \left(1 + \frac{3(n-2)}{3+2r_n} + \frac{2(n-2)(n-3)}{(3+r_n)(3+2r_n)}\right)$$

and

$$\text{Var}(|E_n|) = \frac{r_n n (n-1)(n^2 + 2r_n^2 + 2nr_n + n + 5r_n + 3)}{2(1+r_n)^2(3+r_n)(3+2r_n)}$$

A.1 Proof of Proposition 3.5

The proof of Proposition 3.5 parallels that of Proposition 3.3, but this time the topology of the genealogy of the pairs of vertices has to be taken into account. Indeed, define

$$S_t = \{a_t(i), a_t(j), a_t(k), a_t(\ell)\}$$

and let $\tau_1 < \tau_2 < \tau_3$ be the times of coalescence in the genealogy of $\{i, j, k, \ell\}$, i.e.

$$\tau_p = \inf\{t \geq 0 : |S_t| = 4 - p\}, \quad p = 1, 2, 3.$$

Write $I_1 = [0, \tau_1[$, $I_2 = [\tau_1, \tau_2[$ and $I_3 = [\tau_2, \tau_3[$. Finally, for $m = 1, 2$, let

$$A_{\{uv\}}^{(m)} = \{a_{\tau_m-}(u) \neq a_{\tau_m-}(v)\} \cap \{a_{\tau_m}(u) = a_{\tau_m}(v)\}$$

be the event that the m -th coalescence in the genealogy of $\{i, j, k, \ell\}$ involved the lineages of u and v (note that the third coalescence is uniquely determined by the first and the second, so we do not need $A_{\{uv\}}^{(3)}$).

On $A_{\{ij\}}^{(1)} \cap A_{\{k\ell\}}^{(2)}$, $\{i \leftrightarrow j, k \leftrightarrow \ell\}$ is equivalent to

$$(P_{\{ij\}}^* \cap I_1) \cup (P_{\{k\ell\}}^* \cap I_1) \cup (P_{\{k\ell\}}^* \cap I_2) = \emptyset$$

so that, conditionally on I_1 and I_2 , by Lemma 2.4,

$$\begin{aligned} \mathbb{P}\left(i \leftrightarrow j, k \leftrightarrow \ell \mid A_{\{ij\}}^{(1)} \cap A_{\{k\ell\}}^{(2)}\right) &= \mathbb{P}\left((P_{\{ij\}}^* \cup P_{\{k\ell\}}^*) \cap I_1 = \emptyset\right) \times \mathbb{P}\left(P_{\{k\ell\}}^* \cap I_2 = \emptyset\right) \\ &= \frac{6}{6+2r_n} \times \frac{3}{3+r_n}. \end{aligned}$$

By contrast, on $A_{\{ij\}}^{(1)} \cap A_{\{ik\}}^{(2)}$, $\{i \leftrightarrow j, k \leftrightarrow \ell\}$ is

$$(P_{\{ij\}}^* \cap I_1) \cup (P_{\{k\ell\}}^* \cap I_1) \cup (P_{\{k\ell\}}^* \cap I_2) \cup (P_{\{k\ell\}}^* \cap I_3) = \emptyset$$

and thus

$$\mathbb{P}(i \leftrightarrow j, k \leftrightarrow \ell \mid A_{\{ij\}}^{(1)} \cap A_{\{ik\}}^{(2)}) = \frac{6}{6+2r_n} \times \frac{3}{3+r_n} \times \frac{1}{1+r_n}.$$

Given a realization of the topology of the genealogy of the form $A_{\{u_1v_1\}}^{(1)} \cap A_{\{u_2v_2\}}^{(2)}$, we can always express $\{i \leftrightarrow j, k \leftrightarrow \ell\}$ as a union of intersections of $P_{\{ij\}}^*$ and $P_{\{k\ell\}}^*$ with I_1, I_2 and I_3 . In total, there are $\binom{4}{2} \times \binom{3}{2} = 18$ possible events $A_{\{u_1v_1\}}^{(1)} \cap A_{\{u_2v_2\}}^{(2)}$, each having probability $1/18$. This enables us to compute $\mathbb{P}(i \leftrightarrow j, k \leftrightarrow \ell)$, but in fact the calculations can be simplified by exploiting symmetries, such as the fact that $\{ij\}$ and $\{k\ell\}$ are interchangeable. In the end, it suffices to consider four cases, as depicted in Figure 6.

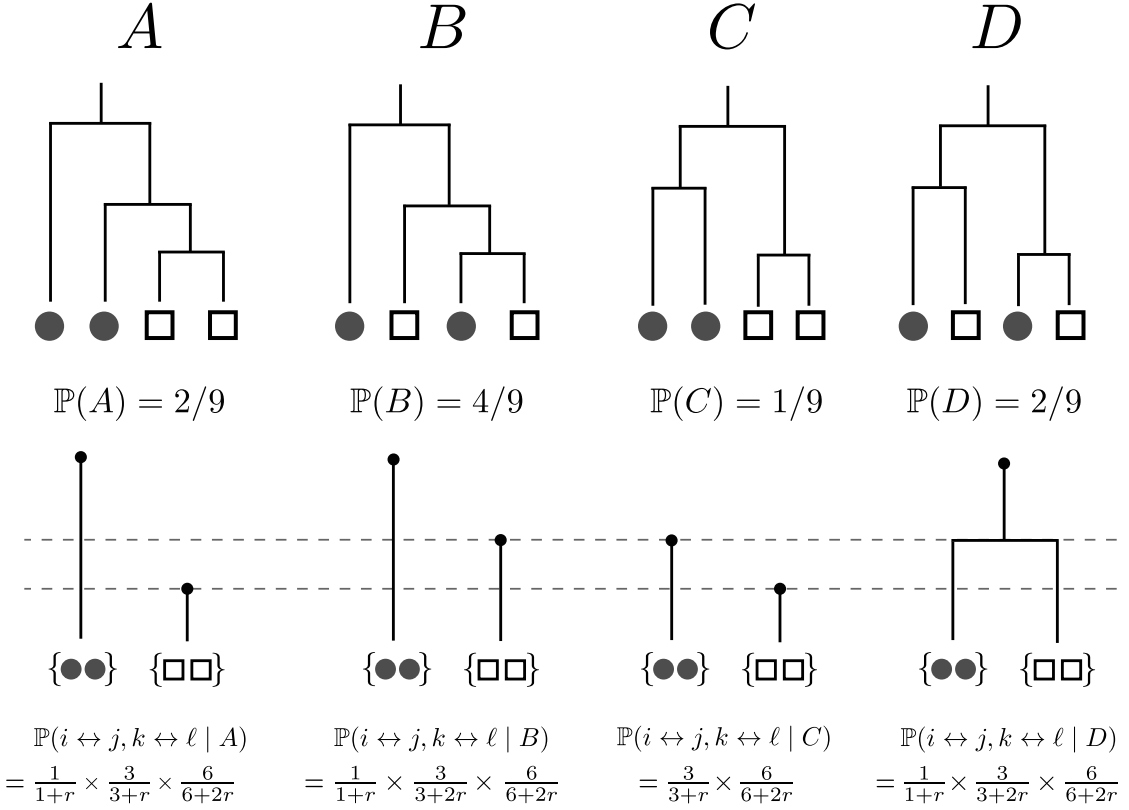


Figure 6: The four cases that we consider to compute $\mathbb{P}(i \leftrightarrow j, k \leftrightarrow \ell)$. Top, the “aggregated” genealogies of vertices and their probability. Each of these correspond to several genealogies on $\{i, j, k, \ell\}$, which are obtained by labeling symbols in such a way that a pair of matching symbols has to correspond to either $\{ij\}$ or $\{k\ell\}$. For instance, $C = (A_{\{ij\}}^{(1)} \cap A_{\{k\ell\}}^{(2)}) \cup (A_{\{k\ell\}}^{(1)} \cap A_{\{ij\}}^{(2)})$ and therefore $\mathbb{P}(C) = 2/18$. Similarly, $A = (A_{\{ij\}}^{(1)} \cap A_{\{ik\}}^{(2)}) \cup (A_{\{ij\}}^{(1)} \cap A_{\{il\}}^{(2)}) \cup (A_{\{k\ell\}}^{(1)} \cap A_{\{ik\}}^{(2)}) \cup (A_{\{k\ell\}}^{(1)} \cap A_{\{jk\}}^{(2)})$ and $\mathbb{P}(A) = 4/18$, etc. Bottom, the associated genealogy of the pairs and the corresponding conditional probability of $\{i \leftrightarrow j, k \leftrightarrow \ell\} \Leftrightarrow \{\square \leftrightarrow \square, \bullet \leftrightarrow \bullet\}$.

Putting the pieces together, we find that

$$\begin{aligned}
\mathbb{P}(i \leftrightarrow k, j \leftrightarrow \ell) &= \frac{6}{9} \times \frac{1}{1+r_n} \times \frac{3}{3+2r} \times \frac{6}{6+2r_n} \\
&+ \frac{2}{9} \times \frac{1}{1+r_n} \times \frac{3}{3+r_n} \times \frac{6}{6+2r_n} \\
&+ \frac{1}{9} \times \frac{3}{3+r_n} \times \frac{6}{6+2r_n} \\
&= \frac{9+2r_n}{(1+r_n)(3+r_n)(3+2r_n)}.
\end{aligned}$$

and Proposition 3.5 follows, since

$$\mathbb{P}(i \leftrightarrow j) \mathbb{P}(k \leftrightarrow \ell) = \left(\frac{1}{1+r_n} \right)^2.$$

A.2 Proof of Corollary 3.6

Corollary 3.6 is proved by standard calculations. First,

$$\begin{aligned}
\text{Cov}(D_n^{(i)}, D_n^{(j)}) &= \text{Cov}\left(\sum_{k \neq i} \mathbb{1}_{\{i \leftrightarrow k\}}, \sum_{\ell \neq j} \mathbb{1}_{\{j \leftrightarrow \ell\}}\right) \\
&= \text{Var}(\mathbb{1}_{\{i \leftrightarrow j\}}) \\
&\quad + 3(n-2) \text{Cov}(\mathbb{1}_{\{i \leftrightarrow k\}}, \mathbb{1}_{\{j \leftrightarrow k\}}) \\
&\quad + (n-2)(n-3) \text{Cov}(\mathbb{1}_{\{i \leftrightarrow k\}}, \mathbb{1}_{\{j \leftrightarrow \ell\}})
\end{aligned}$$

Remembering from Proposition 3.1 that $\text{Var}(\mathbb{1}_{\{i \leftrightarrow j\}}) = r_n/(1+r_n)^2$ and from Proposition 3.3 that $\text{Cov}(\mathbb{1}_{\{i \leftrightarrow k\}}, \mathbb{1}_{\{j \leftrightarrow k\}}) = \frac{r_n}{(1+r_n)^2(3+2r_n)}$, and using Proposition 3.5, we find that

$$\text{Cov}(D_n^{(i)}, D_n^{(j)}) = \frac{r_n}{(1+r_n)^2} \left(1 + \frac{3(n-2)}{3+2r_n} + \frac{2(n-2)(n-3)}{(3+r_n)(3+2r_n)} \right).$$

Finally, to compute $\text{Var}(|E_n|)$, we could do a similar calculation. However, it is easier to note that

$$|E_n| = \frac{1}{2} \sum_{i=1}^n D_n^{(i)}.$$

As a result,

$$\begin{aligned}
\text{Var}(|E_n|) &= \frac{1}{4} \left(n \text{Var}(D_n^{(i)}) + n(n-1) \text{Cov}(D_n^{(i)}, D_n^{(j)}) \right) \\
&= \frac{r_n n (n-1) (n^2 + 2r_n^2 + 2nr_n + n + 5r_n + 3)}{2(1+r_n)^2(3+r_n)(3+2r_n)}.
\end{aligned}$$

B Proof of Theorem 4.2

In this section, we prove Theorem 4.2.

Theorem 4.2 (convergence of the rescaled degree).

- (i) *If $r_n \rightarrow r > 0$, then $\frac{D_n}{n}$ converges in distribution to a Beta(2, 2r) random variable.*
- (ii) *If r_n is both $\omega(1)$ and $o(n)$, then $\frac{D_n}{n/r_n}$ converges in distribution to a size-biased exponential variable with parameter 2.*
- (iii) *If $2r_n/n \rightarrow \rho > 0$, then $D_n + 1$ converges in distribution to a size-biased geometric variable with parameter $\rho/(1 + \rho)$.*

The proof of (iii) is immediate: indeed, by Theorem 4.1,

$$\mathbb{P}(D_n + 1 = k) = \frac{2r_n(2r_n + 1)}{(n + 2r_n)(n - 1 + 2r_n)} k \prod_{i=1}^{k-1} \frac{n - i}{n - i + 2r_n - 1}.$$

If $2r_n/n \rightarrow \rho$, then for any fixed k this goes to $k \left(\frac{\rho}{1+\rho}\right)^2 \left(\frac{1}{1+\rho}\right)^{k-1}$ as $n \rightarrow +\infty$.

B.1 Outline of the proof

To prove (i) and (ii), we show the pointwise convergence of the cumulative distribution function F_n of the rescaled degree. To do so, in both cases,

1. We show that, for any $\varepsilon > 0$, for n large enough,

$$\forall y \geq 0, \quad \int_0^y f_n(x) dx \leq F_n(y) \leq \int_0^{y+\varepsilon} f_n(x) dx$$

for some function f_n to be introduced later.

2. We identify the limit of f_n as a classical probability density f , and use dominated convergence to conclude that

$$\forall y \geq 0, \quad \int_0^y f_n(x) dx \rightarrow \int_0^y f(x) dx.$$

In order to factorize as much of the reasoning as possible, we introduce the rescaling factor N_n :

- When $r_n \rightarrow r$, i.e. when we want to prove (i), $N_n = n$.
- When r_n is both $\omega(1)$ and $o(n)$, i.e. when we want to prove (ii), $N_n = n/r_n$.

Thus, in both cases the rescaled degree is D_n/N_n and its cumulative distribution function is

$$F_n(y) = \sum_{k=0}^{\lfloor N_n y \rfloor} \mathbb{P}(D_n = k).$$

B.2 Step 1

For all $x > 0$, let

$$f_n(x) = N_n \mathbb{P}(D_n = \lfloor N_n x \rfloor),$$

so that

$$\forall k \in \mathbb{N}, \quad \mathbb{P}(D_n = k) = \int_{k/N_n}^{(k+1)/N_n} f_n(x) dx.$$

It follows that

$$F_n(y) = \int_0^{(\lfloor N_n y \rfloor + 1)/N_n} f_n(x) dx.$$

Finally, since $y \leq \frac{\lfloor N_n y \rfloor + 1}{N_n} \leq y + \frac{1}{N_n}$ and f_n is non-negative, for any $\varepsilon > 0$, for n large enough,

$$\forall y \geq 0, \quad \int_0^y f_n(x) dx \leq F_n(y) \leq \int_0^{y+\varepsilon} f_n(x) dx,$$

and the rank after which these inequalities hold is uniform in y , because the convergence of $(\lfloor N_n y \rfloor + 1)/N_n$ to y is.

B.3 Step 2

To identify the limit of f_n , we reexpress it in terms of the gamma function. Using that $\Gamma(z) = z\Gamma(z)$, by induction,

$$\prod_{i=1}^k (n-i) = \frac{\Gamma(n)}{\Gamma(n-k)} \quad \text{and} \quad \prod_{i=1}^k (n-i+2r_n-1) = \frac{\Gamma(n+2r_n-1)}{\Gamma(n-k+2r_n-1)}.$$

Therefore, $f_n(x)$ can also be written

$$f_n(x) = \frac{N_n 2r_n (2r_n + 1)}{(n + 2r_n)(n - 1 + 2r_n)} \left(\lfloor N_n x \rfloor + 1 \right) \times P_n(x), \quad (6)$$

where

$$P_n(x) = \frac{\Gamma(n) \Gamma(n - \lfloor N_n x \rfloor + 2r_n - 1)}{\Gamma(n - \lfloor N_n x \rfloor) \Gamma(n + 2r_n - 1)}. \quad (7)$$

We now turn to the specificities of the proofs of (i) and (ii).

B.3.1 Proof of (i)

In this subsection, $r_n \rightarrow r > 0$ and $N_n = n$.

Limit of f_n . Recall that

$$\forall \alpha \in \mathbb{R}, \quad \frac{\Gamma(n + \alpha)}{\Gamma(n)} \sim n^\alpha.$$

Using this in (7), we see that, for all $x \in [0, 1[$,

$$P_n(x) \rightarrow (1 - x)^{2r-1}.$$

Therefore, for all $x \in [0, 1[$,

$$f_n(x) \rightarrow 2r(2r + 1)x(1 - x)^{2r-1}.$$

Noting that $2r(2r + 1) = 1/B(2, 2r)$, where B denotes the beta function, we can write $f = \lim_n f_n$ as

$$f: x \mapsto \frac{x(1 - x)^{2r-1}}{B(2, 2r)} \mathbb{1}_{[0, 1[}(x)$$

and we recognize the probability density function of the Beta(2, 2r) distribution.

Domination of (f_n) . First note that, for all $x \in [0, 1[$,

$$\frac{1}{n - 1 + 2r_n} \prod_{i=1}^{\lfloor nx \rfloor} \frac{n - i}{n - i + 2r_n - 1} = \frac{1}{n - \lfloor nx \rfloor + 2r_n - 1} \prod_{i=1}^{\lfloor nx \rfloor} \frac{n - i}{n - i + 2r_n},$$

where the empty product is understood to be 1. Since $2r_n > 0$, this enables us to write that, for all $x \in [0, 1[$,

$$f_n(x) = \underbrace{\frac{n 2r (2r + 1)}{n + 2r}}_{\leq (2r+1)^2} \times \frac{\lfloor nx \rfloor + 1}{n - 1 + 2r} \times \underbrace{\frac{1}{n - \lfloor nx \rfloor + 2r - 1}}_{\leq \frac{1}{2r}} \times \underbrace{\prod_{i=1}^{\lfloor nx \rfloor} \frac{n - i}{n - i + 2r}}_{\leq 1}.$$

where, to avoid cluttering the expression, the n index of r_n has been dropped. Since

$$\frac{\lfloor nx \rfloor + 1}{n - 1 + 2r_n} \leq \frac{(n - 1)x + x + 1}{n - 1} \leq x + \frac{2}{n - 1} \xrightarrow[n \rightarrow +\infty]{\text{uniformly}} x,$$

there exists c such that, for all $x \in [0, 1[$ and n large enough,

$$f_n(x) \leq cx$$

Since f_n is zero outside of $[0, 1[$, this shows that (f_n) is dominated by $g: x \mapsto cx \mathbb{1}_{[0, 1[}(x)$.

B.3.2 Proof of (ii)

In this subsection, r_n is both $\omega(1)$ and $o(n)$, and $N_n = n/r_n$. For brevity, we will write k_n for $\lfloor nx/r_n \rfloor$. It should be noted that

- k_n is both $\omega(1)$ and $o(n)$.
- $k_n r_n / n \rightarrow x$ uniformly in x on $[0, +\infty[$.

Limit of f_n . In this paragraph, we will need Stirling's formula for the asymptotics of Γ :

$$\Gamma(t+1) \sim \sqrt{2\pi t} \frac{t^t}{e^t}.$$

Using this in Equation (7),

$$\begin{aligned} P_n(x) &= \frac{\Gamma(n) \Gamma(n - \lfloor N_n x \rfloor + 2r_n - 1)}{\Gamma(n - \lfloor N_n x \rfloor) \Gamma(n + 2r_n - 1)} \\ &\sim \underbrace{\sqrt{\frac{(n-1)(n-2-k_n+2r_n)}{(n-1-k_n)(n-2+2r_n)}}}_{\sim 1} \times \underbrace{\frac{e^{n-1-k_n} e^{n-2+2r_n}}{e^{n-1} e^{n-2-k_n+2r_n}}}_{=1} \times Q_n \end{aligned}$$

where

$$Q_n = \frac{(n-1)^{n-1} (n-2-k_n+2r_n)^{n-2-k_n+2r_n}}{(n-1-k_n)^{n-1-k_n} (n-2+2r_n)^{n-2+2r_n}}.$$

Let us show that $Q_n \rightarrow e^{-2x}$:

$$\begin{aligned} \log Q_n &= (n-1) \log(n-1) \\ &\quad + (n-a+b) \log(n-a+b) \\ &\quad - (n-a) \log(n-a) \\ &\quad - (n-1+b) \log(n-1+b) \end{aligned}$$

where, to avoid cluttering the text, we have written a for $k_n + 1$ and b for $2r_n - 1$. Factorizing, we get

$$\log Q_n = n \log \left(\frac{(n-1)(n-a+b)}{(n-a)(n-1+b)} \right) - a \log \left(\frac{n-a+b}{n-a} \right) + b \log \left(\frac{n-a+b}{n-1+b} \right) - \log \left(\frac{n-1}{n-1+b} \right).$$

Now,

$$\frac{(n-1)(n-a+b)}{(n-a)(n-1+b)} = 1 + \frac{(a-1)b}{\underbrace{n^2 - n + nb - na + a - ab}_{\sim \frac{2k_n r_n}{n^2} = o(1)}}$$

so that

$$n \log \left(\frac{(n-1)(n-a+b)}{(n-a)(n-1+b)} \right) \sim \frac{2k_n r_n}{n} \rightarrow 2x$$

Similarly,

$$-a \log\left(\frac{n-a+b}{n-a}\right) = -a \log\left(1 + \frac{b}{n-a}\right) \sim -\frac{ab}{n} \rightarrow -2x$$

$$b \log\left(\frac{n-a+b}{n-1+b}\right) = b \log\left(1 + \frac{1-a}{n-1+b}\right) \sim -\frac{ab}{n} \rightarrow -2x$$

and, finally, $-\log\left(\frac{n-1}{n-1+b}\right) \rightarrow 0$. Putting the pieces together,

$$\log Q_n \rightarrow -2x.$$

Having done that, we note that

$$\frac{2n(2r_n+1)}{(n+2r_n)(n-1+2r_n)}(k_n+1) \rightarrow 4x.$$

Plugging these results in Equation (6), we see that

$$\forall x \in \mathbb{R}, \quad f_n(x) \rightarrow 4x e^{-2x} \mathbf{1}_{[0,+\infty[}(x)$$

and we recognize the probability density function of a size-biased exponential distribution with parameter 2.

Domination of (f_n) . Recall that, since $N_n = n/r_n$, for all $x \in [0, 1[$,

$$f_n(x) = \frac{2n(2r_n+1)}{(n+2r_n)(n-1+2r_n)}(k_n+1) \prod_{i=1}^{k_n} \frac{n-i}{n-i+2r_n-1}.$$

Next, note that, for all i ,

$$\frac{n-i}{n-i+2r_n-1} = 1 - \frac{2r_n-1}{n-i+2r_n-1} \leq \exp\left(-\frac{2r_n-1}{n-i+2r_n-1}\right)$$

so that

$$\prod_{i=1}^{k_n} \frac{n-i}{n-i+2r_n-1} \leq \exp\left(-\sum_{i=1}^{k_n} \frac{2r_n-1}{n-i+2r_n-1}\right),$$

with

$$\sum_{i=1}^{k_n} \frac{2r_n-1}{n-i+2r_n-1} \geq k_n \frac{2r_n-1}{n-1+2r_n-1}.$$

Because $r_n = \omega(1)$, for all $\varepsilon > 0$, $2r_n-1 \geq (1-\varepsilon)2r_n$ for n large enough. Similarly, since $r_n = o(n)$, $\frac{1}{n+2r_n} \geq \frac{1}{(1+\varepsilon)n}$. As a result, there exists $c > 0$ such that

$$k_n \frac{2r_n-1}{n-1+2r_n-1} \geq c k_n \frac{2r_n}{n} \xrightarrow{\text{uniformly}} 2cx.$$

We conclude that

$$\forall x \geq 0, \quad \prod_{i=1}^{k_n} \frac{n-i}{n-i+2r_n-1} \leq \exp(-2cx)$$

for n large enough. Finally,

$$2 \times \underbrace{\frac{n}{n+2r_n}}_{\leq 1} \times \underbrace{\frac{(2r_n+1)(k_n+1)}{(n-1+2r_n)}}_{\rightarrow 2x, \text{ uniformly}} \leq 4cx$$

and so (f_n) is dominated by $g: x \mapsto 4cx e^{-2cx} \mathbb{1}_{[0,+\infty[}(x)$.