

ASYMPTOTIC PROPERTIES OF A RANDOM GRAPH WITH DUPLICATIONS

ÁGNES BACKHAUSZ AND TAMÁS F. MÓRI

DEPARTMENT OF PROBABILITY THEORY AND STATISTICS,
EÖTVÖS LORÁND UNIVERSITY
Pázmány P. s. 1/C, H-1117 Budapest, Hungary
E-mail address: agnes@math.elte.hu, mori@math.elte.hu

ABSTRACT. We deal with a random graph model evolving in discrete time steps by duplicating and deleting the edges of randomly chosen vertices. We prove the existence of an a.s. asymptotic degree distribution, with stretched exponential decay; more precisely, the proportion of vertices of degree d tends to some positive number $c_d > 0$ almost surely as the number of steps goes to infinity, and $c_d \sim (e\pi)^{1/2}d^{1/4}e^{-2\sqrt{d}}$ holds as $d \rightarrow \infty$.

Keywords: scale free, duplication, deletion, random graphs, martingales.

1. INTRODUCTION

In the last decades, inspired by the examination of large real networks, various types of random graph models with preferential attachment dynamics (meaning that vertices with larger degree have larger chance to get new edges as the graph evolves randomly) were introduced and analysed. After some early work [20, 17, 19], this started with the seminal papers of Barabási and Albert [2, 1999] and Bollobás, Riordan, Spencer and Tusnády [5, 2001]. Among many others, we may mention the model of Cooper and Frieze for the Internet [3, 2003] or that of Sridharan, Yong Gao, Kui Wu and Nastos [18, 2011] for social networks.

An important feature of these graph sequences is the scale-free property: the proportion of vertices of degree d tends to some positive number c_d almost surely as the number of steps goes to infinity, and $c_d \sim Kd^{-\gamma}$ holds as $d \rightarrow \infty$ (throughout this paper, $a_d \sim b_d$ means that $a_d/b_d \rightarrow 1$ as $d \rightarrow \infty$). To put it in another way, the asymptotic degree distribution (c_d) is polynomially decaying. See also [2, 10, 23] and the references therein about the scale-free property of the internet.

Date: 30 August 2013.

2010 *Mathematics Subject Classification.* 60G42, 05C80.

Supported by the Hungarian Scientific Research Fund – OTKA K 108615.

However, scale-free property captures only the behavior of the degrees of vertices, and does not examine other kinds of structures. For example, especially in biological networks, e.g., proteomes, it happens that we can find groups of vertices having a similar neighborhood, that is, most of their neighbors are the same. One can say that these networks are highly clustered; loosely speaking, there are large cliques, in which almost every vertex is connected to almost every other one, and there are only a few edges going between cliques.

A simple way to generate cliques is duplication: when a new vertex is added, we choose an old vertex randomly, and connect the new vertex to the neighbors of the old one. In other words, the new vertex becomes a copy of the old vertex. Note that if the old vertex is chosen uniformly at random, then the probability that a vertex of degree d gets a new edge is just the probability that one of its neighbors is chosen, which is proportional to its actual degree. Hence this model is also driven by a kind of preferential attachment dynamics.

After the duplication, we can add some extra edges randomly, or we can delete some of them to guarantee that the network remains sparse. The graph may still have some large cliques due to the duplication.

Duplication is not only a technical step that proved to be useful: it is inherent. “This may be because duplication of the information in the genome is a dominant evolutionary force in shaping biological networks (like gene regulatory networks and protein–protein interaction networks)” [6].

These kinds of models – where the duplicated vertex is chosen uniformly at random – were examined for example by Kim, Krapivsky, Kahng, and Redner [13, 2002]. In their model the new vertex is connected to each neighbor of the chosen one with probability $1 - \delta$, independently. In addition, the new vertex is connected to each old one independently with probability β/n at the n th step (δ, β are the parameters of the model). Scale free property is claimed for this model. However, Pastor-Satorras, Smith and Solé [14, 2003] stated that, instead of polynomial decay, for the limit c_d of the expected value of the proportion of vertices of degree d we have $c_d \sim Kd^{-\gamma}e^{-\lambda d}$ with some positive constants K, γ, λ ; that is, the degree distribution has a polynomial decay with exponential cut-off. On the contrary, Chung, Lu, Dewey, and Galas [6, 2003] claimed that for $\beta = 0$, when we do not have any extra edges, the asymptotic degree distribution exists, and (c_d) is decaying polynomially. None of these papers contained a mathematically rigorous proof.

Bebek, Berenbrink, Cooper, Friedetzky, Nadeau and Sahinalp [4, 2006] disclaimed the above mentioned results of [14] and [6]. In the latter case, they showed that the fraction of isolated vertices (that have no edges) increases with time in the pure duplication model, where

$\beta = 0$. They modified the model to avoid singletons by adding a fixed number of edges to the new vertex, chosen uniformly at random. They assumed without any proof that the asymptotic degree distribution exists, and they claimed that it is decaying polynomially.

Hamdi, Krishnamurthy and Yin [11, 2013+] present a model where the probabilities of adding a duplicated edge depends on the state of a hidden Markov chain. Polynomial decay is stated for the limit of the mean of degree distribution. We also mention the somewhat different model of Jordan [12, 2011], and the duplication model of Cohen, Jordan and Voliotis [7, 2010], where the duplicated vertex is chosen not uniformly, but with probabilities proportional to the actual degrees.

In our paper, we present a simple random graph model based on the duplication of a vertex chosen uniformly at random, and the erasure of the edges of another vertex also chosen uniformly at random. We prove that for all d , the proportion of vertices of degree d tends to some c_d with probability 1 as the number of steps goes to infinity. Here c_d is a positive number; we will formulate it as an integral, and then we will determine the asymptotics of the sequence (c_d) as $d \rightarrow \infty$, showing that it has a stretched exponential decay. Hence this model does not have the scale free property. We use methods of martingale theory for proving almost sure convergence, and generating function and Taylor series techniques for deriving the integral representation and the asymptotics of the sequence (c_d) .

2. DEFINITION OF THE MODEL AND MAIN RESULTS

Our model has two different versions. Both of them start with a single vertex. The graph evolves in discrete time steps; each step has a duplication and an erasure part. At each step a new vertex will be born; therefore the number of vertices after n steps is $n + 1$. The graph is always a simple graph; it has neither multiple edges nor loops. At each step we do the following.

Version 1. We choose two (not necessarily different) old vertices independently, uniformly at random. Then the new vertex is added to the graph; we connect it to the first vertex and to all its neighbors. After that we delete all edges emanating from the second old vertex we have selected, with the possible exception that edges of the new vertex cannot be deleted.

Version 2. We choose two (not necessarily different) old vertices independently, uniformly at random. The new vertex is connected to the first one and to its neighbors. Then we delete all edges of the second vertex without any exceptions.

That is, the new edges are protected in the erasure part of the same step in version 1, but they might be deleted immediately in version 2.

We will see that the version 2 graph has a simple structure that enables us to describe its asymptotical degree distribution. Then, using this and a coupling of the two models, we can prove similar results for version 1.

Let us remark that the presence of deletion makes the analysis more difficult than in the usual recursive graph models, since it causes intensive fluctuation in the model's behavior.

Our model is a kind of coagulation–fragmentation one: the effect of duplication is coagulation, and deletion results fragmentation. Coagulation–fragmentation models are frequently used in several areas, see e.g. [8]. Ráth and Tóth applied these models for random graph models [15], namely, for the Erdős–Rényi model, which is completely different from ours.

The basic property of version 2 is that the evolving graph always consists of separated complete graphs. That is, it is a disjoint union of cliques. Within a component, every pair of vertices is connected, and there are no edges between the components. Indeed, we start from a single vertex, which is a clique of size one, and both duplication and erasure make cliques from cliques. Moreover, it is easy to see that if we start the model with an arbitrary graph, all edges of the initial configuration are deleted after a while, and after that the graph will consist of separated cliques. So the initial configuration does not make any difference asymptotically.

We may formulate the second version as follows. At each step we choose two components independently such that the probability that a given clique is chosen is proportional to its size. The new vertex is attached to the first clique, so its size is increased by 1; the size of the secondly chosen clique is decreased by 1, and an isolated vertex (the deleted one) comes into existence. Note that if we choose an isolated vertex to be deleted, then it remains isolated.

This structure of version 2 makes it easier to handle, as the number of d -cliques does not vary so vehemently as the number of degree d vertices; the fluctuation is bounded by 2. This will lead to the description of the asymptotic degree distribution of version 1 in an almost sure sense. Our main results are the following.

Theorem 1. *Denote by $X[n, d]$ the number of vertices of degree d after n steps in version 1. Then*

$$\frac{X[n, d]}{n+1} \rightarrow c_d$$

holds almost surely as $n \rightarrow \infty$, where (c_d) is a sequence of positive numbers satisfying

$$(1) \quad c_0 = \frac{1+c_1}{3}; \quad c_d = \frac{d+1}{2d+3}(c_{d-1} + c_{d+1}) \quad (d \geq 2).$$

For the asymptotic analysis we first present an integral representation for the limiting sequence (c_d) . As a corollary, we get that the sum of this sequence is 1; it is really a probability distribution.

Theorem 2. *For the sequence (c_d) of Theorem 1 we have*

$$c_d = (d + 1) \int_0^\infty \frac{y^d e^{-y}}{(1 + y)^{d+2}} dy \quad (d \geq 0),$$

and $\sum_{d=0}^\infty c_d = 1$.

Using this formula we can derive the asymptotics of c_d .

Theorem 3. *For the sequence (c_d) of Theorem 1 we have*

$$c_d \sim (e\pi)^{1/2} d^{1/4} e^{-2\sqrt{d}}, \quad \text{as } d \rightarrow \infty.$$

Our model is invented to ensure high degree clustering. Finally, let us quantify this property.

The local clustering coefficient of a vertex of degree d is defined to be the fraction of connections that exist between the $\binom{d}{2}$ pairs of neighbors (meant 0 when $d < 2$). Watts and Strogatz [22] define the clustering coefficient of the whole graph as the average of the local clustering coefficients of all the vertices. Let us call this quantity the average clustering coefficient. Another possibility for a such a measure is the ratio of 3 times the number of triangles divided by the number of connected triplets (paths of length 2), see [21]. This version is sometimes called transitivity; we will refer to it as the global clustering coefficient.

Since the graph in version 2 is consists of disjoint cliques, its global clustering coefficient is obviously 1, while the average clustering coefficient is equal to the proportion of vertices with degree at least 2. By Theorem 1 it converges to $1 - c_0 - c_1 = 2 - 4c_0$ almost surely, as $n \rightarrow \infty$. We note that the limit is equal to $0.38538\dots$ by Theorem 2. These results can be transferred to version 1.

Theorem 4. *In version 1, the global clustering coefficient converges to 1, and the average clustering coefficient to $1 - c_0 - c_1$, almost surely, as $n \rightarrow \infty$.*

The high clustering property of our model shows that is is a so-called small-world graph [22].

3. PROOFS

Preliminaries. First we formulate the lemma from martingale theory that we will use several times and whose proof can be found in [1].

Lemma 1. *Let (\mathcal{F}_n) be a filtration, (ξ_n) a nonnegative adapted process. Suppose that*

$$(2) \quad E((\xi_n - \xi_{n-1})^2 \mid \mathcal{F}_{n-1}) = O(n^{1-\delta})$$

holds with some $\delta > 0$. Let (u_n) , (v_n) be nonnegative predictable processes such that $u_n < n$ for all $n \geq 1$. Finally, let (w_n) be a regularly varying sequence of positive numbers with exponent $\mu \geq 0$.

(a) Suppose that

$$E(\xi_n \mid \mathcal{F}_{n-1}) \leq \left(1 - \frac{u_n}{n}\right) \xi_{n-1} + v_n,$$

and $\lim_{n \rightarrow \infty} u_n = u$, $\limsup_{n \rightarrow \infty} v_n/w_n \leq v$ with some random variables $u > 0$, $v \geq 0$. Then

$$\limsup_{n \rightarrow \infty} \frac{\xi_n}{nw_n} \leq \frac{v}{u + \mu + 1} \quad a.s.$$

(b) Suppose that

$$E(\xi_n \mid \mathcal{F}_{n-1}) \geq \left(1 - \frac{u_n}{n}\right) \xi_{n-1} + v_n,$$

and $\lim_{n \rightarrow \infty} u_n = u$, $\liminf_{n \rightarrow \infty} v_n/w_n \geq v$ with some random variables $u > 0$, $v \geq 0$. Then

$$\liminf_{n \rightarrow \infty} \frac{\xi_n}{nw_n} \geq \frac{v}{u + \mu + 1} \quad a.s.$$

Asymptotic degree distribution in version 2. Recall that in this case the graph is always disjoint union of complete graphs.

First we prove the following analogue of Theorem 1.

Proposition 5. Denote by $Y[n, k]$ the number of cliques of size k after n steps in version 2. Then for all positive integers k we have

$$\frac{Y[n, k]}{n} \rightarrow y_k \quad \text{almost surely as } n \rightarrow \infty,$$

where (y_k) is a sequence of positive numbers satisfying

$$(3) \quad y_1 = \frac{1 + 2y_2}{3}, \quad y_k = \frac{(k-1)y_{k-1} + (k+1)y_{k+1}}{2k+1} \quad (k \geq 2).$$

Note that (3) (as well as equation (1)) is not a recursion. This prevents us proceeding simply in the usual, direct way, with induction over k .

Proof. For $n = 0$ we have $Y[0, 1] = 1$, all the other ones are equal to zero. The total number of vertices is n after $n - 1$ steps. Let \mathcal{F}_n denote the σ -field generated by the first n steps.

We enumerate the events that can happen to the cliques of different sizes at a step.

At the n th step an isolated vertex may become

- a clique of size 2 (increased but not decreased) with probability $\frac{1}{n}(1 - \frac{1}{n})$;
- an isolated vertex (any other cases).

A clique of size $k \geq 2$ may become a clique of size

- $k - 1$ (not increased but decreased) with probability $\frac{k}{n}\left(1 - \frac{k}{n}\right)$;
- $k + 1$ (increased but not decreased) with probability $\frac{k}{n}\left(1 - \frac{k}{n}\right)$;
- k (any other cases).

The deleted vertex will be a new isolated one unless one of them is chosen for erasure but not for duplication, which has probability $\frac{1}{n}\left(1 - \frac{1}{n}\right)$ for each of them.

Putting this together with the fact that the random choices are independent and probabilities are proportional to clique sizes, we can compute the conditional expectation of $Y[n, k]$ with respect to \mathcal{F}_{n-1} , which is the σ -field generated by the first $n - 1$ steps.

$$\begin{aligned} E(Y[n, 1]|\mathcal{F}_{n-1}) &= Y[n-1, 1] \left[1 - \frac{1}{n} \left(1 - \frac{1}{n} \right) - \frac{1}{n} \left(1 - \frac{1}{n} \right) \right] \\ &\quad + 1 + Y[n-1, 2] \cdot \frac{2}{n} \left(1 - \frac{2}{n} \right); \\ E(Y[n, k]|\mathcal{F}_{n-1}) &= Y[n-1, k] \left[1 - 2 \cdot \frac{k}{n} \left(1 - \frac{k}{n} \right) \right] \\ &\quad + Y[n-1, k-1] \cdot \frac{k-1}{n} \left(1 - \frac{k-1}{n} \right) \\ &\quad + Y[n-1, k+1] \cdot \frac{k+1}{n} \left(1 - \frac{k+1}{n} \right) \quad (k \geq 2). \end{aligned}$$

Let $A_k = \liminf_{n \rightarrow \infty} \frac{Y[n, k]}{n}$ and $B_k = \limsup_{n \rightarrow \infty} \frac{Y[n, k]}{n}$ for $k \geq 1$. It is clear that $0 \leq A_k \leq B_k \leq 1$ holds for these random variables.

We will give a sequence of lower bounds for (A_k) , and similarly, a sequence of upper bounds for (B_k) ; then we will show that their limits are equal to each other. First, let $a_k^{(0)} = 0$ for $k \geq 1$. Having constructed the sequence $(a_k^{(j)})_{k \geq 1}$, we define

$$(4) \quad a_1^{(j+1)} = \frac{1 + 2a_2^{(j)}}{3}, \quad a_k^{(j+1)} = \frac{(k-1)a_{k-1}^{(j)} + (k+1)a_{k+1}^{(j)}}{2k+1} \quad (k \geq 2).$$

We get $a_k^{(j)}$ recursively for every $k \geq 1$ and $j \geq 1$.

We prove by induction on j that $a_k^{(j)} \leq A_k$ ($k \geq 1$). Since $Y[n, k] \geq 0$, this is clear for $j = 0$. Suppose that this is satisfied for some j for every k . For $k = 1$ we apply Lemma 1 with

$$\xi_n = Y[n, 1], \quad u_n = 2 - \frac{2}{n} \rightarrow 2, \quad v_n = 1 + Y[n-1, 2] \cdot \frac{2}{n} \left(1 - \frac{2}{n} \right).$$

Now (ξ_n) is nonnegative adapted. (u_n) and (v_n) are clearly nonnegative predictable sequences; we can choose $w_n = 1$, $\mu = 0$, $u = 2 > 0$

and finally, $v = 1 + 2a_2^{(j)} \geq 0$ due to the induction hypothesis. Note that at each step at most one of the isolated points vanishes and at most two may appear. Thus (2) is clearly satisfied. Lemma 1 implies that

$$A_1 = \liminf_{n \rightarrow \infty} \frac{Y[n, 1]}{n} = \liminf_{n \rightarrow \infty} \frac{\xi_n}{n} \geq \frac{v}{u+1} = \frac{1 + 2a_2^{(j)}}{3} = a_1^{(j+1)}$$

almost surely.

Similarly, for $k \geq 2$, if we have $A_k \geq a_k^{(j)}$ for some $j \geq 1$, we can choose

$$\begin{aligned} \xi_n &= Y[n, k], \quad u_n = 2k - \frac{2k^2}{n} \rightarrow 2k, \\ v_n &= Y[n-1, k-1] \cdot \frac{k-1}{n} \left(1 - \frac{k-1}{n}\right) \\ &\quad + Y[n-1, k+1] \cdot \frac{k+1}{n} \left(1 - \frac{k+1}{n}\right), \\ v &= (k-1)a_{k-1}^{(j)} + (k+1)a_{k+1}^{(j)}. \end{aligned}$$

At each step at most three cliques are changed, which implies that (2) holds. Thus in this case from Lemma 1 we obtain that

$$A_k = \liminf_{n \rightarrow \infty} \frac{Y[n, k]}{n} \geq \frac{v}{u+1} = \frac{(k-1)a_{k-1}^{(j)} + (k+1)a_{k+1}^{(j)}}{2k+1} = a_k^{(j+1)}$$

almost surely.

By induction on j we get that $A_k \geq a_k^{(j)}$ holds almost surely for $k \geq 1$ and $j \geq 0$.

Now we verify that for fixed k the sequence $(a_k^{(j)})$ is monotone increasing in j . Since $a_k^{(0)} = 0$ for every k , from equations (4) it is clear that $a_k^{(1)} \geq a_k^{(0)}$. Suppose that for some $j \geq 1$ we have $a_k^{(j)} \geq a_k^{(j-1)}$ for every k . Then

$$\begin{aligned} a_1^{(j+1)} &= \frac{1 + 2a_2^{(j)}}{3} \geq \frac{1 + 2a_2^{(j-1)}}{3} = a_1^{(j)}; \\ a_k^{(j+1)} &= \frac{(k-1)a_{k-1}^{(j)} + (k+1)a_{k+1}^{(j)}}{2k+1} \\ &\geq \frac{(k-1)a_{k-1}^{(j-1)} + (k+1)a_{k+1}^{(j-1)}}{2k+1} = a_k^{(j)} \end{aligned}$$

follows from equations (4). Thus by induction on j we get that $a_k^{(j)} \geq a_k^{(j-1)}$ for $k, j \geq 1$.

It is clear that the sequence $(a_k^{(j)})_{j \geq 0}$ is uniformly bounded from above by 1. Using monotonicity we can define

$$a_k = \lim_{j \rightarrow \infty} a_k^{(j)} \quad (k \geq 1).$$

From equation (4) it follows that (a_k) satisfies (3), that is,

$$a_1 = \frac{1 + 2a_2}{3}, \quad a_k = \frac{(k-1)a_{k-1} + (k+1)a_{k+1}}{2k+1} \quad (k \geq 2).$$

On the other hand, since $A_k \geq a_k^{(j)}$ for $k \geq 1$ and $j \geq 0$, we have $A_k \geq a_k$ almost surely.

Similarly, we define $b_k^{(0)} = 1$ for every k , and then

$$b_1^{(j+1)} = \frac{1 + 2b_2^{(j)}}{3}, \quad b_k^{(j+1)} = \frac{(k-1)b_{k-1}^{(j)} + (k+1)b_{k+1}^{(j)}}{2k+1} \quad (k \geq 2).$$

Using part (a) of Lemma 1 it follows by induction on j that $B_k \leq b_k^{(j)}$ holds almost surely.

In this case, for fixed k the sequence $b_k^{(j)}$ is decreasing, and for the limits $b_k = \lim_{j \rightarrow \infty} b_k^{(j)}$ we also have

$$b_1 = \frac{1 + 2b_2}{3}, \quad b_k = \frac{(k-1)b_{k-1} + (k+1)b_{k+1}}{2k+1} \quad (k \geq 2).$$

In addition, $B_k \leq b_k$ almost surely.

By definition, $0 \leq A_k \leq B_k \leq 1$ and $0 \leq a_k \leq b_k \leq 1$ hold. Let $d_k = b_k - a_k \geq 0$ for all k . We have the same equations for (a_k) and (b_k) . This yields

$$d_1 = \frac{2d_2}{3}, \quad d_k = \frac{(k-1)d_{k-1} + (k+1)d_{k+1}}{2k+1} \quad (k \geq 2).$$

By rearranging we get that

$$(5) \quad d_2 = \frac{3}{2}d_1, \quad d_{k+1} = \frac{(2k+1)d_k - (k-1)d_{k-1}}{k+1} \quad (k \geq 2).$$

Suppose that $d_k \geq \frac{k+1}{k}d_{k-1}$ holds for some $k \geq 2$. (For $k = 2$ this is true with equality.) Since d_{k-1} is nonnegative, $d_k \geq d_{k-1}$ also follows from this assumption. We obtain from equation (5) that

$$d_{k+1} \geq \frac{(k+2)d_k}{k+1}.$$

Therefore this inequality holds for every k .

This implies that $d_k \geq (k+1)d_1$ for every k . Since $0 \leq d_k = b_k - a_k \leq 1$, it follows that $d_1 = 0$.

From (5) we obtain that $d_k = 0$ for all k , which implies that $a_k = b_k$. Since these were the lower and upper bounds for the limit inferior and limit superior of $\frac{Y[n,k]}{n}$, we get that the latter must converge almost surely as $n \rightarrow \infty$, and the limits satisfy (3). \square

Corollary 6. *In version 2, the proportion of vertices of degree d tends to c_d satisfying (1) almost surely as $n \rightarrow \infty$.*

Proof. For a fixed d we have $d + 1$ vertices of degree d in each clique of size $k = d + 1$. Therefore for the proportion of vertices of degree d tends to $(d + 1)y_{d+1}$ by Proposition 5. From equations (3) we obtain that

$$c_0 = y_1 = \frac{1 + 2y_2}{3} = \frac{1 + c_1}{3};$$

$$c_d = (d + 1)y_{d+1} = \frac{d + 1}{2d + 3}(c_{d-1} + c_{d+1}) \quad (d \geq 2).$$

□

Asymptotic degree distribution in version 1. When proving the results for version 2 we essentially used the property that the graphs consists of disjoint union of cliques: at most three of the cliques may change at a step, but the number of vertices whose degree is changed is not bounded uniformly. However, we can push through the results by a kind of coupling of versions 1 and 2.

Proof of Theorem 1. Both in versions 1 and 2 two old vertices are selected with replacement, independently, uniformly at random. Thus we can couple the models such that the selected vertices are the same in all steps. The duplication part is the same in the two versions. The difference is in the deletion: in version 1, the edges of the new vertex cannot be deleted. So in version 1, we do the following. In the deletion part, we colour an edge red if it is saved in version 2. That is, if it connects the new vertex with the old vertex to be deleted. In the duplication part, copies of red edges are also red: if there is a red edge between the duplicated vertex and one of its neighbors, then the new edge connecting this neighbor to the new vertex is also red. All other new edges are originally black, but they may turn red in the deletion part of the same step.

The colouring is defined in such a way that the graph sequence of the black edges is a realization of version 2. Indeed, edges turning red are deleted and hence the copies of them does not appear in this model, but all other edges are black.

Our goal is to prove that the number of vertices having red edges divided by n tends to zero almost surely. This implies that the results of Corollary 6 holds for version 1 as well.

First we need an upper bound for the total number of edges.

Lemma 2. *Denote by S_n the number of edges (both black and red ones) after n steps in version 1. Then for all $\varepsilon > 0$ we have $S_n = O(n \log^{1+\varepsilon} n)$ with probability 1.*

Proof. Let $\delta_n = S_n - S_{n-1}$. As before, \mathcal{F}_n denotes the σ -field generated by the first n steps, and $X[n, d]$ is the number of vertices of degree d after n steps. Let U_n , resp. V_n , denote the degree of the old vertex selected for duplication, resp. deletion, at step n . The new vertex is connected to the duplicated one with an edge that cannot be deleted; this increases the number of edges by 1 for sure. Thus, $\delta_n = U_n - V_n + 1$. Clearly, U_n and V_n are conditionally i.i.d. with respect to \mathcal{F}_{n-1} , hence $S_n - n = \sum_{j=1}^n (\delta_j - 1)$ is a zero mean martingale. Consequently, $ES_n = n$ for every n .

Clearly,

$$E(|\delta_n - 1| \mid \mathcal{F}_{n-1}) \leq 2E(U_n \mid \mathcal{F}_{n-1}) = \sum_{d=0}^n \frac{X[n-1, d]}{n} d = \frac{2S_{n-1}}{n}.$$

Hence

$$E\left(\sum_{n=2}^{\infty} \frac{|\delta_n - 1|}{n \log^{1+\varepsilon} n}\right) < \infty,$$

therefore the series

$$\sum_{n=2}^{\infty} \frac{\delta_n - 1}{n \log^{1+\varepsilon} n}$$

is convergent with probability 1. Then Kronecker's lemma [16, Lemma IV.3.2] implies that

$$\frac{S_n - n}{n \log^{1+\varepsilon} n} \rightarrow 0 \quad \text{a.s.}$$

as $n \rightarrow \infty$. □

Now we will colour some of the vertices red in such a way that the remaining black vertices cannot have any red edges. We will be able to give an upper bound for the number of red vertices.

At the duplication step the new vertex becomes red if and only if the duplicated vertex is red. If this old vertex is black and has no red edges, the same holds for the new vertex at the moment. After that, if there is an edge between the new vertex and the deleted one, this edge may turn red, as we defined before. We colour both endpoints of this new red edge red. On the other hand, if the old vertex chosen for deletion loses all its edges, then its new colour will be black. Note that black vertices still have only black edges, but it may happen that an old vertex has only one red edge which is deleted, because its other endpoint is chosen for deletion; in this case the vertex stays red without having any red edges.

The proof continues with giving an upper bound for the number of red vertices.

Lemma 3. *Denote by Z_n the number of red vertices after n steps. Then for all $\varepsilon > 0$ we have $Z_n = O(\log^{2+\varepsilon} n)$ almost surely.*

Proof. At each step, every old vertex has the same probability to be duplicated or deleted. If a red vertex is duplicated, then the new vertex becomes red; if it is deleted, then Z_n decreases by 1 unless the deleted vertex is connected to the new one which turns this edge red. Therefore without the exceptional new red edge, the conditional expectation of Z_n with respect to \mathcal{F}_{n-1} would be equal to Z_{n-1} . The deleted vertex and the new one are connected if and only if the deleted and duplicated vertices are the same or they are connected to each other. Since we did sampling with replacement, the probability of the first event is $1/n$; while the probability of the second event is $2S_{n-1}/n^2$. In the first case, the new vertex is red originally, but the other one stays red instead of turning back to black when deleted; Z_n is increased by an extra 1. In the other case, both endpoints of the edge turning red may be red vertices in addition. To sum up, we obtain that

$$E(Z_n|\mathcal{F}_{n-1}) \leq Z_{n-1} + \frac{1}{n} + 4 \cdot \frac{S_{n-1}}{n^2}.$$

We set $\eta_n = Z_n - Z_{n-1}$. With this notation

$$(6) \quad E(\eta_n|\mathcal{F}_{n-1}) \leq \frac{1}{n} + 4 \cdot \frac{S_{n-1}}{n^2}$$

We have already shown that $ES_{n-1} = n - 1$, hence $E\eta_n \leq 5/n$, and $EZ_n = O(\log n)$.

Note that the number of red vertices cannot change by more than three at a single step, because if an old vertex is neither deleted, nor duplicated, it cannot be coloured red. Hence $|\eta_n| \leq 3$ for all n . Moreover, we can give an upper bound on the probability that the number of red vertices changes at step n . Namely, it can change only if

- we duplicate and delete the same vertex; this has (conditional) probability $1/n$.
- the duplicated and the deleted vertices are connected to each other; this has probability $2S_{n-1}/n^2$, because there are S_{n-1} edges.
- a red vertex is duplicated; this has probability Z_{n-1}/n .
- a red vertex is deleted; this has probability Z_{n-1}/n .

Thus

$$(7) \quad P(Z_n \neq Z_{n-1}|\mathcal{F}_{n-1}) \leq \frac{1}{n} + 2\frac{S_{n-1}}{n^2} + 2\frac{Z_{n-1}}{n},$$

therefore

$$E|\eta_n| \leq 3P(Z_n \neq Z_{n-1}) = O\left(\frac{\log n}{n}\right),$$

which implies that

$$E\left(\sum_{n=2}^{\infty} \frac{|\eta_n|}{\log^{2+\varepsilon} n}\right) < \infty.$$

The proof can be completed by the help of Kronecker's lemma, just like in the proof of Lemma 2. \square

Now we can finish the proof of Theorem 1.

The total number of vertices is $n + 1$ after n steps, hence the proportion of red vertices converges to 0 almost surely as $n \rightarrow \infty$. Since we defined the colours in such a way that red edges are exactly the edges that are present in version 1 but are not present in version 2, and only red vertices may have red edges, it follows that the proportion of vertices having different degree in the two versions converges to 0. Corollary 6 states that for every d the proportion of vertices of degree d in version 2 converges almost surely to c_d . Now the same follows for version 1, which is the statement of Theorem 1. \square

Remark 1. *We could have given an upper bound for the conditional expectation of the number of red edges. The advantage of using red vertices is the uniform bound on the total change in their number; there is no such bound for the change in the number of red edges.*

Remark 2. *It follows that version 1 has a quite specific structure: it consists of cliques that are connected with relatively few edges (those are coloured red). An edge can be red only if both its endpoints are red, hence Lemma 3 gives an $O(\log^{4+\varepsilon} n)$ bound for the number of red edges.*

This is not sharp; however, the estimates of Lemmas 2 and 3 can be further improved, which might be, as pointed out above, of independent interest. Thus, before turning to the proof of Theorem 2, we present the following improvement.

Proposition 7. *$S_n \sim n$, and $Z_n = O(\log^{1+\varepsilon} n)$ for every $\varepsilon > 0$ almost surely, as $n \rightarrow \infty$.*

Proof. First we give a crude bound for the maximal degree $M_n = \max\{d : X[n, d] > 0\}$. According to Lemma 2, $S_n = O(n \log^{1+\varepsilon} n)$ also holds for the number of edges in version 2. Since a clique of size k contains $\binom{k}{2}$ edges, it follows that the size of the maximal clique is $O(n^{1/2+\varepsilon})$. The same holds for the maximal degree in version 2; and, by Lemma 3, in version 1, too. Thus $M_n = O(n^{1/2+\varepsilon})$ for every $\varepsilon > 0$.

Next, consider the martingale $S_n - n = \sum_{j=1}^n (\delta_j - 1)$ from the proof of Lemma 2. In order to prove that $S_n - n = o(\gamma_n)$ for a positive increasing predictable sequence (γ_n) it is sufficient to show that

$$\sum_{n=1}^{\infty} \gamma_n^{-2} E((\delta - 1)^2 \mid \mathcal{F}_{n-1}) < \infty$$

with probability 1 [16, Theorem VII.5.4]. To this end we need to estimate the conditional variance of the martingale differences.

$$\begin{aligned} \text{Var}(\delta_n - 1 | \mathcal{F}_{n-1}) &= 2\text{Var}(U_n | \mathcal{F}_{n-1}) \leq 2E(U_n^2 | \mathcal{F}_{n-1}) \\ &= 2 \sum_{d=1}^n \frac{X[n-1, d]}{n} d^2 \leq \frac{2}{n} M_{n-1} \sum_{d=1}^n X[n-1, d] d \\ &= \frac{2}{n} M_{n-1} S_{n-1} = O(n^{1/2+\varepsilon}), \end{aligned}$$

for every positive ε . Hence

$$\sum_{n=1}^{\infty} \frac{E((\delta - 1)^2 | \mathcal{F}_{n-1})}{n^{3/2+\varepsilon}} < \infty,$$

implying

$$S_n - n = o(n^{3/4+\varepsilon})$$

Thus $S_n \sim n$ a.s., indeed.

Finally, let us consider the martingale $\zeta_n = \sum_{j=1}^n (\eta_j - E(\eta_j | \mathcal{F}_{j-1}))$, where $\eta_n = Z_n - Z_{n-1}$, and derive an upper bound for the conditional variance of the differences. Keeping in mind that $|\eta_n| \leq 3$ and using (7) we have

$$\begin{aligned} E((\zeta_n - \zeta_{n-1})^2 | \mathcal{F}_{n-1}) &= \text{Var}(\eta_n | \mathcal{F}_{n-1}) \\ &\leq E((Z_n - Z_{n-1})^2 | \mathcal{F}_{n-1}) \leq 9P(Z_n \neq Z_{n-1} | \mathcal{F}_{n-1}) \\ &\leq 9 \left(\frac{1}{n} + 2 \frac{S_{n-1}}{n^2} + 2 \frac{Z_{n-1}}{n} \right) = O\left(\frac{1 + Z_{n-1}}{n} \right). \end{aligned}$$

Now suppose that $Z_n = O(\log^\alpha n)$ is satisfied for some $\alpha > 0$. Then

$$E((\zeta_n - \zeta_{n-1})^2 | \mathcal{F}_{n-1}) = O\left(\frac{\log^\alpha n}{n} \right),$$

hence

$$\sum_{n=2}^{\infty} \frac{E((\zeta_n - \zeta_{n-1})^2 | \mathcal{F}_{n-1})}{\log^{\alpha+1+\varepsilon} n} < \infty$$

with probability 1. Again, by [16, Theorem VII.5.4] we have

$$(8) \quad \zeta_n = o(\log^{(\alpha+1)/2+\varepsilon}) \quad \text{a.s.}$$

for every positive ε .

Clearly,

$$Z_n = \sum_{j=1}^n \eta_j = \zeta_n + \sum_{j=1}^n E(\eta_j | \mathcal{F}_{j-1}),$$

where the last sum can be estimated by the help of (6) in the following way. Since $S_{n-1} \sim n$, we have $E(\eta_n | \mathcal{F}_{n-1}) = O(1/n)$, hence

$$\sum_{j=1}^n E(\eta_j | \mathcal{F}_{j-1}) = O(\log n).$$

This, combined with (8) gives that $Z_n = O(\log^{(\alpha+1)/2+\varepsilon})$ holds almost surely for all $\varepsilon > 0$. By Lemma 3 we can start from $\alpha = 2 + \varepsilon$, and repeating the argument we finally end up with the a.s. estimation $Z_n = O(\log^{1+\varepsilon})$, for all $\varepsilon > 0$. \square

Proof of Theorem 2. Let $G(z)$ denote the generating function of the sequence (c_d) , that is,

$$G(z) = \sum_{d=0}^{\infty} c_d z^d, \quad |z| \leq 1.$$

Multiplying equation $(d+1)(c_{d-1} + c_{d+1}) = (2d+3)c_d$ by z^d , then summing up from $d=1$ to ∞ and using that $c_0 = (1+c_1)/3$, we obtain an inhomogeneous linear differential equation for $G(z)$.

$$(1-z)^2 G'(z) = (3-2z)G(z) - 1, \quad G(0) = c_0.$$

Solving this equation we get the following expression

$$G(z) = \frac{c(z)}{(1-z)^2} \exp\left(\frac{z}{1-z}\right),$$

where

$$c(z) = c_0 - \int_0^z \exp\left(-\frac{y}{1-y}\right) dy.$$

Since $G(1) = \sum_{d=0}^{\infty} c_d \leq 1$, it follows that

$$c_0 = \int_0^1 \exp\left(-\frac{y}{1-y}\right) dy,$$

hence, via the substitution $x = 1 - y$,

$$c(z) = \int_z^1 \exp\left(-\frac{y}{1-y}\right) dy = \int_0^{1-z} \exp\left(1 - \frac{1}{x}\right) dx.$$

Thus we have

$$G(z) = \int_0^{1-z} \exp\left(1 - \frac{1}{x}\right) dx \frac{1}{(1-z)^2} \exp\left(\frac{z}{1-z}\right),$$

from which, by substituting $y = \frac{1}{x} - \frac{1}{1-z}$, we obtain

$$\begin{aligned}
(9) \quad G(z) &= \int_0^\infty \frac{e^{-y}}{(1+(1-z)y)^2} dy \\
&= \int_0^\infty \frac{e^{-y}}{(1+y)^2(1-z\frac{y}{1+y})^2} dy \\
&= \int_0^\infty \sum_{d=0}^\infty (d+1) \frac{z^d y^d e^{-y}}{(1+y)^{d+2}} dy \\
&= \sum_{d=0}^\infty z^d (d+1) \int_0^\infty \frac{y^d e^{-y}}{(1+y)^{d+2}} dy,
\end{aligned}$$

completing the proof of the first statement of the theorem.

In addition, note that the first equality of (9) immediately implies that $\sum_{d=0}^\infty c_d = G(1) = 1$. \square

Proof of Theorem 3. In order to approximate the integral of Theorem 2 we first analyse the behavior of the integrand around the point where it attains its maximum. Let

$$y_d = \arg \max \frac{y^d e^{-y}}{(1+y)^{d+2}} = \arg \max f(y),$$

where

$$f(y) = d \log y - (d+2) \log(1+y) - y.$$

Clearly,

$$\begin{aligned}
f'(y) &= \frac{d}{y} - \frac{d+2}{y+1} - 1 = -\frac{y^2 + 3y - d}{y(y+1)}, \\
f''(y) &= -\frac{d}{y^2} + \frac{d+2}{(y+1)^2} = \frac{2y^2 - 2dy - d}{y^2(y+1)^2}, \\
f'''(y) &= \frac{2d}{y^3} - \frac{2(d+2)}{(y+1)^3}.
\end{aligned}$$

Since y_d satisfies $f'(y_d) = 0$, we get that

$$y_d = -\frac{3}{2} + \sqrt{d + \frac{9}{4}} = \sqrt{d} - \frac{3}{2} + o(1).$$

Let us write y in the form $y = y_d + y_d^{1/2}t$. Then

$$g(t) := f(y) - f(y_d) = \frac{y_d}{2} f''(y_d + \theta y_d^{1/2}t) t^2,$$

where $\theta = \theta(d, t)$ belongs to the interval $[0; 1]$. For every fixed t

$$f''(y_d + \theta y_d^{1/2}t) \sim -2y_d^{-1},$$

thus $g(t) \rightarrow -t^2$ as $d \rightarrow \infty$. Moreover, for $y \leq y_d$, that is, for $y_d^{1/2} \leq t \leq 0$ we have $f'(y) \geq 0$. Thus $d/y - (d+2)/(y+1) > 0$ holds, and

after rearranging we get that $(d+2)/d < (y+1)/y$. This yields that $(d+2)/d < (y+1)^3/y^3$ is satisfied, which implies that $f'''(y) \geq 0$. Hence

$$g(t) \leq \frac{y_d}{2} f''(y_d) t^2 = a_d t^2,$$

where $a_d \rightarrow -1$, as $d \rightarrow \infty$. On the other hand, let $y_d \leq y \leq \frac{3}{2} y_d$, that is, $0 \leq t \leq \frac{1}{2} y_d^{1/2}$. In this domain f''' is increasing, hence $f'''(y) \leq f'''(y_d) \sim 6d y_d^{-4} \sim 6d^{-1}$. Thus,

$$\begin{aligned} g(t) &\leq \frac{y_d}{2} f''(y_d) t^2 + \frac{1}{6} y_d^{3/2} f'''(y_d) t^3 \\ &\leq \left(\frac{y_d}{2} f''(y_d) + \frac{y_d^2}{12} f'''(y_d) \right) t^2 = b_d t^2, \end{aligned}$$

where $b_d \rightarrow -1/2$, as $d \rightarrow \infty$.

Thus, by the dominated convergence theorem,

$$\begin{aligned} \int_0^{3y_d/2} e^{f(y)} dy &= y_d^{1/2} \int_{-y_d^{1/2}}^{\frac{1}{2} y_d^{1/2}} \exp(f(y_d) + g(t)) dt \\ &\sim y_d^{1/2} \exp(f(y_d)) \int_{-\infty}^{+\infty} \exp(-t^2) dt = \sqrt{\pi} y_d^{1/2} \exp(f(y_d)). \end{aligned}$$

Here

$$f(y_d) = -2 \log y_d - (d+2) \log \left(1 + \frac{1}{y_d} \right) - y_d,$$

and

$$\begin{aligned} (d+2) \log \left(1 + \frac{1}{y_d} \right) &= (d+2) \left(\frac{1}{y_d} - \frac{1}{2y_d^2} \right) + o(1) \\ &= y_d + \frac{(d+2)(2y_d-1) - 2y_d^3}{2y_d^2} + o(1) \\ &= y_d + \frac{(y_d^2 + 3y_d + 2)(2y_d-1) - 2y_d^3}{2y_d^2} + o(1) \\ &= y_d + \frac{5y_d^2 + y_d - 2}{2y_d^2} + o(1), \end{aligned}$$

where we used that $y_d^2 + 3y_d = d$. Thus,

$$f(y_d) = -2 \log y_d - 2y_d - \frac{5}{2} + o(1) = -2 \log y_d - 2\sqrt{d} + \frac{1}{2} + o(1).$$

Finally,

$$\begin{aligned} \int_{3y_d/2}^{\infty} e^{f(y)} dy &\leq (2y_d)^{-2} \int_{3y_d/2}^{\infty} \left(1 - \frac{1}{1+y} \right)^d e^{-y} dy \\ &\leq (2y_d)^{-2} \int_{3y_d/2}^{\infty} \exp\left(-\frac{d}{y+1} - y\right) dy. \end{aligned}$$

The exponent on the right-hand side can be estimated with the help of the AM–GM inequality as follows.

$$-\frac{d}{y+1} - y = -\frac{d}{y+1} - \frac{y+1}{2} - \frac{y-1}{2} \leq -\sqrt{2d} - \frac{y-1}{2},$$

hence

$$\int_{3y_d/2}^{\infty} e^{f(y)} dy \leq (2y_d)^{-2} \exp\left(-\sqrt{2d} + \frac{1}{2} - \frac{3}{4} y_d\right) = o\left(y_d^{-2} \exp(-2\sqrt{d})\right).$$

From all these we obtain that

$$c_d = (d+1) \int_0^{\infty} e^{f(y)} dy \sim (e\pi)^{1/2} d^{1/4} e^{-2\sqrt{d}},$$

as claimed. \square

Proof of Theorem 4. Black vertices have the same local clustering coefficient in both versions. Since the proportion of red vertices tends to be negligible as $n \rightarrow \infty$, the limit of the average clustering coefficient is also the same in both versions. The global clustering coefficient of version 2 is identically equal to 1. In its defining fraction the numerator and the denominator are proportional to n . When turning to version 1 the denominator have to be increased by the number of triplets containing at least one red edge. Such a triplet must have a red central vertex and at least one more red vertex. Hence the increment of the denominator cannot exceed $M_n Z_n^2$, where M_n denotes the maximal degree, and Z_n the number of red vertices. In the proof of Proposition 7 we have shown that $M_n = O(n^{1/2+\varepsilon})$ and $Z_n = O(\log^{1+\varepsilon} n)$, thus the increment of the denominator is asymptotically negligible with respect to n . Hence the global clustering coefficient of version 1 must converge to 1. \square

REFERENCES

- [1] Backhausz, Á., and Móri, T. F., A random model of publication activity, *Discrete Appl. Math.* **162** (2014), 78–89.
- [2] Barabási, A-L., and Albert, R., Emergence of scaling in random networks, *Science* **286** (1999), 509–512.
- [3] Cooper, C., and Frieze, A., A general model of web graphs, *Random Structures Algorithms*, **22** (2003), 311–335.
- [4] Bebek, G., Berenbrink, P., Cooper, C., Friedetzky, T., Nadeau, J. and Sahinalp, S. C., The degree distribution of the generalized duplication model. *Theor. Comput. Sci.*, **369** (2006), 234–249.
- [5] Bollobás, B., Riordan, O., Spencer, J., and Tusnády, G., The degree sequence of a scale-free random graph process, *Random Structures Algorithms*, **18** (2001), 279–290.
- [6] Chung, F., Lu, L., Dewey, T. G., and Galas, D. J., Duplication models for biological networks, *J. Comput. Biol.*, **16** (2003), 677–687.
- [7] Cohen, N., Jordan, J., and Voliotis, M., Preferential duplication graphs, *J. Appl. Probab.*, **47** (2010), 572–585.

- [8] Dong, R., Goldschmidt, C., and Martin, J. B., Coagulation-fragmentation duality, Poisson–Dirichlet distributions and random recursive trees. *Ann. Appl. Probab.* **16** (2006), 1733–1750.
- [9] Durrett, R., *Random graph dynamics*, Cambridge University Press, 2006.
- [10] Faloutsos, M., Faloutsos, P., and Faloutsos, C., On power-law relationships of the internet topology, in *Proceedings of the Conference on Applications, Technologies, Architectures, and Protocols for Computer Communication, SIGCOMM '99*. ACM, New York, 1999, 251–262.
- [11] Hamdi, M., Krishnamurthy, V., Yin, G. G., Tracking a Markov-modulated stationary degree distribution of a dynamic random graph, *IEEE Trans. Inform. Theory* **60** (2014), no. 10, 6609–6625.
- [12] Jordan, J., Randomised reproducing graphs. *Electron. J. Probab.*, **16** (2011), 1549–1562.
- [13] Kim, J., Krapivsky, P. L., Kahng, B. and Redner, S., Infinite-order percolation and giant fluctuations in a protein interaction network. *Phys. Rev.*, E66: 055101(R), 2002.
- [14] Pastor-Satorras, R., Smith, E., and Solé, R. V., Evolving protein interaction networks through gene duplication. *J. Theor. Biol.*, **222** (2003), 199–210.
- [15] Ráth, B. and Tóth, B., Erdős–Rényi random graphs + forest fires = self-organized criticality. *Electron. J. Probab.*, **14** (2009), 1290–1327.
- [16] Shiryaev, A. N., *Probability*, 2nd ed., Springer, New York, 1996.
- [17] Simon, H. A., On a class of skew distribution functions. *Biometrika*, **42** (1955), 425–440.
- [18] Sridharan, A., Gao, Y., Wu, K., and Nastos, J., Statistical behavior of embeddedness and communities of overlapping cliques in online social networks. In: *2011 Proceedings IEEE INFOCOM*, 546–550.
- [19] Szymański, J., On a nonuniform random recursive tree, in *Random graphs '85, Poznań 1985*, North-Holland, Amsterdam, 1985, 297–306.
- [20] Yule, G. U., A mathematical theory of evolution, based on the conclusions of Dr. J. C. Willis, F.R.S., *Philos. Trans. R. Soc. Lond. Ser. B.*, **213** (1925), 402–410.
- [21] van der Hofstad, R., *Random graphs and complex networks*. Preprint, <http://www.win.tue.nl/~rhofstad/NotesRGCN.pdf>.
- [22] Watts, D. J., Strogatz, S. H., Collective dynamics of 'small-world' networks. *Nature* **393** (1998), 440–442.
- [23] Willinger, W., Alderson, D., and Doyle, J. C., Mathematics and the Internet: A source of enormous confusion and great potential, *Notices of the AMS* **56(5)** (2009), 586–599.

DEPARTMENT OF PROBABILITY THEORY AND STATISTICS, EÖTVÖS LORÁND UNIVERSITY, PÁZMÁNY P. S. 1/C, H-1117 BUDAPEST, HUNGARY

E-mail address: agnes@math.elte.hu

DEPARTMENT OF PROBABILITY THEORY AND STATISTICS, EÖTVÖS LORÁND UNIVERSITY, PÁZMÁNY P. S. 1/C, H-1117 BUDAPEST, HUNGARY

E-mail address: mori@math.elte.hu