



HAL
open science

Variations on depth-first search performance in random digraphs

Philippe Jacquet, Svante Janson

► **To cite this version:**

Philippe Jacquet, Svante Janson. Variations on depth-first search performance in random digraphs. AofA 2022 - 33rd International Conference on Probabilistic, Combinatorial, and Asymptotic Methods in the Analysis of Algorithms, Jun 2022, Philadelphia, PA, United States. hal-03537124

HAL Id: hal-03537124

<https://hal.science/hal-03537124>

Submitted on 20 Jan 2022

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.

VARIATIONS ON DEPTH-FIRST SEARCH PERFORMANCE IN RANDOM DIGRAPHS

PHILIPPE JACQUET AND SVANTE JANSON

ABSTRACT. We present an analysis of the depth-first search algorithm in a random digraph model with geometric outdegree distribution. This problem posed by Don Knuth in his next to appear volume of *The Art of Computer Programming* gives interesting insight in one of the most elegant and efficient algorithm for graph analysis due to Tarjan.

1. INTRODUCTION

The motivation of this paper is a new section in Donald Knuth's *The Art of Computer Programming* [1], which is dedicated to Depth-First Search (DFS) in a digraph. We refer to [1] for the definition of DFS as well as for historical notes. Note that the digraphs in [1] and here are multi-digraphs, where loops and multiple arcs are allowed. The DFS algorithm generates a spanning forest (the *depth-first forest*) in the digraph, with all arcs in the forest directed away from the roots. Our main purpose is to study the distribution of the depth of vertices in the depth-first forest, starting with a random digraph G .

Furthermore, the DFS algorithm in [1] classifies the arcs in the digraph into the following five types, see Figure 1 for examples:

- *loops*;
- *tree arcs*, the arcs in the resulting depth-first forest;
- *back arcs*, the arcs which point to an ancestor of the current vertex in the current tree;
- *forward arcs*, the arcs which point to an already discovered descendant of the current vertex in the current tree;
- *cross arcs*, all other arcs (these point to an already discovered vertex which is neither a descendant nor an ancestor of the current vertex, and might be in another tree).

We will discuss the numbers of arcs of different types. (See further the exercises in [1].)

The random digraph model that we consider has n vertices and a given outdegree distribution \mathbf{P} . The outdegrees (number of outgoing arcs) of the n vertices are independent random numbers with this distribution. The endpoint of each arc is uniformly selected at random among the n vertices, independently of all other arcs. (Therefore, an arc can loop back to the starting vertex, and multiple arcs can occur.) We consider asymptotics as $n \rightarrow \infty$ for a fixed outdegree distribution.

We will focus on the case of a geometric outdegree distribution; the lack-of-memory property in this case leads to interesting features and a simpler analysis. The paper will study the following outdegree distribution in the following order:

- a geometric distribution;
- a shifted geometric distribution (starting from integer 1 instead of zero);
- a general distribution.

Date: 11 January 2022.

Key words and phrases. Combinatorics, Depth-First Search, Random Digraphs.

Supported by the Knut and Alice Wallenberg Foundation.

We thank Donald Knuth for posing us questions and conjectures that led to the present paper.

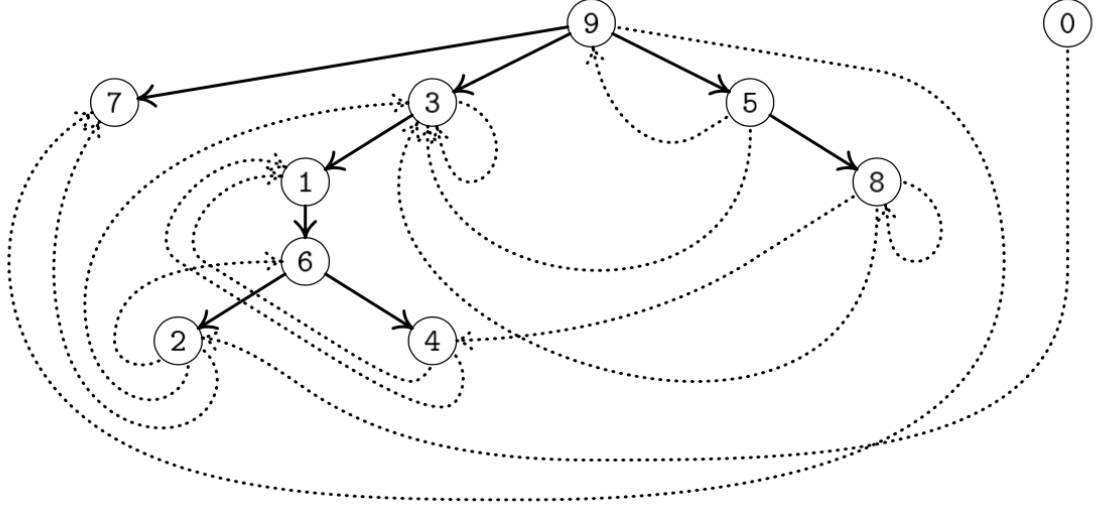


FIGURE 1. Example of a depth-first forest (jungle) from [1], by courtesy of Donald Knuth. Tree arcs are solid (e.g. $\textcircled{9} \rightarrow \textcircled{3}$). For example, $\textcircled{3} \dashrightarrow \textcircled{3}$ is a loop, $\textcircled{2} \dashrightarrow \textcircled{3}$ is a back arc, $\textcircled{9} \dashrightarrow \textcircled{7}$ is a forward arc, $\textcircled{8} \dashrightarrow \textcircled{4}$ and $\textcircled{0} \dashrightarrow \textcircled{2}$ are cross arcs.

fig:forest

1.1. **Some notation.** As usual, w.h.p. means *with high probability*, i.e., with probability $1 - o(1)$.

We let $O_{L^2}(a_n)$, where (a_n) is a sequence of positive numbers, denote a sequence of random variables X_n such that $\mathbb{E}[|X_n/a_n|^2] = O(1)$. Note that this implies $|X_n| \leq a_n \omega(n)$ w.h.p. for any sequence $\omega(n) \rightarrow \infty$.

The mean outdegree, i.e., the expectation of \mathbf{P} , is denoted by λ .

2. DEPTH ANALYSIS WITH GEOMETRIC OUTDEGREE DISTRIBUTION

Sgeo

In this section we assume that the outdegree distribution is geometric $\text{Ge}(1-p)$ for some fixed $0 < p < 1$, and thus has mean $\lambda := p/(1-p)$. Note that in the DFS, when we find a new vertex v , we do not have to immediately reveal its outdegree. Instead, we only check whether there is at least one outgoing arc (probability p), and if so, we find its endpoint and explores this endpoint if it has not already been visited; eventually, we return to v , and then we check whether there is another outgoing arc (again probability p , by the lack-of-memory property of the geometric distribution), and so on.

In the following, by a *future arc* from some vertex, we mean an arc that at the current time has not yet been seen by the DFS.

2.1. **Depth Markov chain.** Our aim is to track the evolution of the search depth as a function of the number of t of discovered vertices. Let v_t be the t -th vertex discovered by the DFS ($t = 1, \dots, n$), and let $d(t)$ be the depth of v_t in the resulting depth-first forest, i.e., the number of tree edges that connect the root of the current tree to v_t . The first found vertex v_1 is a root, and thus $d(1) = 0$.

The quantity $d(t)$ follows a Markov chain with transitions ($1 \leq t < n$):

- (i) $d(t+1) = d(t) + 1$.

This happens if, for some $k \geq 1$, v_t has at least k outgoing arcs, the first $k-1$ arcs lead to vertices already visited, and the k th arc leads to a new vertex (which then

becomes v_{t+1}). The probability of this is

$$\sum_{k=1}^{\infty} p^k \left(\frac{t}{n}\right)^{k-1} \left(1 - \frac{t}{n}\right) = \frac{(1-t/n)p}{1-pt/n}. \quad (1) \quad \boxed{\text{new}}$$

(ii) $d(t+1) = d(t)$, assuming $d(t) > 0$.

This holds if all arcs from v_t lead to already visited vertices, *i.e.*, (i) does not happen, and furthermore, the parent of v_t has at least one further arc leading to an unvisited vertex. These two events are independent. Moreover, by the lack-of-memory property, the number of further arcs from the parent of v_t has the same distribution $\text{Ge}(1-p)$. Hence, the probability that one of these further arcs leads to an unvisited vertex equals the probability in (1). The probability of (ii) is thus

$$\left(1 - \frac{(1-t/n)p}{1-pt/n}\right) \frac{(1-t/n)p}{1-pt/n}. \quad (2) \quad \boxed{\text{ii}}$$

(iii) $d(t+1) = d(t) - \ell$, assuming $d(t) > \ell \geq 1$.

This happens if all arcs from v_t lead to already visited vertices, and so do all further arcs from the ℓ closest ancestors of v_t , while the $(\ell+1)$ th ancestor has at least one further arc leading to an unvisited vertex. The argument in (ii) generalizes and shows that this has probability

$$\left(1 - \frac{(1-t/n)p}{1-pt/n}\right)^{\ell+1} \frac{(1-t/n)p}{1-pt/n}. \quad (3) \quad \boxed{\text{iii}}$$

(iv) $d(t+1) = d(t) - \ell$, assuming $d(t) = \ell \geq 0$.

By the same argument as in (ii) and (iii), except that the $(\ell+1)$ th ancestor does not exist and we ignore it, we obtain the probability

$$\left(1 - \frac{(1-t/n)p}{1-pt/n}\right)^{\ell+1}. \quad (4) \quad \boxed{\text{iv}}$$

Note that (iv) is the case when $d(t+1) = 0$ and thus v_{t+1} is the root of a new tree in the depth-first forest.

We can summarize (i)–(iv) in the formula

$$d(t+1) = (d(t) + 1 - \xi_t)^+, \quad (5) \quad \boxed{\text{dt+}}$$

where $x^+ := \max\{x, 0\}$, and ξ_t is a random variable, independent of the history, with the distribution

$$\mathbb{P}(\xi_t = k) = (1 - \pi_t)^k \pi_t, \quad k \geq 0, \quad \text{with} \quad \pi_t := \frac{(1-t/n)p}{1-pt/n} = 1 - \frac{1-p}{1-pt/n}. \quad (6) \quad \boxed{\text{xi}}$$

In other words, ξ_t has the geometric distribution $\text{Ge}(\pi_t)$. Define

$$\tilde{d}(t) := \sum_{i=1}^{t-1} (1 - \xi_i), \quad (7) \quad \boxed{\text{td}}$$

and note that (7) is a sum of independent random variables. Then (5) and induction yield

$$d(t) = \tilde{d}(t) - \min_{1 \leq j \leq t} \tilde{d}(j), \quad 1 \leq t \leq n. \quad (8) \quad \boxed{\text{tdt}}$$

We can also express these relations using generating functions. Let $p(t, z)$ be the probability generating function $\mathbb{E} z^{\xi_t}$ of ξ_t , *i.e.*,

$$p(t, z) := \frac{(1-t/n)p}{1-pt/n} \sum_{\ell \geq 0} \left(1 - \frac{(1-t/n)p}{1-pt/n}\right)^{\ell} z^{\ell} = \frac{(1-t/n)p}{1-pt/n - (1-p)z}, \quad (9) \quad \boxed{\text{ptz}}$$

and let $f(t, z) := E[z^{d(t)}]$. We then have the identity, equivalent to (5),

$$f(t+1, z) = \mathcal{N}[R(t, z)f(t, z)] \quad (10)$$

where $R(t, z) := p(t, 1/z)z$ and \mathcal{N} is the operator on power series in $z^{\pm 1}$:

$$\mathcal{N}g(z) = \Pi^+g(z) + \Pi^-g(1) \quad (11)$$

where Π^+ is the operator which removes the strictly negative powers of z and Π^- is the operator which removes the non-negative powers of z . Thus we have, since $f(1, z) = 1$,

$$f(t+1, z) = \mathcal{N}R(t, z)\mathcal{N}R(t-1, z)\mathcal{N}\cdots\mathcal{N}R(1, z). \quad (12)$$

2.2. Main result for depth analysis. Note first that (7) implies that the expectation of $\tilde{d}(t)$ is

$$\mathbb{E}[\tilde{d}(t)] = \sum_{i=1}^{t-1} (1 - \mathbb{E}\xi_i) = \sum_{i=1}^{t-1} \left(1 - \frac{1 - \pi_i}{\pi_i}\right) = \sum_{i=1}^{t-1} \left(1 - \frac{1-p}{p(1-i/n)}\right). \quad (13) \quad \boxed{\text{Etd1}}$$

Let $\theta := t/n$. We fix a $\theta^* < 1$ and obtain that, uniformly for $\theta \leq \theta^*$,

$$\mathbb{E}[\tilde{d}(t)] = \int_0^\theta \left(1 - \frac{1-p}{p(1-x/n)}\right) dx + O(1) = n\tilde{\ell}(\theta) + O(1), \quad (14) \quad \boxed{\text{Etd}}$$

where

$$\tilde{\ell}(\theta) := \int_0^\theta \left(1 - \frac{1-p}{p(1-\tau)}\right) d\tau = \theta + \frac{1-p}{p} \log(1-\theta). \quad (15) \quad \boxed{\text{t1}}$$

Note that the derivative $\tilde{\ell}'(\theta) = 1 - (1-p)/(p(1-\theta))$ is (strictly) decreasing on $(0, 1)$, *i.e.*, $\tilde{\ell}$ is concave. Moreover, if $p > \frac{1}{2}$, which we call the *supercritical* case, then $\tilde{\ell}'(0) > 0$, and (15) shows that $\tilde{\ell}(\theta)$ is positive and increasing for $\theta < \theta_0 := (2p-1)/p$. After the maximum at θ_0 , $\tilde{\ell}(\theta)$ decreases and tends to $-\infty$ as $\theta \nearrow 1$. Hence, there exists a $0 < \theta_1 < 1$ such that $\tilde{\ell}(\theta_1) = 0$; we then have $\tilde{\ell}(\theta) > 0$ for $0 < \theta < \theta_1$ and $\tilde{\ell}(\theta) < 0$ for $\theta > \theta_1$. We will see that in this case the depth-first forest w.h.p. contains a giant tree, of order and height both linear in n , while all other trees are small.

On the other hand, if $p \leq \frac{1}{2}$ (the *subcritical* and *critical* cases), then $\tilde{\ell}'(0) \leq 0$ and $\tilde{\ell}(\theta)$ is negative and decreasing for all $\theta \in (0, 1)$. In this case, we define $\theta_0 := \theta_1 := 0$ and note that the properties just stated for $\tilde{\ell}$ still hold (rather trivially). We will see that in this case w.h.p. all trees in the depth-first forest are small.

Note that in all cases, θ_1 is the largest solution in $[0, 1)$ to, recalling $\lambda = p/(1-p)$,

$$\log(1-\theta_1) = -\lambda\theta_1. \quad (16) \quad \boxed{\text{gth1}}$$

Rgth1 Remark 2.1. The equation (16) may also be written $1-\theta_1 = \exp(-\lambda\theta_1)$, which shows that θ_1 is the survival probability of a Galton-Watson process with $\text{Po}(\lambda)$ offspring distribution.

We define $\tilde{\ell}^+(\theta) := [\tilde{\ell}(\theta)]^+$. Thus, by (15) and the comments above,

$$\tilde{\ell}^+(\theta) = \begin{cases} \theta + \frac{1-p}{p} \log(1-\theta), & 0 \leq \theta \leq \theta_1, \\ 0, & \theta_1 \leq \theta \leq 1. \end{cases} \quad (17) \quad \boxed{\text{t1p}}$$

We can now state one of our main results. Proofs are given in the next subsection.

T1 Theorem 2.2. *We have*

$$\max_{1 \leq t \leq n} |d(t) - n\ell^+(t/n)| = O_{L^2}(n^{1/2}). \quad (18) \quad \boxed{\text{t1}}$$

CH Corollary 2.3. *The height Υ of the depth-first forest is*

$$\Upsilon := \max_{1 \leq t \leq n} d(t) = vn + O_{L^2}(n^{1/2}), \quad (19) \quad \boxed{\text{gU}}$$

where

$$v = v(p) := \ell^+(\theta_0) = \begin{cases} 0, & 0 < p \leq 1/2, \\ \frac{2p-1}{p} - \frac{1-p}{p} \log \frac{p}{1-p}, & 1/2 < p < 1. \end{cases} \quad (20) \quad \boxed{\text{gU}}$$

Moreover, we can show that the height Υ is asymptotically normally distributed. Details will be given in the full paper.

CA **Corollary 2.4.** *The average depth \bar{d} in the depth-first forest is*

$$\bar{d} := \frac{1}{n} \sum_{t=1}^n d(t) = \alpha n + O_{L^2}(n^{1/2}), \quad (21) \quad \text{ca}$$

where

$$\alpha = \alpha(p) := \frac{1}{2}\theta_1^2 - \frac{1-p}{p} \left((1-\theta_1) \log(1-\theta_1) + \theta_1 \right) = \frac{2p-1}{p}\theta_1 - \frac{1}{2}\theta_1^2 \quad (22) \quad \text{ga}$$

2.3. Proofs.

Proof of Theorem 2.2. Since (7) is a sum of independent random variables, $\tilde{d}(t) - \mathbb{E}\tilde{d}(t)$ ($t = 1, \dots, n$) is a martingale, and Doob's inequality yields, for all $T \leq n$,

$$\mathbb{E} \left[\max_{t \leq T} |\tilde{d}(t) - \mathbb{E}\tilde{d}(t)|^2 \right] \leq 4 \mathbb{E} [|\tilde{d}(T) - \mathbb{E}\tilde{d}(T)|^2] = 4 \sum_{i=1}^{T-1} \text{Var}(\xi_i). \quad (23) \quad \text{emma}$$

As above, fix $\theta^* < 1$, and assume, as we may, that $\theta^* > \theta_1$. Let $T^* := \lfloor n\theta^* \rfloor$, and consider first $t \leq T^*$. For $i < T^*$, we have $\text{Var} \xi_i = O(1)$, and thus, for $T = T^*$, the sum in (23) is $O(T^*) = O(n)$. Consequently, (23) yields

$$\max_{t \leq T^*} |\tilde{d}(t) - \mathbb{E}\tilde{d}(t)| = O_{L^2}(n^{1/2}). \quad (24)$$

Hence, by (14),

$$M^* := \max_{t \leq T^*} |\tilde{d}(t) - n\tilde{\ell}(t/n)| = O_{L^2}(n^{1/2}). \quad (25) \quad \text{m*}$$

For $t \leq T^*$, the definition of M^* in (25) implies

$$\left| \min_{1 \leq j \leq t} \tilde{d}(j) - n \min_{1 \leq j \leq t} \tilde{\ell}(j/n) \right| \leq M^*. \quad (26) \quad \text{gal}$$

Moreover, for $t/n \leq \theta_1$, we have $\min_{1 \leq j \leq t} \tilde{\ell}(j/n) = O(1/n)$, while for $t/n \geq \theta_1$, we have $\min_{1 \leq j \leq t} \tilde{\ell}(j/n) = \tilde{\ell}(t/n)$. Hence, for all $t \leq T^*$,

$$\min_{1 \leq j \leq t} \tilde{\ell}(j/n) = \tilde{\ell}(t/n) - \tilde{\ell}^+(t/n) + O(1/n), \quad (27)$$

and thus, by (26),

$$\left| \min_{1 \leq j \leq t} \tilde{d}(j) - n\tilde{\ell}(t/n) + n\tilde{\ell}^+(t/n) \right| \leq M^* + O(1/n). \quad (28) \quad \text{ew}$$

Finally, by (8), (25) and (28),

$$|d(t) - n\tilde{\ell}^+(t/n)| \leq 2M^* + O(1/n). \quad (29) \quad \text{jb}$$

This holds uniformly for $t \leq T^*$, and thus, by (25),

$$\max_{1 \leq t \leq T^*} |d(t) - n\tilde{\ell}^+(t/n)| = O_{L^2}(n^{1/2}). \quad (30) \quad \text{jesp}$$

It remains to consider $T^* < t \leq n$. Then the argument above does not quite work, because $\pi_t \searrow 0$ and thus $\text{Var} \xi_t \nearrow \infty$ as $t \nearrow n$. We therefore modify ξ_t . We define $\hat{\pi}_t := \max\{\pi_t, \pi_{T^*}\}$; thus $\hat{\pi}_t = \pi_t$ for $t \leq T^*$ and $\hat{\pi}_t > \pi_t$ for $t > T^*$. We may thus define independent random variables $\hat{\xi}_t$ such that $\hat{\xi}_t \sim \text{Ge}(\hat{\pi}_t)$ and $\hat{\xi}_t \leq \xi_t$ for all $t < n$. (Thus, $\hat{\xi}_t = \xi_t$ for $t \leq T^*$.)

In analogy with (7)–(8), we further define

$$\widehat{d}(t) := \sum_{i=1}^{t-1} (1 - \widehat{\xi}_i), \quad (31) \quad \boxed{\text{htd}}$$

$$\widehat{d}(t) := \widehat{d}(t) - \min_{1 \leq j \leq t} \widehat{d}(j) = \max_{1 \leq j \leq t} \sum_{i=j}^{t-1} (1 - \widehat{\xi}_i). \quad (32) \quad \boxed{\text{hd}}$$

Since $\widehat{\xi}_i \leq \xi_i$, (32) implies that $\widehat{d}(t) \geq d(t)$ for all t .

We have $\text{Var}[\widehat{\xi}_i] = O(1)$, uniformly for all $t < n$, and thus the argument above yields

$$\max_{1 \leq t \leq n} |\widehat{d}(t) - n[\widehat{\ell}(t/n)]^+| = O_{L^2}(n^{1/2}), \quad (33) \quad \boxed{\text{kasp}}$$

where

$$\widehat{\ell}(\theta) := \int_0^\theta \min \left\{ \left(1 - \frac{1-p}{p(1-\tau)}\right), \left(1 - \frac{1-p}{p(1-\theta^*)}\right) \right\} d\tau. \quad (34) \quad \boxed{\text{ht1}}$$

We have $\widehat{\ell}(\theta) = \widetilde{\ell}(\theta)$ for $\theta \leq \theta^*$, and for $\theta \geq \theta^*$, $\widehat{\ell}(\theta)$ is negative and decreasing (since $\theta^* > \theta_1$). Hence, $[\widehat{\ell}(\theta)]^+ = \widetilde{\ell}^+(\theta)$ for all $0 < \theta \leq 1$. In particular, $[\widehat{\ell}(\theta)]^+ = \widetilde{\ell}^+(\theta) = 0$ for all $\theta \geq \theta^*$, and (33) implies

$$\max_{T^* < t \leq n} \widehat{d}(t) = O_{L^2}(n^{1/2}). \quad (35)$$

Recalling $d(t) \leq \widehat{d}(t)$, we thus have

$$\max_{T^* < t \leq n} |d(t) - n\widetilde{\ell}^+(t/n)| = \max_{T^* < t \leq n} d(t) \leq \max_{T^* < t \leq n} \widehat{d}(t) = O_{L^2}(n^{1/2}), \quad (36)$$

which completes the proof. \square

Proof of Corollary 2.3. Immediate from Theorem 2.2 and (15), since we have $\max_t \widetilde{\ell}^+(t/n) = \max_\theta \widetilde{\ell}^+(\theta) + O(1/n)$ and $\max_\theta \widetilde{\ell}^+(\theta) = \widetilde{\ell}^+(\theta_0) = \widetilde{\ell}(\theta_0)$. \square

Proof of Corollary 2.4. By Theorem 2.2,

$$\frac{1}{n} \sum_{t=1}^n d(t) = \sum_{t=1}^n \widetilde{\ell}^+(t/n) + O_{L^2}(n^{1/2}) = n\alpha + O_{L^2}(n^{1/2}), \quad (37) \quad \boxed{\text{ba11}}$$

where

$$\begin{aligned} \alpha &:= \int_0^1 \widetilde{\ell}^+(\tau) d\tau = \int_0^{\theta_1} \widetilde{\ell}(\tau) d\tau = \int_0^{\theta_1} \left(\tau + \frac{1-p}{p} \log(1-\tau) \right) d\tau \\ &= \frac{1}{2}\theta_1^2 - \frac{1-p}{p} \left((1-\theta_1) \log(1-\theta_1) + \theta_1 \right), \end{aligned} \quad (38) \quad \boxed{\text{ba12}}$$

which yields (22), using (16). \square

2.4. The trees in the forest.

$\boxed{\text{TT}}$ **Theorem 2.5.** *Let N be the number of trees in the depth-first forest. Then*

$$N = \rho n + O_{L^2}(n^{1/2}), \quad (39) \quad \boxed{\text{tt}}$$

where

$$\rho = \rho(p) := 1 - \theta_1 - \frac{p}{2(1-p)} (1 - \theta_1)^2. \quad (40) \quad \boxed{\text{rho}}$$

Proof. Let $I_t := \mathbf{1}\{d(t) = 0\}$, the indicator that vertex t is a root and thus starts a new tree. Thus $N = \sum_1^n I_t$.

If $\theta_1 > 0$ (i.e., $p > \frac{1}{2}$), then Theorem 2.2 shows that w.h.p. $d(t) > 0$ in the interval $(1, n\theta_1)$, except possibly close to the endpoints. Thus the DFS will find one giant tree of order $\approx \theta_1 n$, possibly preceded by a few small trees, and, as we will see later, followed by many small trees. To obtain a precise estimate, we note that there exists a constant $c > 0$ such that $\tilde{\ell}(\theta) \geq \min\{c\theta, c(\theta_1 - \theta)\}$ for $\theta \in [0, \theta_1]$. Hence, if $t \leq n\theta_1$ and $d(t) = 0$, then $\tilde{d}(t) \leq d(t) = 0$ by (8) and, recalling (25),

$$M^* \geq n\tilde{\ell}(t/n) \geq c \min\{t, n\theta_1 - t\}. \quad (41)$$

Consequently, $d(t) = 0$ with $t \leq n\theta_1$ implies $t \in [1, c^{-1}M^*] \cup [n\theta_1 - c^{-1}M^*, n\theta_1]$. The number of such t is thus $O(M^* + 1) = O_{L^2}(n^{1/2})$, using (25).

Let $T_1 := \lceil n\theta_1 \rceil$. We have just shown that (the case $\theta_1 = 0$ is trivial)

$$\sum_{t=1}^{T_1-1} I_t = O_{L^2}(n^{1/2}). \quad (42) \quad \boxed{\text{mma}}$$

It remains to consider $t \geq T_1$. Let

$$\mu_t := \mathbb{E} \xi_t = \frac{1 - \pi_t}{\pi_t} = \frac{1 - p}{p(1 - t/n)}. \quad (43) \quad \boxed{\text{mut}}$$

For any integer $k \geq 0$, the conditional distribution of $\xi_t - k$ given $\xi_t \geq k$ equals the distribution of ξ_t . Hence,

$$\mathbb{E}[(\xi_t - k)^+] = \mathbb{E}[\xi_t - k \mid \xi_t \geq k] \mathbb{P}(\xi_t \geq k) = \mu_t \mathbb{P}(\xi_t - k \geq 0). \quad (44) \quad \boxed{\text{erika}}$$

We use again the stochastic recursion (5). Let \mathcal{F}_t be the σ -field generated by ξ_1, \dots, ξ_{t-1} . Then $d(t)$ is \mathcal{F}_t -measurable, while ξ_t is independent of \mathcal{F}_t . Hence, (5) and (44) yield

$$\begin{aligned} \mathbb{E}[d(t+1) \mid \mathcal{F}_t] &= \mathbb{E}[d(t) + 1 - \xi_t \mid \mathcal{F}_t] + \mathbb{E}[(\xi_t - 1 - d(t))^+ \mid \mathcal{F}_t] \\ &= d(t) + 1 - \mu_t + \mu_t \mathbb{P}[\xi_t - 1 - d(t) \geq 0 \mid \mathcal{F}_t] \\ &= d(t) + 1 - \mu_t + \mu_t \mathbb{P}[d(t+1) = 0 \mid \mathcal{F}_t] \\ &= d(t) + 1 - \mu_t + \mu_t \mathbb{E}[I_{t+1} \mid \mathcal{F}_t]. \end{aligned} \quad (45) \quad \boxed{\text{ele}}$$

We write $\Delta d(t) := d(t+1) - d(t)$ and $\bar{I}_t := 1 - I_t$. Then (45) yields

$$\mathbb{E}[\Delta d(t) - 1 + \mu_t \bar{I}_{t+1} \mid \mathcal{F}_t] = 0. \quad (46) \quad \boxed{\text{win}}$$

Define

$$\mathcal{M}_t := \sum_{i=1}^{t-1} \mu_i^{-1} (\Delta d(i) - 1 + \mu_i \bar{I}_{i+1}) = \sum_{i=1}^{t-1} (\mu_i^{-1} \Delta d(i) - \mu_i^{-1} + \bar{I}_{i+1}). \quad (47) \quad \boxed{\text{cm}}$$

Then \mathcal{M}_t is \mathcal{F}_t -measurable, and (46) shows that \mathcal{M}_t is a martingale. We have, with $\Delta \mathcal{M}_t := \mathcal{M}_{t+1} - \mathcal{M}_t$, using (5),

$$|\Delta \mathcal{M}_t| \leq \mu_t^{-1} |d(t+1) - d(t) - 1| + \bar{I}_{t+1} \leq \mu_t^{-1} \xi_t + 1, \quad (48) \quad \boxed{\text{cm1}}$$

and thus, since $\pi_t \leq p < 1$ for all t by (6),

$$\mathbb{E}|\Delta \mathcal{M}_t|^2 \leq 2\mu_t^{-2} \mathbb{E} \xi_t^2 + 2 = 2 \left(\frac{\pi_t}{1 - \pi_t} \right)^2 \frac{1 - \pi_t + (1 - \pi_t)^2}{\pi_t^2} + 2 = O(1). \quad (49) \quad \boxed{\text{cm2}}$$

Hence, uniformly for all $T \leq n$,

$$\mathbb{E} \mathcal{M}_T^2 = \sum_{t=1}^{T-1} \mathbb{E} |\Delta \mathcal{M}_t|^2 = O(T) = O(n). \quad (50) \quad \boxed{\text{cm3}}$$

The definition (47) yields

$$\mathcal{M}_n - \mathcal{M}_{T_1} = \sum_{t=T_1}^{n-1} \mu_t^{-1} \Delta d(t) - \sum_{t=T_1}^{n-1} \mu_t^{-1} + \sum_{t=T_1}^{n-1} \bar{I}_{t+1}. \quad (51) \quad \boxed{\text{cm4}}$$

By a summation by parts, and interpreting $\mu_n^{-1} := 0$,

$$\sum_{t=T_1}^{n-1} \mu_t^{-1} \Delta d(t) = \sum_{t=T_1+1}^n (\mu_{t-1}^{-1} - \mu_t^{-1}) d(t) - \mu_{T_1}^{-1} d(T_1). \quad (52) \quad \boxed{\text{cm5}}$$

As t increases, μ_t increases by (43), and thus $\mu_{t-1}^{-1} - \mu_t^{-1} > 0$. Hence, (52) implies

$$\begin{aligned} \left| \sum_{t=T_1}^{n-1} \mu_t^{-1} \Delta d(t) \right| &\leq \sum_{t=T_1+1}^n (\mu_{t-1}^{-1} - \mu_t^{-1}) \sup_{i>T_1} |d(t)| + \mu_{T_1}^{-1} |d(T_1)| \leq 2\mu_{T_1}^{-1} \sup_{i \geq T_1} |d(t)| \\ &= O_{L^2}(n^{1/2}) \end{aligned} \quad (53) \quad \boxed{\text{cm6}}$$

by (18), since $\tilde{\ell}^+(t/n) = 0$ for $t \geq T_1 \geq n\theta_1$. Furthermore, (50) shows that $\mathcal{M}_n, \mathcal{M}_{T_1} = O_{L^2}(n^{1/2})$. Hence, (51) yields

$$\sum_{t=T_1+1}^n I_t = n - T_1 - \sum_{t=T_1+1}^n \bar{I}_t = n - T_1 - \sum_{t=T_1}^{n-1} \mu_t^{-1} + O_{L^2}(n^{1/2}) = n\rho + O_{L^2}(n^{1/2}), \quad (54) \quad \boxed{\text{cm7}}$$

where, recalling (43),

$$\rho := 1 - \theta_1 - \int_{\theta_1}^1 \frac{p(1-\tau)}{1-p} d\tau = 1 - \theta_1 - \frac{p}{2(1-p)}(1-\theta_1)^2. \quad (55) \quad \boxed{\text{rho=}}$$

The result follows by (54) and (42). \square

The argument in the proof of Theorem 2.5 shows also the following; we omit the details.

TT2 **Theorem 2.6.** *If $p > \frac{1}{2}$, then the largest tree in the depth-first forest has order $\theta_1 n + O_{L^2}(n^{1/2})$.*

Rslow **Remark 2.7.** When $p > \frac{1}{2}$, the height is thus linear in n , unlike many other types of random trees. This might imply a rather slow performance of algorithms that operate on the depth-first forest.

Conjecture 2.8. *If $p = \frac{1}{2}$, then the largest tree has order roughly $n^{2/3}$. If $p < \frac{1}{2}$, then the largest tree has order roughly $\log n$.*

2.5. Types of arcs. Recall from the introduction the classification of the arcs in the digraph G . Since we assume that the outdegrees are $\text{Ge}(1-p)$ and independent, the total number of arcs, M say, has a negative binomial distribution with mean λn , and, by a weak version of the law of large numbers,

$$M = \lambda n + O_{L^2}(n^{1/2}). \quad (56) \quad \boxed{\text{sw}}$$

In the following theorem, we give the asymptotics of the expected number of arcs of each type. We conjecture that $\text{Var } B, \text{Var } F, \text{Var } C = O(n)$, so that (59)–(60) can be improved to asymptotics for the random numbers similar to (58).

Tarcs **Theorem 2.9.** *Let L, T, B, F and C be the numbers of loops, tree arcs, back arcs, forward arcs, and cross arcs in the random digraph G .*

$$L = O_{L^2}(1), \quad (57) \quad \boxed{\text{tal}}$$

$$T = \chi n + O_{L^2}(n^{1/2}), \quad (58) \quad \boxed{\text{tat}}$$

$$\mathbb{E} F = \mathbb{E} B = \lambda \mathbb{E} \bar{d} = \beta n + O(n^{1/2}), \quad (59) \quad \boxed{\text{taf}}$$

$$\mathbb{E} C = \chi n + O(n^{1/2}), \quad (60) \quad \boxed{\text{tac}}$$

where

$$\beta := \lambda\alpha = \frac{2p-1}{1-p}\theta_1 - \frac{\lambda}{2}\theta_1^2, \quad (61) \quad \boxed{\text{gb}}$$

$$\chi := 1 - \rho = \theta_1 + \frac{\lambda}{2}(1 - \theta_1)^2. \quad (62) \quad \boxed{\text{chi}}$$

Remark 2.10. Part of the theorem is the exact equality $\mathbb{E}B = \mathbb{E}F$ for any n and p . Knuth [1] conjectures, based on exact formulas for small n , that, much more strongly, B and F have the same distribution for every n .

Moreover, Knuth [1] conjectures that $\mathbb{E}C = \mathbb{E}T$ for every n . (The theorem above shows only that this holds asymptotically as $n \rightarrow \infty$.)

Proof. L: A simple argument with generating functions shows that the number of loops at a given vertex v is $\text{Ge}(1 - p/(n - np + p))$; these numbers are independent, and thus $L \sim \text{NegBin}(n, 1 - p/(n - np + p))$ with $\mathbb{E}L = p/(1 - p) = O(1)$ and $\text{Var}(L) = p(1 - p + p/n)/(1 - p)^2 = O(1)$ [1]. Moreover, it is easily seen that asymptotically, L has a Poisson distribution, $L \xrightarrow{d} \text{Po}(\lambda)$

T: This follows immediately from Theorem 2.5, since $T = n - N$.

B, F: Let v, w be two distinct vertices. If the DFS finds w as a descendant of v , then there will later be $\text{Ge}(1 - p)$ arcs from w , and each has probability $1/n$ of being a back arc to v . Similarly, there will be $\text{Ge}(1 - p)$ future arcs from v , and each has probability $1/n$ of being a forward arc to w . Hence, if I_{vw} is the indicator that w is a descendant of v , and B_{vw} [F_{vw}] is the number of back [forward] arcs vw , then

$$\mathbb{E}B_{vw} = \mathbb{E}F_{vw} = \frac{\lambda}{n} \mathbb{E}I_{vw}. \quad (63)$$

Summing over all pairs of distinct v and w , we obtain

$$\mathbb{E}B = \mathbb{E}F = \frac{\lambda}{n} \mathbb{E} \sum_w \sum_{v \neq w} I_{vw} = \frac{\lambda}{n} \mathbb{E} \sum_w d(w), \quad (64)$$

and (59) follows by Corollary 2.4.

C: We have $L + T + B + F + C = M$, and thus (60) follows from (56) and (57)–(59), noting that (61)–(62) imply $\beta + \chi = \lambda/2$, and thus $\lambda - (\chi + 2\beta) = \chi$. \square

3. DEPTH, TREES AND ARC ANALYSIS IN THE SHIFTED GEOMETRIC OUTDEGREE DISTRIBUTION

In this section, the outdegree distribution is $\text{Ge}_1(1 - p) = 1 + \text{Ge}(1 - p)$. Thus its mean $\lambda = 1/(1 - p)$. As in Section 2, the depth $d(t)$ is a Markov chain given by (5), but the distribution of ξ_t is now different. The probability (1) is replaced by $(1 - t/n)/(1 - pt/n)$, but the number of future arcs from an ancestor is still $\text{Ge}(1 - p)$, and, with $\theta := t/n$,

$$\mathbb{P}(\xi_t = k) = \begin{cases} \bar{\pi}_t := \frac{1-\theta}{1-p\theta}, & k = 0, \\ (1 - \bar{\pi}_t)(1 - \pi_t)^{k-1}\pi_t, & k \geq 1, \end{cases} \quad (65) \quad \boxed{\text{x11}}$$

where $\pi_t = p\bar{\pi}_t$ is as in (6). The probability generating function of ξ_t is, instead of (9),

$$p(t, z) = \bar{\pi}_t + (1 - \bar{\pi}_t) \frac{\pi_t z}{1 - (1 - \pi_t)z} = (1 - \theta) \frac{1 - (1 - p)z}{1 - p\theta - (1 - p)z}. \quad (66)$$

We now have $\mathbb{E}\xi_i = \frac{(1-p)\theta}{p(1-\theta)}$ and instead of (14) we have $\mathbb{E}\tilde{d}(t) = n\tilde{\ell}(\theta) + O(1)$ where now $\tilde{\ell}(\theta)$ takes the new value

$$\tilde{\ell}(\theta) := \frac{1}{p}\theta + \frac{1-p}{p} \log(1 - \theta) \quad (67) \quad \boxed{\text{t11}}$$

The rest of the analysis does not change but we get different values for the constants. Note that $\tilde{\ell}(\theta_1) = 0$ still gives (16), now with $\lambda = 1/(1-p)$, and that $\lambda > 1$ for every p . Differentiating (67) shows that the maximum point $\theta_0 = p > 0$.

The results in Theorems 2.2–2.6 thus hold, with, by straightforward calculations,

$$v := \tilde{\ell}(p) = 1 + \frac{1-p}{p} \log(1-p), \quad (68)$$

$$\alpha := \frac{1}{p} \left(\frac{\theta_1^2}{2} - \frac{1}{\lambda} ((1-\theta_1) \log(1-\theta_1) - \frac{1}{\lambda} \theta_1) \right) = \theta_1 - \frac{\theta_1^2}{2p}, \quad (69)$$

$$\rho := 1 - \theta_1 - \frac{1}{2(1-p)} (1 - \theta_1)^2. \quad (70)$$

Figure 2 shows $\tilde{\ell}(\theta)$ for both geometric distributions.

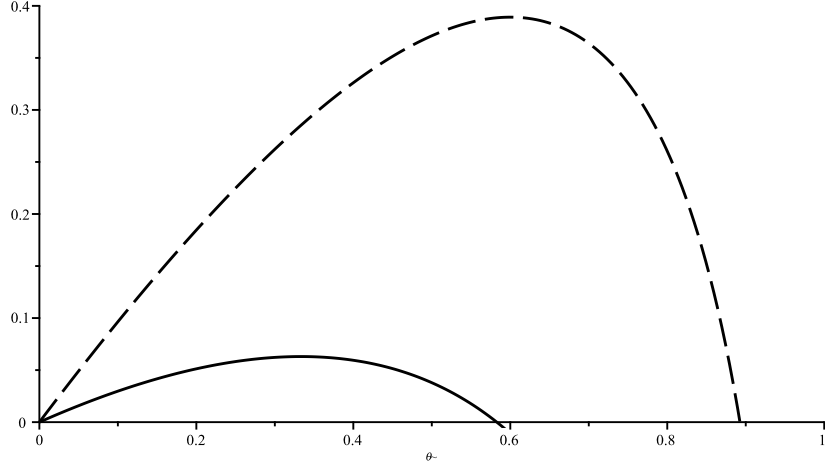


FIGURE 2. $\tilde{\ell}(\theta)$, the asymptotic search depth, for geometric distribution (solid) and shifted geometric distribution (dashed) with $p = 0.6$.

fig:depth

Now the expected numbers of back and forward arcs differ since $\mathbb{E}B = \lambda \mathbb{E}\bar{d} \sim \lambda \alpha n$ and $\mathbb{E}F = (\lambda - 1) \mathbb{E}\bar{d} \sim (\lambda - 1) \alpha n$ because the average number of future arcs at a vertex after a descendant have been created is $\lambda - 1$. Thus the equality of the number of backward and forward arcs in Theorem 2.9 was an artefact of the geometric degree distribution.

The number of tree arcs is still

$$T = \chi n + O_{L^2}(n^{1/2}) \quad (71)$$

with $\chi = 1 - \rho = \theta_1 + \frac{\lambda}{2}(1 - \theta_1)^2$. But the number of cross arcs differs. We get

$$\mathbb{E}C = \chi' n + O(n^{1/2}) \quad (72)$$

with

$$\chi' := \lambda - \chi - \lambda \alpha - (\lambda - 1) \alpha = \frac{\lambda}{2} \left(1 - 2\theta_1 + \frac{1}{p} \theta_1^2 \right) = \frac{\lambda}{2} (1 - \theta_1)^2 + \frac{1}{2p} \theta_1^2. \quad (73)$$

Again the (asymptotic) identity of $\mathbb{E}T$ and $\mathbb{E}C$ was an effect of the geometric distribution.

4. STACK INDEX ANALYSIS AND FOREST SIZE FOR A GENERAL OUTDEGREE DISTRIBUTION

In this section, we consider a general outdegree distribution, with mean λ and finite variance. When the outdegree distribution is general, the depth does not longer follow an easy Markov chain, since we should keep track of the number of children seen so far at each level of the branch of the tree toward the current vertex.

We denote m the degree of a random vertex. When the distribution of m is general, the depth does not longer follow an easy Markov chain, since we should keep track of the number of children seen so far at each level of the branch of the tree toward the current vertex. However we can set the following theorem.

Theorem 4.1. *In the general outdegree distribution when $\lambda > 1$ we have the following expansion for small values of t :*

$$\mathbb{E}[d(t)] = (1 - T(1))t + O(t^2/n) \quad (74)$$

with $T(z)$ the p.g.f. of the finite tree of the Galton Watson tree generated by the outdegree m .

Remark 4.2. It turns out that the theorem also holds for $\lambda < 1$, but in this case $T(1) > 1$.

Proof. The p.g.f. $T(z)$ satisfies $T(z) = z \mathbb{E}[T(z)^m]$. For the geometric distribution, when $p > 0.5$ we have

$$T(z) = \frac{1 - \sqrt{1 - 4(1-p)pz}}{2p} \quad (75)$$

and $T(1) = \frac{1-p}{p}$ is the probability to have a finite tree. For the shifted geometric distribution we simply have $T(z) = 0$. If we omit the variation of the fraction of discovered vertices θ , the tree search experience starts by a sequence of finite trees, before entering into the infinite tree. Said in other words, the depth reaches the level 1 and will never come back to zero. But it will hit level 1 a finite time before leaving it for ever. If the degree of the root is m , the search is limited to m subtrees among them there are one or more infinite trees. The probability that the first infinite sub-tree is the k th sub-tree is $T(1)^{k-1}(1 - T(1))$, therefore the cumulated size of the finite trees before the first infinite subtree is $\frac{1-T(z)^m}{1-T(z)}(1 - T(1))$, multiplying by z we get the p.g.f. of the duration to reach the ultimate level 2 (*i.e.* the level which will never be decremented afterward). Averaging over the distribution of degrees and normalizing on the fact that there is always an infinite sub-tree we get $\frac{z - z \mathbb{E}[T(z)^m]}{1 - T(z)}$. Since $z \mathbb{E}[T(z)^m] = T(z)$ we get $\frac{z - T(z)}{1 - T(z)}$ whose average obtained by the first derivative at $z = 1$ is exactly $\frac{1}{1 - T(1)}$.

The average duration from the ultimate level 2 to the ultimate is also $\frac{1}{1 - T(1)}$. And so forth to reach the ultimate level k we need an average duration $\frac{k}{1 - T(1)}$. Therefore we have

$$\mathbb{E}[d(t)] = (1 - T(1))t(1 + O(\theta)) \quad (76)$$

the term $O(\theta)$ comes from the impact of the variation of θ when t is small which have neglected in the first place. \square

However we don't know how to go further. If we want to get back to a Markov chain we must replace the depth by the stack index $I(t)$. The DFS can be regarded as keeping a stack of unexplored arcs, for which we have seen the start vertex but not the end. The stack evolves as follows:

stack1

- (1) If the stack is empty, pick a new vertex v that has not been seen before (if there is no such vertex, we have finished). Otherwise, pop the last arc from the stack, and reveal its endpoint v (which is uniformly random over all vertices). If v already is seen, repeat.
- (2) (v is now a new vertex) Reveal the outdegree m of v and add to the stack m new arcs, with unspecified endpoints. GOTO 1.

Let again v_t be the t th vertex seen by the DFS, and let $I(t)$ be the size of the stack when $v(t)$ is found (but before we add the arcs from v_t). Also let η_t be the outdegree of v_t . Then $I(1) = 0$ and, in analogy with (5),

$$I(t+1) = (I(t) + \eta_t - 1 - \xi_t)^+, \quad 1 \leq t < n, \quad (77) \quad \boxed{\text{It}}$$

where ξ_t is the number of arcs leading to already seen vertices before we find a new one; we have $\mathbb{P}(\xi = k) = (1 - \frac{t}{n}) (\frac{t}{n})^k$ and thus $\xi_t \sim \text{Ge}(1 - t/n)$.

In analogy with (7), we define also

$$\tilde{I}(t) := \sum_{i=1}^{t-1} (\eta_i - 1 - \xi_i) = \sum_{i=1}^{t-1} \zeta_i, \quad (78) \quad \boxed{\text{tIt}}$$

where we define $\zeta_t := \eta_t - \xi_t - 1$. Then, as in (8),

$$I(t) = \tilde{I}(t) - \min_{1 \leq j \leq t} \tilde{I}(j). \quad (79) \quad \boxed{\text{ItI}}$$

Note that

$$\mathbb{E} \zeta_t = \mathbb{E} \eta_t - \mathbb{E} \xi_t - 1 = \lambda - \frac{t/n}{1 - t/n} - 1 = \lambda - \frac{1}{1 - t/n}. \quad (80)$$

Hence, uniformly in $t/n \leq \theta^*$ for any fixed $\theta^* < 1$,

$$\mathbb{E} \tilde{I}(t) = \sum_{i=1}^{t-1} \mathbb{E} \zeta_i = (t-1)\lambda - \sum_{i=1}^{t-1} \frac{1}{1 - t/n} = n\tilde{\iota}(t/n) + O(1), \quad (81)$$

where

$$\tilde{\iota}(\theta) := \int_0^\theta \left(\lambda - \frac{1}{1 - \tau} \right) d\tau = \lambda\theta + \log(1 - \theta). \quad (82)$$

Let

$$\tilde{\iota}^+(\theta) := [\tilde{\iota}(\theta)]^+ = \begin{cases} \lambda\theta + \log(1 - \theta), & 0 \leq \theta \leq \theta_1, \\ 0, & \theta_1 \leq \theta \leq 1, \end{cases} \quad (83)$$

where again θ_1 is the largest root in $[0, 1)$ of (16), now with $\lambda = \mathbb{E} \eta_1$, the mean of \mathbf{P} . The proof of Theorem 2.2 applies with very minor differences, and yields:

TG1 **Theorem 4.3.** *Suppose that the outdegree distribution has finite variance. Then*

$$\max_{1 \leq t \leq n} |I(t) - n\tilde{\iota}^+(t/n)| = O_{L^2}(n^{1/2}). \quad (84)$$

Moreover, v_{t+1} is a root if and only if $I(t) + \zeta_t = I(t) + \eta_t - 1 - \xi_t < 0$, cf. (77). The arguments in the proof of Theorem 2.5 apply with minor differences, and show:

Theorem 4.4. *Theorems 2.5 and 2.6 hold for any outdegree distribution with finite variance, with $\rho := 1 - \theta_1 - \frac{\lambda}{2}(1 - \theta_1)^2$.*

The previous results suggest that the $\mathbb{E} \tilde{I}(t)$ and $\mathbb{E} \tilde{d}(t)$ asymptotically are just proportional by a fixed factor independent of t ; note that both have the same root θ_1 . We make the following conjecture.

Conjecture 4.5. *If $\lambda > 1$, then, for all $t \in [0, 1]$, we have $\mathbb{E} \tilde{d}(t) = \frac{1-z_1}{\lambda-1} \mathbb{E} \tilde{I}(t) + o(n)$ where $z_1 \in [0, 1)$ satisfies $z_1 = \mathbb{E}[z_1^\eta]$.*

Remark 4.6. The quantity z_1 is the probability that a Galton–Watson tree generated by the outdegree distribution η is finite. The reason to invoke this probability is that $\mathbb{E} \tilde{d}(t)$ should be asymptotically equivalent to θ times the probability of an infinite tree $(1 - z_1)$ when $\theta \rightarrow 0$. For the geometric distribution the probability is $\frac{1-p}{p}$ thus with $\lambda = \frac{p}{1-p}$ the proportional factor is $\frac{1-p}{p}$ and we find again our first result. When $\lambda < 1$ the fixed point

$z_1 > 1$ and the proportional factor remains positive. For the shifted geometric distribution the Galton–Watson tree is always in finite, thus $z_1 = 0$ and the proportional factor is again $\frac{1-p}{p}$.

REFERENCES

Knuth12A

- [1] Donald E. Knuth, *The Art of Computer Programming*, Section 7.4.1.2 (Preliminary draft, 13 December 2021). <http://cs.stanford.edu/~knuth/fasc12a.ps.gz>

INRIA SACLAY ÎLE DE FRANCE, FRANCE
Email address: `philippe.jacquet@inria.fr`

DEPARTMENT OF MATHEMATICS, UPPSALA UNIVERSITY, PO Box 480, SE-751 06 UPPSALA, SWEDEN
Email address: `svante.janson@math.uu.se`
URL: <http://www2.math.uu.se/~svante/>