## Persistence of centrality in random growing trees

Varun Jog varunjog@wharton.upenn.edu

Departments of Statistics & CIS Warren Center for Network and Data Sciences University of Pennsylvania Philadelphia, PA 19104 Po-Ling Loh loh@wharton.upenn.edu

Department of Statistics The Wharton School University of Pennsylvania Philadelphia, PA 19104

November 2015

#### Abstract

We investigate properties of node centrality in random growing tree models. We focus on a measure of centrality that computes the maximum subtree size of the tree rooted at each node, with the most central node being the tree centroid. For random trees grown according to a preferential attachment model, a uniform attachment model, or a diffusion processes over a regular tree, we prove that a single node persists as the tree centroid after a finite number of steps, with probability 1. Furthermore, this persistence property generalizes to the top  $K \geq 1$ nodes with respect to the same centrality measure. We also establish necessary and sufficient conditions for the size of an initial seed graph required to ensure persistence of a particular node with probability  $1 - \epsilon$ , as a function of  $\epsilon$ : In the case of preferential and uniform attachment models, we derive bounds for the size of an initial hub constructed around the special node. In the case of a diffusion process over a regular tree, we derive bounds for the radius of an initial ball centered around the special node. Our necessary and sufficient conditions match up to constant factors for preferential attachment and diffusion tree models.

# 1 Introduction

Heterogeneity is a common phenomenon arising naturally in many network datasets. Although some networks exist in which connections form at random between approximately exchangeable individuals, it is usually more realistic to assume that certain nodes occupy a more favorable position in the network than others. This could arise because particular nodes possess attributes that increase their likelihood of connectivity in relation to other nodes in the network. It could also be due to strategic network formation, which—even when nodes are indistinguishable—may settle on an equilibrium position where one node is in a more powerful position than the others, due to tradeoffs between the cost and utility of maintaining pairwise connections [9]. A third possibility is that the network is formed over a period of time, and older nodes are more likely to possess a higher degree of connectivity than newer nodes in the network. In order to quantify the amount of heterogeneity present in a network, various summary statistics have been proposed, including degree distributions, average path lengths, clustering coefficients, and different measures of centrality [8].

The Barábasi-Albert model, also known as the preferential attachment model, is one popular probabilistic framework for modeling the dynamics of a random growing graph [2]. In this model, each new node connects to existing nodes with probability proportional to the degrees of the nodes in the previous time step. Galashin [6] recently showed that with probability 1, a single node emerges as a persistent hub in a preferential attachment network, meaning it remains the highest-degree node in the graph after a finite number of time steps. In contrast, such a phenomenon does

not occur for the uniform attachment model, in which each new node connects to existing nodes uniformly at random. Intuitively, this is due to the fact that newly created nodes have a relatively high probability of replacing the current node of highest degree after the graph evolves further. Although older nodes in a uniform attachment model may not have a substantial lead in terms of degree, it is nonetheless reasonable to expect older nodes in the network to exhibit a higher level of connectivity according to some suitable measure. We confirm this intuition by tracking the dynamics of a different summary statistic, the *centrality* of a node in a random growing network, and prove that a single persistent node of highest centrality emerges almost surely in the case of uniform attachment trees, preferential attachment trees, and another related random growing tree arising from a diffusion process over a *d*-regular tree.

Numerous notions of centrality have been introduced in the literature on social networks, including degree centrality (also known as the maximum degree), distance centrality, betweenness centrality, and eigenvalue centrality (see, e.g., [3, 8]). In the case of trees, many popular notions of centrality conveniently coincide in terms of the most central node, which we will refer to as the tree centroid. The notion of a tree centroid was first introduced by Jordan [10], where it was originally called the branch weight centroid, and was subsequently studied by various authors using equivalent characterizations such as the distance center [8], rumor center [18], median vertex of a graph [21], security center, [19], accretion center [20], and telephone center [16]. In a random growing tree, we will call a vertex a "persistent centroid" if it is the tree centroid for all but finite moments in time. Our first main contribution therefore establishes the existence of a persistent centroid in each of the random growth models described above.

Centrality measures in random graphs have also been analyzed recently in the probability theory literature for devising root-finding algorithms in growing networks [4, 12]. In such settings, selecting the top K nodes with respect to an appropriate centrality measure yields a confidence set for the initial node in the random graph, where K is only required to be a function of the error probability, and not the total number of nodes. Motivated by these findings, our next contribution is to generalize the result on the persistence of a single centroid to the case of the top K central nodes. Consequently, the confidence set generated by a root-finding algorithm based on this measure of centrality is guaranteed to stabilize after a finite amount of time, which is a desirable property from the point of view of robustness.

As a final contribution, we address the following natural question: Suppose an individual wants to ensure that he or she is the persistent centroid of the network. The individual may boost his or her probability of becoming the persistent centroid by creating a large number of initial links to other nodes (i.e., forming a large "seed hub," for which it is the center node). How large should the initial hub be in order to ensure that the individual becomes the unique persistent centroid, with probability  $1 - \epsilon$ ? We answer this question for each of the random growing tree models. In the case of preferential and uniform attachment, we establish necessary and sufficient conditions for the initial hub size k. In a d-regular tree, we instead surround the first individual with a seed graph consisting of all nodes within radius r of that node, and derive necessary and sufficient conditions for the size of r. Our results are summarized in the following table:

model	necessary condition	sufficient condition
preferential attachment	$k \ge c \log(1/\epsilon)$	$k = C \log(1/\epsilon)$
uniform attachment	$k \ge c' \frac{\log(1/\epsilon)}{\log\log(1/\epsilon)}$	$k = C' \log(1/\epsilon)$
d-regular diffusion	$r \ge c'' \log \log(1/\epsilon)$	$r = C'' \log \log(1/\epsilon)$

Note that the necessary and sufficient conditions match up to constant factors for preferential attachment and *d*-regular diffusion trees, implying the existence of a threshold at  $k = \Theta(\log(1/\epsilon))$ 

and  $r = \Theta(\log \log(1/\epsilon))$ , respectively. In the case of uniform attachment, our bounds differ by a factor of  $\log \log(1/\epsilon)$ .

The remainder of the paper is organized as follows: We begin in Section 2 by defining tree centrality and establishing basic properties of centroids. We also define the random growth models that we will discuss in the paper. In Section 3, we establish persistence of a unique centroid for each of the random growing tree models, with probability 1, and then extend the result to the set of top K central nodes in Section 4. In Section 5, we explore the problem of ensuring persistent centrality of the root node by initializing the random growth model by an appropriate seed tree. We establish upper and lower bounds on the size of the initial seed as a function of the error probability of persistence. We conclude in Section 6 by discussing several interesting open problems.

#### 2 Preliminaries

We begin by introducing the notion of centrality in trees, as well as the probabilistic models of random growing trees that we will study in this paper.

#### 2.1 Centrality

Let the set of vertices of a tree T be denoted by V(T). A rooted tree is denoted by (T, u), with  $u \in V(T)$ . The subtree starting from v in a rooted tree T is denoted by  $T_{v\downarrow}$ . Define the function  $\psi_T : V(T) \to \mathbb{N}$  by

$$\psi_T(u) = \max_{v \in V(T) \setminus \{u\}} |(T, u)_{v\downarrow}|.$$
(1)

Thus,  $\psi_T(u)$  is size of the largest subtree of the rooted tree (T, u).

**Definition 2.1.** Given a tree T, a vertex  $u \in V(T)$  is called a *centroid* if

$$\psi_T(u) \le \psi_T(v), \quad \text{for all } v \in V(T).$$

For any two nodes u and v, if  $\psi_T(u) \leq \psi_T(v)$ , we say that u is at least as central as v.

The first lemma provides a characterization of tree centroids. Similar results have been discovered and rediscovered in a number of papers [13, 11, 18, 21]. We include a proof here for completeness.

**Lemma 2.1.** For a tree T on n vertices, the following statements hold:

(i) If  $v^*$  is a centroid, then

$$\psi_T(v^*) \le \frac{n}{2}.$$

- (ii) T can have at most two centroids.
- (iii) If  $u^*$  and  $v^*$  are two centroids, then  $u^*$  and  $v^*$  are adjacent vertices. Furthermore,

$$\psi_T(u^*) = |(T, u^*)_{v^*\downarrow}|, \text{ and } \psi_T(v^*) = |(T, v^*)_{u^*\downarrow}|.$$

*Proof.* It is easy to check that the results hold for n = 2, so we assume that  $n \ge 3$  for the rest of the proof. Let  $v^*$  be a centroid of T. Let the neighbors of  $v^*$  be the vertices  $\{a_1, \ldots, a_k\}$ . Note that if k = 1, then  $\psi_T(v^*) = n - 1$ , and one can check that  $\psi_T(a_1) < n - 1$ . This contradicts the assumption that  $v^*$  is a centroid. Hence, we must have  $k \ge 2$ . Denote

$$|(T, v^*)_{a_i\downarrow}| = r_i, \text{ for } 1 \le i \le k.$$

Without loss of generality, assume  $r_1 \ge r_2 \ge \cdots \ge r_k$ . Thus, we have  $\psi_T(v^*) = r_1$ . The key step is to look at the subtrees of  $(T, a_1)$ . If  $b \ne v^*$  is any neighbor of  $a_1$ , we have  $(T, a_1)_{b\downarrow} \subset (T, v^*)_{a_1\downarrow}$ . Thus,  $|(T, a_1)_{b\downarrow}| < r_1$ . Therefore, to ensure that  $\psi_T(v^*) \le \psi_T(a_1)$ , we must have

$$|(T,a_1)_{v^*\downarrow}| \ge r_1$$

which simplifies to

$$1 + \sum_{i=2}^{k} r_i \ge r_1.$$
 (2)

Adding  $r_1$  to both sides, and noting that  $\sum_{i=1}^{k} r_i + 1 = n$ , we conclude that  $r_1 \leq n/2$ , which is part (i).

To show part (ii), note that none the vertices in the set  $\bigcup_{i=2}^{k} (T, v^*)_{a_i \downarrow}$  can be centroids, since for any  $u \in \bigcup_{i=2}^{k} (T, v^*)_{a_i \downarrow}$ , we have

$$\psi_T(u) > |(T, u)_{a_1\downarrow}| = |(T, v^*)_{a_1\downarrow}| = \psi_T(v^*).$$

Thus, any centroids apart from  $v^*$  must lie in  $(T, v^*)_{a_1\downarrow}$ . For any node  $u \in (T, v^*)_{a_1\downarrow}$  such that  $u \neq a_1$ , we have

$$\psi_T(u) > |(T, u)_{v^*\downarrow}| = 1 + \sum_{i=2}^k r_i \ge r_1,$$

where the second inequality follows from inequality (2). Thus, the only potential centroid apart from  $v^*$  is the node  $a_1$ , which proves part (ii). Note that  $a_1$  can be a centroid if and only if

$$\psi_T(a_1) = |(T, a_1)_{v^*\downarrow}| = 1 + \sum_{i=2}^k r_i = r_1 = |(T, v^*)_{a_1\downarrow}| = \psi_T(v^*).$$

This proves part (iii) and concludes the proof.

We now turn our attention to growing trees. We have the following definition:

**Definition 2.2.** A collection of trees  $\{T_n\}_{n\geq 1}$  is called a sequence of growing trees if  $T_n$  has n nodes, and  $T_{n+1}$  is obtained from  $T_n$  by adding a single vertex that is attached to a vertex of  $T_n$  by a single new edge.

The next lemma concerns the evolution of centroids in sequences of growing trees.

**Lemma 2.2.** Consider a sequence of growing trees  $\{T_n\}_{n\geq 1}$ , with vertices labeled in order of appearance, so  $V(T_n) = \{v_1, v_2, \ldots, v_n\}$ . Let  $v^*(n)$  be a centroid of  $T_n$ , and let n > 2. If at some time N > n, the node  $v_{n+1}$  becomes at least as central as  $v^*(n)$ ; i.e., if

$$\psi_{T_N}(v_{n+1}) \le \psi_{T_N}(v^*(n)),$$

then for some  $n+1 \leq M \leq N$ , we must have

$$\psi_{T_M}(v^*(n)) = \psi_{T_M}(v_{n+1}),\tag{3}$$

and

$$|(T_M, v^*(n))_{v_{n+1}\downarrow}| = |(T_M, v_{n+1})_{v^*(n)\downarrow}|.$$
(4)

*Proof.* Note that for any fixed vertex v, the size of the largest subtree of  $(T_n, v)$  either increases by 1 or remains constant when the new vertex  $v_{n+1}$  joins  $T_n$ . Thus,  $\psi_{T_n}(v)$  increases by at most 1 at for each time step. At time n + 1, we have  $\psi_{T_{n+1}}(v_{n+1}) = n > \psi_{T_{n+1}}(v^*(n))$ , where the inequality follows from Lemma 2.1. Note that the difference  $\psi_T(v_{n+1}) - \psi_T(v^*(n))$  changes by at most 1 as the tree T evolves at each time step. Hence, if the difference becomes nonpositive at some time n = N, there must exist a time  $M \leq N$  when the difference is exactly zero. This implies that there exists an M such that  $k + 1 \leq M \leq N$ , so equation (3) holds.

Now consider the subtrees of  $(T_M, v^*(n))$  and  $(T_M, v_{n+1})$ . Let  $(v^*(n), u_1, u_2, \ldots, u_\ell, v_{n+1})$  denote the path from  $v^*(n)$  to  $v_{n+1}$ , where  $\ell \ge 0$ . Suppose the largest subtree of  $(T_M, v_{n+1})$  is  $(T_M, v_{n+1})_{w\downarrow}$ , for some  $w \ne u_\ell$ . It is easy to see that

$$\psi_{T_M}(v^*(n)) \ge |(T_M, v^*(n))_{u_1\downarrow}| \ge \ell + 1 + |(T_M, v^*(n))_{w\downarrow}| = \ell + 1 + |(T_M, v_{n+1})_{w\downarrow}| = \ell + 1 + \psi_{T_M}(v_{n+1}),$$

which contradicts equation (3). Thus, the largest subtree of  $(T_M, v_{n+1})$  must be  $(T_M, v_{n+1})_{u_\ell \downarrow}$ . Using the same argument for  $v^*$ , we conclude that the largest subtree of  $(T_M, v^*(n))$  must be  $(T_M, v^*(n))_{u_1 \downarrow}$ . By equation (3), we then have  $|(T_M, v_{n+1})_{u_\ell \downarrow}| = |(T_M, v^*(n))_{u_1 \downarrow}|$ . It is then easy to see that equation (4) holds, as well.

We also have the following useful result:

**Lemma 2.3.** Let  $\{T_n\}_{n\geq 1}$  be a sequence of growing trees, with  $V(T_n) = \{v_1, \ldots, v_n\}$ . At time n+1, we have the inequality

$$|(T_{n+1}, v_{n+1})_{v^*(n)\downarrow}| \ge \frac{n}{2}$$

*Proof.* As before, let  $(v^*(n), u_1, u_2, \ldots, u_\ell, v_{n+1})$  denote the path from  $v^*(n)$  to  $v_{n+1}$ . We have the equality

$$|(T_{n+1}, v_{n+1})_{v^*(n)\downarrow}| = (n+1) - |(T_{n+1}, v^*(n))_{u_1\downarrow}|$$

From Lemma 2.1, we have

$$|(T_{n+1}, v^*(n))_{u_1\downarrow}| \le 1 + \psi_{T_n}(v^*(n)) \le 1 + \frac{n}{2}$$

Substituting, we arrive at

$$|(T_{n+1}, v_{n+1})_{v^*(n)\downarrow}| = (n+1) - |(T_{n+1}, v^*(n))_{u_1\downarrow}| \ge (n+1) - \frac{n}{2} - 1 = \frac{n}{2}.$$

#### 2.2 Random growing trees

We now describe the probabilistic models generating the sequences of growing trees to be considered in this paper. Accordingly, we have the following definitions:

**Definition 2.3** (Uniform attachment). A sequence of growing trees  $\{T_n\}_{n\geq 1}$  is generated by a *uniform attachment process* if

(a)  $T_1$  consists of a single vertex  $v_1$ , and

(b)  $T_{n+1}$  is created from  $T_n$  by introducing a new vertex  $v_{n+1}$  and attaching it to a vertex in  $T_n$  uniformly at random; i.e., with probability 1/n to each existing node.

**Definition 2.4** (Preferential attachment). A sequence of growing trees  $\{T_n\}_{n\geq 1}$  is generated by a *preferential attachment process* if

- (a)  $T_1$  consists of a single vertex  $v_1$ , and
- (b)  $T_{n+1}$  is created from  $T_n$  by introducing a new vertex  $v_{n+1}$  and attaching it to a random vertex in  $T_n$ , with probability  $\frac{\deg(v_i)}{\sum_{j=1}^n \deg(v_j)}$  for vertex  $v_i \in V(T_n)$ .

**Definition 2.5** (*d*-regular tree diffusion). For  $d \ge 2$ , let  $\mathcal{G}$  be an infinite *d*-regular tree; i.e., a tree where each vertex has degree *d*. A sequence of growing trees  $\{T_n\}_{n\ge 1}$  is generated by a *d*-regular diffusion process if

- (a)  $T_1$  consists of a single vertex  $v_1 \in \mathcal{G}$ , and
- (b) if  $\mathcal{N}(T_n)$  denotes the set of neighbors of vertices in  $T_n$  not contained in  $V(T_n)$ , the tree  $T_{n+1}$  is created from  $T_n$  by picking  $v_{n+1} \in \mathcal{N}(T_n)$  uniformly at random, and adding it to  $T_n$  together with its connecting edge.

The models described above are well-studied [1, 2, 8] and are also examples of plane-oriented recursive trees [5].

#### 3 Existence of a persistent centroid

In this section, we show that with probability 1, a single centroid emerges for each sequence of random growing trees described in the previous section. We have the following definition:

**Definition 3.1.** A vertex  $v^* \in \bigcup_{n=1}^{\infty} V(T_n)$  is a *persistent centroid* for the sequence of growing trees  $\{T_n\}_{n\geq 1}$  if there exists N such that for all  $n \geq N$ , the vertex  $v^*$  is a centroid of  $T_n$ .

For a tree  $T_n$  on n vertices, let  $\mathcal{C}(T_n)$  denote the set of centroids of  $T_n$ . Note that by Lemma 2.1, we have  $|\mathcal{C}(T_n)| \in \{1, 2\}$ . Define

$$\mathcal{C}_{\text{tot}} = \bigcup_{n=1}^{\infty} \mathcal{C}(T_n),$$

so  $C_{tot}$  is the set of all vertices that are centroids at any point in time.

**Remark 3.1.** Throughout this section, we will assume that  $d \ge 3$  in the case of d-regular trees. Indeed, for d = 2, the set  $C_{tot}$  is infinite with probability 1. This is because diffusion on a 2-regular tree produces a sequence of line graphs, so the midpoint is the unique centroid if the number of vertices is odd, and the middle two nodes constitute the centroid set if the number of vertices is even. Since the number of vertices alternates between odd and even, a unique centroid cannot exist. Moreover, it is impossible for any node v to be a centroid for all but finitely many time steps: If v becomes the centroid at some time N, it will be a centroid at time  $n \ge N$  if and only if the number of additional nodes added to the left of v differs from the number of additional nodes added to the right of v by at most 1. Since nodes are added to the left or to the right with equal probability, it follows from properties of a simple random walk that with probability 1, centrality cannot persist.

We first show that the total number of vertices that have ever been centroids is finite with probability 1.

**Lemma 3.1.** For the preferential and uniform attachment models, and for the d-regular diffusion tree with  $d \geq 3$ , we have  $|C_{tot}| < \infty$ , with probability 1.

*Proof.* We aim to show that any node joining the tree "sufficiently late" has a very small chance of becoming a centroid at some future time. We first explain how Lemma 2.2 may be leveraged to substantially simplify the proof of this fact.

Let  $v^*(k)$  be a centroid of  $T_k$ . Suppose the node  $v_{k+1}$ , which joins  $T_k$  at time k+1, becomes a centroid of  $T_N$  for some large enough N. Then  $v_{k+1}$  must be at least as central as  $v^*(k)$  at time N. Consider the evolution of the vector  $(\psi_{T_n}(v_{k+1}), \psi_{T_n}(v^*(k)))$  with n. At time n = k+1, this vector is equal to  $(k, \psi_{T_{k+1}}(v^*(k)))$ , which is a point below the diagonal in  $\mathbb{N} \times \mathbb{N}$ . At each time step, this vector may perform one of four moves: move one step to the right, move one step above, move one step diagonally, or remain stationary. At time N, this walk is either on or above the diagonal, since  $v_{k+1}$  is at least as central as  $v^*(k)$ . To bound the probability of that event, we must keep track of the largest subtrees of  $(T_n, v^*(k))$  and  $(T_n, v_{k+1})$ , as well as the location of the new node  $v_{n+1}$ . However, Lemma 2.2 makes it possible to ignore complicated tree dynamics: First, the lemma indicates that we may bound the probability of the random walk crossing the diagonal by the probability of it reaching the diagonal at some time M. Second, the random walk reaches the diagonal at time M if and only if  $|(T_M, v^*(k))_{v_{k+1}\downarrow}| = |(T_M, v_{k+1})_{v^*(k)\downarrow}|$ . Thus, we may simply keep track of the random vector  $(|(T_n, v^*(k))_{v_{k+1}\downarrow}|, |(T_n, v_{k+1})_{v^*(k)\downarrow}|)$ , for  $n \ge k+1$ , and bound the probability of it reaching the diagonal. The evolution of the latter vector is significantly easier to track, since the dynamics of the tree are largely ignored. This random walk may either move one step to the right or one step up (it can also stay in the same place, but we may simply ignore those time steps).

For a point (i, j), let the probability of moving up be U(i, j) and of moving right be R(i, j). For the growing random trees we consider, these probabilities are given by:

1. **Preferential attachment:** The probability of a new node joining either  $(T_n, v^*(k))_{v_{k+1}\downarrow}$  or  $(T_n, v_{k+1})_{v^*(k)\downarrow}$  is proportional to the total number of edges incident upon the vertices in the corresponding subtrees. Thus, the probabilities governing the random walk are given by

$$R(i,j) = \frac{2i-1}{2(i+j-1)}$$
, and  $U(i,j) = \frac{2j-1}{2(i+j-1)}$ .

2. Uniform attachment: Here, the probability of a new node joining either  $(T_n, v^*(k))_{v_{k+1}\downarrow}$  or  $(T_n, v_{k+1})_{v^*(k)\downarrow}$  is proportional to the sizes of these subtrees. Thus, the probabilities are given by

$$R(i,j) = \frac{i}{i+j}$$
, and  $U(i,j) = \frac{j}{i+j}$ 

3. Diffusion on a *d*-regular tree: In this model, the probability of a new node joining either  $(T_n, v^*(k))_{v_{k+1}\downarrow}$  or  $(T_n, v_{k+1})_{v^*(k)\downarrow}$  is proportional to the respective neighborhood sizes  $\mathcal{N}((T_n, v^*(k))_{v_{k+1}\downarrow})$  and  $\mathcal{N}((T_n, v_{k+1})_{v^*(k)\downarrow})$ . These numbers depend only on the size of the corresponding subtrees, and we can write the probabilities as

$$R(i,j) = \frac{(d-2)i+1}{(d-2)(i+j)+2}, \quad \text{and} \quad U(i,j) = \frac{(d-2)j+1}{(d-2)(i+j)+1}.$$

Note that in all the examples above, the probability of joining a subtree is proportional to an affine function of the size the subtree. These are precisely the types of random walks discussed in

Lemma A.1 in Appendix A. Consider the events

 $H_k = \{v_{k+1} \text{ becomes at least as central as } v^*(k) \text{ at some future time}\}.$ 

It is enough to show that only finitely many events  $H_k$  occur, since this ensures that new vertices are added to  $C_{tot}$  only finitely many times.

As described in Lemma 2.2, the probability of event  $H_k$  is the probability that the random walk  $(|(T_n, v_{k+1})_{v^*(k)\downarrow}|, |(T_n, v^*(k))_{v_{k+1}\downarrow}|)$  reaches the diagonal at some point. Note that at n = k + 1, Lemma 2.3 gives

$$|(T_{k+1}, v_{k+1})_{v^*(k)\downarrow}| \ge k/2,$$

whereas  $|(T_{k+1}, v^*(k))_{v_{k+1}\downarrow}| = 1$ . By Lemma A.1 in Appendix A we then have

$$\mathbb{P}(H_k) \le \max_{A \ge k/2} \frac{P(A)}{2^A} \stackrel{(a)}{=} \frac{P(k/2)}{2^{k/2}},$$

where P(A) is a fixed polynomial, and equality (a) holds for all large enough  $k \ge K_0$ . We then have

$$\sum_{k=1}^{\infty} \mathbb{P}(H_k) = \sum_{k=1}^{K_0 - 1} \mathbb{P}(H_k) + \sum_{k=K_0}^{\infty} \mathbb{P}(H_k)$$
$$\leq \sum_{k=1}^{K_0 - 1} \mathbb{P}(H_k) + \sum_{k=K_0}^{\infty} \frac{P(k/2)}{2^{k/2}}$$
$$\leq \infty.$$

Using the Borel-Cantelli lemma, we conclude that with probability 1, only finitely many events  $H_k$  occur, completing the proof.

To establish the existence of a persistent centroid, we still need to show that the elements in  $C_{tot}$  do not keep replacing each another as centroids. Our next lemma establishes this fact. The result of the lemma may clearly be extended to any finite collection of vertices, showing that the centrality of all the vertices in the set will eventually separate. For any two vertices u and v, we define

$$\mathcal{D}_{\psi}(u,v) := \{ n \mid \psi_{T_n}(v) = \psi_{T_n}(u) \}.$$

**Lemma 3.2.** For each of the models described in Lemma 3.1, and for any two distinct vertices u and v, we have  $|\mathcal{D}_{\psi}(u,v)| < \infty$ , with probability 1.

Proof. By Lemma 2.2, it suffices to show that with probability 1, the random walk defined by  $(X_n, Y_n) := (|(T_n, v)_{u\downarrow}|, |(T_n, u)_{v\downarrow}|)$  touches the diagonal only finitely many times. Without loss of generality, we assume that vertex v is born after vertex u. Thus, the random walk starts when vertex v is born, and the starting point is  $(|(T_n, v)_{u\downarrow}|, 1) := (A, 1)$ . As in Lemma A.1 in Appendix A, let the probability that a vertex is added to a subtree of size i be proportional to  $i + \beta/\alpha$ , where  $1 + \beta/\alpha > 0$ . The evolution of the vector  $(X_n, Y_n)$  then follows a standard Pólya urn model, and by almost sure convergence of martingale sequences, combined with standard distributional convergence results [17], we have

$$\frac{X_n}{X_n + Y_n} \xrightarrow{\text{a.s.}} \xi \sim \text{Beta}(A + \beta/\alpha, 1 + \beta/\alpha).$$

By absolute continuity of the Beta distribution, we have  $\mathbb{P}(\xi = 1/2) = 0$ . Since the fraction converges to  $\xi \neq 1/2$ , with probability 1 it can equal 1/2 only finitely many times. This proves the lemma.

Lemmas 3.1 and 3.2 together imply the existence of a single persistent centroid. This is summarized in the following theorem:

**Theorem 1.** For the preferential and uniform attachment models, and for the d-regular diffusion tree with  $d \ge 3$ , there exists a time N and a node  $v^* \in T_N$  such that  $v^*$  is the unique centroid of  $T_n$  for all  $n \ge N$ , with probability 1.

*Proof.* By Lemma 3.1, the set of vertices that are ever centroids is finite. Clearly, if a single centroid does not persist, there exist at least two vertices that surpass each other infinitely often in terms of centrality. However, Lemma 3.2 rules out such a scenario, implying the persistence of a single centroid.  $\Box$ 

### 4 Persistence of the top K central nodes

We now extend the result of the previous section to establish persistence of the top K central nodes. The main theorem of this section has an important consequence concerning root-finding algorithms that generate a confidence set for the initial vertex of the random growing tree [4, 12]. As discussed in more detail following the statement of Theorem 2, the theorem implies the eventual stability of the confidence set selected according to the function  $\psi$ .

For  $n \geq K$ , let  $\mathcal{K}_n = \{\nu_1(n), \dots, \nu_K(n)\}$  denote the set of vertices of  $T_n$  that are most central in the following sense: For every vertex  $v \notin \mathcal{K}_n$ , we have the inequality

$$\psi_{T_n}(v) \ge \max_{\nu_i \in \mathcal{K}_n} \psi_{T_n}(\nu_i(n)).$$

The set  $\mathcal{K}_n$  contains the K vertices of  $T_n$  having the smallest largest subtrees, with ties being broken arbitrarily. We assume without loss of generality that

$$\psi_{T_n}(\nu_1(n)) \leq \psi_{T_n}(\nu_i(n)) \leq \cdots \leq \psi_{T_n}(\nu_K(n)).$$

The main result of this section is to show that with probability 1, the set  $\mathcal{K}_n$  also has the persistence property. In other words, there exist vertices  $\{v_1^*, \ldots, v_K^*\}$  and some N such that for all  $n \geq N$ , the  $v_i^*$ 's are the unique top K central nodes in  $T_n$ .

Our first lemma establishes that even the least central vertex in  $\mathcal{K}_n$  has its largest subtree size "not too large"—i.e., of size bounded by a linear function not identically equal to n. The proof requires a Pólya urn analysis that tracks the number of vertices in the subtrees connected to the first K nodes in each of the random growth models. We will again restrict our attention in the d-regular diffusion case to  $d \geq 3$ , since as discussed in Remark 3.1, persistence cannot occur in the case d = 2.

**Lemma 4.1.** For the preferential and uniform attachment models, and for the d-regular diffusion tree with  $d \ge 3$ , there exists a continuous random variable  $\xi$  satisfying  $\mathbb{P}(\xi < 1) = 1$  and

$$\psi_{T_n}(\nu_K(n)) \le \xi n,$$

almost surely, for all  $n \ge K$ .

*Proof.* Let  $\{v_1, \ldots, v_K\}$  denote the first K vertices, i.e., the vertices of  $T_K$ . For any  $n \ge K$ , we have

$$\max_{1 \le i \le K} \psi_{T_n}(v_i) \ge \max_{1 \le i \le K} \psi_{T_n}(\nu_i(n)) = \psi_{T_n}(\nu_K(n)).$$

Thus, it suffices to derive an upper bound for  $\max_{1 \le i \le K} \psi_{T_n}(v_i)$ .

For  $1 \leq i \leq n$ , let  $T_{i,n}$  be the tree in the forest formed by removing all the edges between  $\{v_1, \ldots, v_K\}$  in  $T_n$ . Clearly,

$$\psi_{T_n}(v_i) \le \max(|T_{i,n}|, n - |T_{i,n}|), \text{ for } 1 \le i \le K.$$

Thus,

$$\max_{1 \le i \le K} \psi_{T_n}(v_i) \le \max(|T_{1,n}|, n - |T_{1,n}|, \dots, |T_{K,n}|, n - |T_{K,n}|) = \max(\max_{1 \le i \le K} |T_{i,n}|, n - \min_{1 \le i \le K} |T_{i,n}|) \le n - \min_{1 \le i \le K} |T_{i,n}|.$$
(5)

Thus, an appropriate lower bound on  $\min_{1 \le i \le K} |T_{i,n}|$  will provide the desired upper bound. We establish a random linear lower bound for each of the growing graphs separately, beginning with uniform attachment. Apart from being easier to analyze, it will illustrate the idea that we will use in the other two cases.

1. Uniform attachment: The vector  $(|T_{1,n}|, \ldots, |T_{K,n}|)$  evolves according to a standard Pólya urn process with replacement matrix  $I_K$  and starting state  $(1, 1, \ldots, 1)$ . Thus,

$$\left(\frac{|T_{1,n}|}{n},\ldots,\frac{|T_{K,n}|}{n}\right) \xrightarrow{a.s.} (C_1,\ldots,C_K) \sim \text{Dirichlet}(1,1,\ldots,1).$$

By the continuous mapping theorem, we conclude that

$$\frac{1}{n}\min_{1\leq i\leq K}|T_{i,n}| \stackrel{a.s.}{\to} C = \min(C_1,\ldots,C_K),$$

where C is a continuous random variable taking values in [0, 1]. Taking inverses, we then have

$$\frac{n}{\min_{1 \le i \le K} |T_{i,n}|} \stackrel{a.s.}{\to} \frac{1}{C}$$

Note that 1/C does not have a point mass at infinity, since C is a continuous random variable. This almost sure convergence implies the existence of a random variable  $\hat{\xi}$  such that  $\mathbb{P}(\hat{\xi} < \infty) = 1$ , and which bounds  $\frac{n}{\min_{1 \le i \le K} |T_{i,n}|}$  almost surely, for all n. Hence,

$$\min_{1 \le i \le K} |T_{i,n}| \ge \frac{n}{\hat{\xi}},$$

for all  $n \geq K$ . Substituting into inequality (5), we then obtain

$$\max_{1 \le i \le K} \psi_{T_n}(v_i) = n - \min_{1 \le i \le K} |T_{i,n}| \le n(1 - 1/\hat{\xi}) := n\xi.$$

Since  $\hat{\xi} < \infty$  with probability 1, we have  $\xi < 1$  with probability 1. This concludes the proof.

2. **Preferential attachment:** At time K, the number of possible structures of  $T_K$  is finite. We denote the set of all possible trees at time K by  $\text{Trees}_K = \{t_1, t_2, \ldots, t_\kappa\}$ , where  $\kappa = |\text{Trees}_K|$ . Let  $\mathbb{P}(T_K = t_i) = p_i$ . Also let  $S_{\ell,n}$  denote the degree sum of the vertices in  $T_{\ell,n}$ . Conditioned on  $T_K = t_i$ , the vector  $(S_{1,n}, \ldots, S_{K,n})$  evolves according to a Pólya urn process with replacement matrix  $2I_K$  and initial configuration  $(\deg(v_1), \ldots, \deg(v_K))$ , corresponding to the degrees of the vertices  $(v_1, \ldots, v_K)$  in  $t_i$ . Hence, conditioned on  $T_K = t_i$ , we have

$$\left(\frac{S_{1,n}}{2(n-1)},\ldots,\frac{S_{K,n}}{2(n-1)}\right) \stackrel{a.s.}{\to} \text{Dirichlet}\left(\frac{\deg(v_1)}{2},\ldots,\frac{\deg(v_K)}{2}\right)$$

Furthermore,

$$S_{\ell,n} = 2(|T_{\ell,n}| - 1) + \deg(v_\ell), \qquad \forall 1 \le \ell \le K,$$

so it is easy to see that

$$\left(\frac{|T_{1,n}|}{n}, \dots, \frac{|T_{K,n}|}{n}\right) \stackrel{a.s.}{\to} (C_1^i, \dots, C_K^i) \sim \text{Dirichlet}\left(\frac{\deg(v_1)}{2}, \dots, \frac{\deg(v_K)}{2}\right), \tag{6}$$

as well. By the continuous mapping theorem, equation (6) implies the almost sure convergence

$$\frac{1}{n}\min_{1\leq i\leq K}|T_{i,n}| \stackrel{a.s.}{\to} \min(C_1^i,\ldots,C_K^i),$$

 $\mathbf{SO}$ 

$$\frac{n}{\min_{1\leq i\leq K}|T_{i,n}|} \xrightarrow{a.s.} \frac{1}{\min(C_1^i,\ldots,C_K^i)}$$

Thus, there exists a continuous random variable  $\hat{\xi}^i$  that bounds  $\frac{n}{\min_{1 \le i \le K} |T_{i,n}|}$ , almost surely, for all n, so

$$\min_{1 \le i \le K} |T_{i,n}| \ge \frac{n}{\hat{\xi}^i},$$

for all  $n \geq K$ . Substituting into inequality (5), we then obtain

$$\max_{1 \le i \le K} \psi_{T_n}(v_i) = n - \min_{1 \le i \le K} |T_{i,n}| \le n(1 - 1/\hat{\xi}^i).$$

Define the random variable  $\xi$  to equal  $(1 - 1/\hat{\xi}^i)$  on the event  $\{T_K = t_i\}$ . Using a similar argument as in the case of uniform attachment, we have  $\hat{\xi}^i < \infty$  for each i with probability 1, so  $\xi < 1$  with probability 1.

3. *d*-regular diffusion: As in the case of the preferential attachment model, we define the set of all possible trees at time K by  $\text{Trees}_K = \{t_1, t_2, \ldots, t_\kappa\}$ , where  $\kappa = |\text{Trees}_K|$  and  $\mathbb{P}(T_K = t_i) = p_i$ . Let  $U_{\ell,n}$  denote the number of uninfected neighbors of vertices in  $T_{\ell,n}$ . Conditioned on  $T_K = t_i$ , the vector  $(U_{1,n}, \ldots, U_{K,n})$  evolves according to a Pólya urn process with replacement matrix  $(d-2)I_K$  and initial configuration  $(d - \deg(v_1), \ldots, d - \deg(v_K))$ , where  $(\deg(v_1), \ldots, \deg(v_K))$  again denotes the degrees of  $(v_1, \ldots, v_K)$  in  $t_i$ . Then

$$\left(\frac{U_{1,n}}{(d-2)n},\ldots,\frac{U_{K,n}}{(d-2)n}\right) \stackrel{a.s.}{\to} \text{Dirichlet}\left(\frac{d-\deg(v_1)}{d-2},\ldots,\frac{d-\deg(v_K)}{d-2}\right).$$

implying that

$$\left(\frac{|T_{1,n}|}{n},\ldots,\frac{|T_{K,n}|}{n}\right) \stackrel{a.s.}{\to} (C_1^i,\ldots,C_K^i) \sim \text{Dirichlet}\left(\frac{d-\deg(v_1)}{d-2},\ldots,\frac{d-\deg(v_K)}{d-2}\right).$$

The remainder of the analysis proceeds exactly as in the case of the preferential attachment model.

This completes the proof of the lemma.

We now define the collection of vertices that ever enter the set of top K most central nodes. Let  $\mathcal{K}'_n$  denote the set  $\mathcal{K}_n$  augmented with any additional vertices that are at least as central as  $\nu_K(n)$  in  $T_n$ , and let

$$\mathcal{K}_{\text{tot}} := \cup_{n=1}^{\infty} \mathcal{K}'_n.$$

We have the following lemma, the analog of Lemma 3.1:

**Lemma 4.2.** In the same setting as Lemma 4.1, we have  $|\mathcal{K}_{tot}| < \infty$ , with probability 1.

Proof. Consider the set of events

$$B_M = \bigcap_{n \ge K} \{ \psi_{T_n}(\nu_K(n)) \le nM \},\$$

for any real  $M \in (0,1)$ . Thus,  $B_M$  is the event that the least central node in  $\mathcal{K}_n$ , i.e.,  $\nu_K(n)$ , has its largest subtree upper-bounded by nM at every time n. Now consider the event

$$H_n = \{ \exists \ell : v_{n+1} \in \mathcal{K}'_\ell \}$$

Thus,  $H_n$  is the event that  $v_{n+1}$  becomes at least as central as one of the top K central nodes at some point in the future. On the event  $H_n$ , we must have

$$\psi_{T_{\ell}}(v_{n+1}) \leq \max_{1 \leq i \leq K} \psi_{T_{\ell}}(\nu_i(n)).$$

Now define the event

 $E_i = \{v_{n+1} \text{ becomes at least as central as } \nu_i(n) \text{ at some future point}\}.$ 

We have the bound

$$\mathbb{P}(B_M \cap H_n) \le \mathbb{P}\left(B_M \cap \left(\bigcup_{i=1}^K E_i\right)\right) \le \sum_{i=1}^K \mathbb{P}(B_M \cap E_i).$$

By Lemma 2.2, we may control the probability  $\mathbb{P}(B_M \cap E_i)$  by bounding the probability that the random walk  $(|(T_\ell, v_{n+1})_{\nu_i(n)\downarrow}|, |(T_\ell, \nu_i(n))_{v_{n+1}\downarrow}|)$  reaches the diagonal. Note that this walk starts from the point  $(|(T_{n+1}, v_{n+1})_{\nu_i(n)\downarrow}|, 1)$  at time  $\ell = n + 1$ . If  $(\nu_i(n), u_1, \ldots, v_{n+1})$  is the path from  $\nu_i(n)$  to  $v_{n+1}$ , then on the event  $B_M$ , we have

$$|(T_{n+1}, v_{n+1})_{\nu_i(n)\downarrow}| = n - |(T_n, \nu_i(n))_{u_1\downarrow}|$$
  

$$\geq n - \psi_{T_n}(\nu_i(n))$$
  

$$\geq n - \psi_{T_n}(\nu_K(n))$$
  

$$\geq n - Mn$$
  

$$= n(1 - M).$$

Thus, the starting point lies below the diagonal and to the right of the point ((1 - M)n, 1). Lemma A.1 in Appendix A then implies that

$$\mathbb{P}(B_M \cap H_n) \le K \cdot \max_{A \ge (1-M)n} \frac{P(A)}{2^A} \stackrel{(a)}{=} K \cdot \frac{P((1-M)n)}{2^{(1-M)n}},$$

where (a) holds for all large enough n. The expression on the right-hand side form a convergent series in n. Applying Borel-Cantelli lemma, we conclude that for all M, the event  $H_n \cap B_M$  occurs finitely often, with probability 1. Furthermore, Lemma 4.1 implies that  $\mathbb{P}(B_M) \to 1$  as  $M \to 1$ , since the random variable  $\xi$  appearing in the lemma does not have a point mass at 1. Therefore, with probability 1, the events  $H_n$  can can occur only finitely often, which implies the desired statement.

The stability of the set of top K central nodes then follows by combining Lemmas 4.2 and 3.2, as in the proof of Theorem 1:

**Theorem 2.** For the preferential and uniform attachment models, and for the d-regular diffusion tree with  $d \ge 3$ , with probability 1, there exists a time N and a collection  $\{\nu_1^*, \ldots, \nu_K^*\} \subseteq T_N$  such that  $\{\nu_1^*, \ldots, \nu_K^*\}$  are the K most central nodes of  $T_n$ , for all  $n \ge N$ .

As mentioned at the beginning of the section, Theorem 2 has important implications for rootfinding algorithms in random growing trees: One may obtain confidence sets for the root node in uniform and preferential attachment models [4] and *d*-regular diffusion trees [12] by selecting the nodes that minimize the maximum subtree estimator  $\psi$ . Furthermore, the size of the confidence set may be taken as a fixed function  $K(\epsilon)$  of the error probability  $\epsilon$ , and does not need to grow with *n*. Theorem 2 implies that the confidence sets constructed in this manner will almost surely stabilize after a finite time, showing that the confidence set construction is in some sense robust.

## 5 Ensuring centrality of the root node

The results from the earlier sections indicate that any fixed node has some finite probability of eventually becoming the persistent centroid of a random growing tree. We consider the special case of the root node, i.e., the first vertex  $v_1$ , and ask the question: Can we ensure that  $v_1$  is the persistent centroid of the random growing tree? Note that the probability of the complementary event is at least 1/2, since there is no way to distinguish nodes  $v_1$  and  $v_2$ . However, in the preferential and uniform attachment graphs, we may boost the probability of  $v_1$  being the persistent centroid by initializing the tree with a "hub" centered at  $v_1$  of size k. In other words, the graph begins with a star configuration in which the nodes  $\{v_2, \ldots, v_{k+1}\}$  are all attached to  $v_1$ . In the case of a d-regular tree diffusion, the bounded degree makes it impossible to create a large hub at  $v_1$ . Hence, we instead begin with the subtree consisting of all nodes at a distance at most r from  $v_1$ . As a function of  $\epsilon$ , we derive bounds on the necessary and sufficient size hub size k (for preferential and uniform attachment) and the radius r (for a d-regular diffusion) to ensure the persistent centrality of  $v_1$  with probability  $1 - \epsilon$ .

We begin by deriving necessary conditions.

**Theorem 3.** The following conditions are necessary to ensure that  $v_1$  is the persistent central node, where C, C', and C'' are appropriate constants.

- (i) Preferential attachment: The hub size k is at least  $C \log(1/\epsilon)$ .
- (ii) Uniform attachment: The hub size k is at least  $C' \frac{\log(1/\epsilon)}{\log\log(1/\epsilon)}$ .
- (iii) d-regular tree diffusion: Suppose  $d \ge 3$ . The radius r is at least  $C'' \log \log(1/\epsilon)$ .

*Proof.* We begin by analyzing the preferential and uniform attachment models. Let  $P_k$  denote the probability that the next k-1 vertices  $\{v_{k+2}, \ldots, v_{2k}\}$  all join vertex  $v_2$ . Since the graph will then

be symmetric with respect to  $v_1$  and  $v_2$ , the probability of  $v_1$  not being the persistent centroid is at least  $P_k/2$ , which must in turn be less that  $\epsilon$ . This implies a bound on the required size of k. The value of  $P_k$  and the corresponding bound on k are developed in the following calculations.

1. **Preferential attachment:** We have

$$P_{k} = \frac{1}{2k} \cdot \frac{2}{2k+2} \cdot \dots \cdot \frac{k-1}{4k-4} = \frac{(k-1)!(k-1)!}{2^{k-1}(2k-2)!} = \frac{1}{2^{k-1}\binom{2k-2}{k-1}}.$$

Hence,

$$2\epsilon \ge P_k = \frac{1}{2^{k-1}\binom{2k-2}{k-1}} \ge \frac{1}{2^{k-1}4^{k-1}} = \frac{1}{2^{3k-3}},$$

using the fact that  $\binom{2k-2}{k-1} \leq 4^{k-1}$ . Thus, a hub size of  $k \geq C \log(1/\epsilon)$  is necessary.

#### 2. Uniform attachment: We have

$$P_k = \frac{1}{k+1} \cdot \frac{1}{k+2} \cdot \dots \cdot \frac{1}{2k-1} = \frac{k!}{(2k-1)!}$$

Hence,

$$2\epsilon \ge \frac{k!}{(2k-1)!} \ge \frac{k!}{(2k)!} \ge \frac{c\sqrt{k}(k/e)^k}{\sqrt{2k}(2k/e)^{2k}} = \frac{\tilde{c}e^k}{2^{2k}k^k}$$

where c and  $\tilde{c}$  are suitable constants. Taking logarithms and simplifying, we obtain

$$\log(1/\epsilon) \le k \log k + o(k \log k).$$

Thus, a hub size of  $k \ge C' \frac{\log(1/\epsilon)}{\log\log(1/\epsilon)}$  is necessary.

3. *d*-regular tree diffusion: In a *d*-regular tree, we create an *r*-ball around vertex  $v_1$  and derive bounds on the radius *r* of the ball. Starting from the *r*-ball centered at  $v_1$ , we calculate the probability that vertices added in such a manner will make  $v_1$  and  $v_2$  symmetric and indistinguishable. To ensure this, the next  $(d-1)^r$  vertices must be added to fill in the  $r^{\text{th}}$ level in the subtree  $(T, v_1)_{v_2 \downarrow}$ . This probability is equal to

$$P = \frac{(d-1)^r!}{\prod_{i=0}^{(d-1)^r-1} \left( d(d-1)^r + i(d-2) \right)}.$$

Taking  $\tau = \frac{d(d-1)^r}{d-2} \leq d^{r+1}$ , we simplify this as

$$2\epsilon > P = \frac{((d-1)^r)!}{(d-2)^{(d-1)^r} \prod_{i=0}^{(d-1)^r-1} (\tau+i)}$$
  
$$\geq \frac{((d-1)^r)!}{(d-2)^{(d-1)^r} \prod_{i=0}^{(d-1)^r-1} (d^{r+1}+i)}$$
  
$$= \frac{((d-1)^r)!(d^{r+1})!}{(d-2)^{(d-1)^r} (d^{r+1}+(d-1)^r-1)!}$$

$$\geq \frac{1}{(d-2)^{(d-1)r} \binom{d^{r+1}+(d-1)^r}{(d-1)^r}}$$
  
$$\geq \frac{1}{(d-2)^{(d-1)r} 2^{d^{r+1}}+(d-1)^r}.$$

Taking logarithms and simplifying, we then have

$$(d-1)^r \log(d-2) + \left(d^{r+1} + (d-1)^r\right) \log 2 \ge \log(1/2\epsilon).$$

Since the left-hand side is  $\Theta(d^{r+1})$ , we obtain that a radius of size  $r \geq C'' \log \log(1/\epsilon)$  is necessary.

The next result provides sufficient conditions on the size of the initial hub ensuring persistence of the root node.

**Theorem 4.** The following conditions are sufficient to ensure that  $v_1$  is the persistent central node, where  $\tilde{C}$  and  $\tilde{C}'$  are appropriate constants:

- (i) For preferential and uniform attachment, the hub size k satisfies  $k \ge \tilde{C} \log(1/\epsilon)$ .
- (ii) For diffusion over a d-regular tree, with  $d \ge 3$ , the radius r satisfies  $r \ge \tilde{C}' \log \log(1/\epsilon)$ .

*Proof.* In the case of preferential or uniform attachment, suppose we start with a hub of size K, so  $\{v_2, \ldots, v_{K+1}\}$  are all connected to vertex  $v_1$ . Let  $\mathcal{F}$  be the event that  $v_1$  becomes the persistent centroid, and let  $\mathcal{F}^*$  be the event that  $v_1$  is a tree centroid at all time points. Clearly,  $\mathbb{P}(\mathcal{F}) \geq \mathbb{P}(\mathcal{F}^*)$ . We will select the hub size to ensure that the latter probability is at least  $1 - \epsilon$ . Define the events

 $\hat{H}_i = \{v_i \text{ becomes a centroid at some time step}\}.$ 

Then

$$\left(\mathcal{F}^*\right)^c = \bigcup_{i=2}^\infty \hat{H}_i.$$

Note that

 $\hat{H}_i \subseteq H_i := \{v_i \text{ becomes at least as central as } v_1 \text{ at some time step}\},\$ 

and

 $\hat{H}_i \subseteq G_i := \{v_i \text{ becomes at least as central as the centroid of } T_{i-1} \text{ at some time step}\}.$ 

Thus,

$$\mathbb{P}\left(\left(\mathcal{F}^*\right)^c\right) \leq \mathbb{P}\left(\left(\bigcup_{i=2}^{K+1} H_i\right) \bigcup \left(\bigcup_{i=K+2}^{\infty} G_i\right)\right).$$

Using the bound from Lemma A.1 in Appendix A, and for K greater than an appropriate constant, we then have

$$P((\mathcal{F}^*)^c) \leq K \cdot \frac{P(K)}{2^K} + \sum_{i=K+2}^{\infty} \frac{P(i/2)}{2^{i/2}}$$

$$\stackrel{(a)}{\leq} \frac{2^{K/2}}{2^K} + \sum_{i=K+2}^{\infty} \frac{2^{i/4}}{2^{i/2}}$$

$$= 2^{-K/2} + \frac{2^{-(K+2)/4}}{1 - 2^{-1/4}}$$
  

$$\leq 5 \cdot 2^{-K/4}, \tag{7}$$

where in (a), we have used the fact that for large enough K,

$$P(K) < KP(K) < 2^{K/2}.$$

We can make  $5 \cdot 2^{-K/4} < \epsilon$  by choosing  $K \ge \tilde{C} \log(1/\epsilon)$  for a large enough constant  $\tilde{C}$ .

For diffusion over a d-regular tree, the same calculation (7) holds, except now

$$K = \frac{(d-1)^{r+1} - 1}{d-2} = \Theta(d^r)$$

for the initial seed graph. Thus,  $r \geq \tilde{C}' \log \log(1/\epsilon)$  is a sufficient condition for the size of the radius.

**Remark 5.1.** Comparing the necessary and sufficient conditions in Theorems 3 and 4, we see that a threshold occurs at hub size  $k = \Theta(\log(1/\epsilon))$  in the case of the preferential attachment model, and at radius  $r = \Theta(\log\log(1/\epsilon))$  in the case of a d-regular diffusion. However, the bounds on the hub size in the case of uniform attachment disagree by a factor of  $\log \log(1/\epsilon)$ . It is a topic of future work to determine the exact threshold in this case, which must lie between  $\Omega\left(\frac{\log(1/\epsilon)}{\log\log(1/\epsilon)}\right)$  and  $O(\log(1/\epsilon))$  by our results.

## 6 Discussion

We have established the persistence of a unique centroid (or set of top K central nodes) in three types of random growing trees: Uniform attachment, preferential attachment, and diffusion processes over *d*-regular trees. Furthermore, we have derived necessary and sufficient conditions for the size of the initial seed graph required to ensure that the first node is the persistent centroid in the network with probability  $1 - \epsilon$ . A number of related open questions remain:

- (i) We believe that the results in this paper regarding persistence of the centroid should hold in more general preferential attachment models, where the probability of attaching to a node is proportional to a function f of the vertex degree. In Galashin [6], it was shown that the degree-central node persists when f is a convex function. Results concerning nonlinear Pólya urns indicate that for concave f, degree-centrality cannot persist [14]. However, the centroid persists when f is either a linear or a constant function, and we conjecture that the persistence of the centroid holds for a larger class of functions, if not for all functions.
- (ii) Our results and those from Galashin [6] show that the top central node "stabilizes" after a finite time, but we are unable to provide estimates on the expected time or the distribution of the time when stabilization occurs. This is particularly relevant for practical purposes, when one may wish to guarantee that the current centroid is the persistent centroid.
- (iii) As mentioned in the remark after Theorem 4, the problem of determining the hub-size threshold for the case of uniform attachment trees is still an open question. A related topic concerning degree centrality would be to provide necessary and sufficient conditions on the hub size in order to ensure degree centrality of  $v_1$  (as in Section 5) in the convex preferential attachment model. It would be interesting to compare these conditions to the bounds required for the form of centrality studied in this paper.

- (iv) In general, the degree-central node in a tree need not be the same as the centroid. For the preferential attachment model, one might ask whether the *persistent* degree-central node is the same as the *persistent* centroid with probability 1. It is tempting to think that such a result should hold; however, it probably does not. A heuristic argument is as follows: Consider a tree T rooted at  $v_1$  with neighbors  $\{v_2, v_3, v_4, v_5\}$ . Assume  $v_2$  has a large number of children (say, 10<sup>6</sup>) and no grandchildren, and assume  $(T, v_1)_{v_i\downarrow}$ , for  $i \in \{2, 3, 4\}$ , is simply a line graph with, say, 10<sup>10</sup> nodes. A preferential attachment process starting from such a tree would likely have  $v_1$  as the persistent centroid, and  $v_2$  as the persistent degree-central node. Since one can obtain T with a finite probability starting from  $v_1$ , the persistent degree-central node study what additional constraints could ensure the agreement of both persistent nodes.
- (v) Our results show that the top K central nodes obtained according to the centrality measure  $\psi$  stabilizes after a finite number of steps. However, a confidence set constructed according to  $\psi$  may be sub-optimal in terms of the size of the set required as a function of the error probability  $\epsilon$  [4, 12]. It would be interesting to see whether other centrality measures such as those corresponding to the maximum likelihood estimator are also "robust" in the sense that they produce a stable output after some finite time.
- (vi) The problem of establishing persistence of centrality in non-trees (for example, in preferential or uniform attachment models where more than one node is added at each step) appears to be very challenging. It is not clear what notion of centrality, if any, would persist in such cases. Even for trees, the problems of establishing persistent centrality in alternative models such as preferential or uniform attachment with choice [15, 7], or random tree branching processes, are worth considering.

# Acknowledgments

The authors thank Elchanan Mossel for suggesting the counterexample to the conjecture that the persistent degree-central node agrees with the persistent centroid in the general preferential attachment model.

### A Weighted 2-dimensional random walks

In this section, we consider a class of random walks on  $\mathbb{N} \times \mathbb{N}$ . From position (i, j), at each time step the random walk can move either one step up with probability U(i, j), or one step to the right with probability R(i, j). The probabilities of these movements depend on (i, j) according to

$$R(i,j) \propto \alpha i + \beta$$
, and  $U(i,j) \propto \alpha j + \beta$ ,

for some  $\alpha > 0$  and  $\alpha + \beta \ge 0$ . Our next lemma pertains to such random walks:

**Lemma A.1.** Consider a 2-dimensional random walk on the  $\mathbb{N} \times \mathbb{N}$  lattice, where the location of the walk at time n is denoted by  $W_n$ , and the probabilities of movement are given by

$$\mathbb{P}(W_{n+1} = (i+1,j)|W_n = (i,j)) = R(i,j) \propto \alpha i + \beta, \text{ and} \\ \mathbb{P}(W_{n+1} = (i,j+1)|W_n = (i,j)) = U(i,j) \propto \alpha j + \beta.$$

For A > 2, let f(A) be the probability that the random walk reaches the diagonal at some future time when it starts at  $W_0 = (A, 1)$ . Then there exists a fixed polynomial P such that

$$f(A) \le \frac{P(A)}{2^A},$$

and  $\frac{P(A)}{2^A}$  is monotonically decreasing for sufficiently large A.

*Proof.* Let f(A, m) be the probability that the random walk lies entirely below the diagonal before reaching (m, m) on the diagonal. Clearly,

$$f(A) = \sum_{m=A}^{\infty} f(A,m).$$

We will now bound f(A, m) for  $m \ge A$ . Let  $\Theta((A, B) \to (m, m))$  be the number of paths from (A, B) to (m, m), such that every point on the path lies strictly below the diagonal, except for the endpoint (m, m). Using the reflection principle, Lemma 2 of Galashin [6] shows that the total number of such paths from (A, B) to (m, m) is given by the expression

$$\Theta((A,B) \to (m,m)) = \frac{(2m-1-A-B)!(A-B)}{(m-A)!(m-B)!}.$$

Substituting B = 1, we have

$$\Theta((A,1) \to (m,m)) = \frac{(2m-2-A)!(A-1)}{(m-A)!(m-1)!} = \frac{\Gamma(2m-A-1)(A-1)}{\Gamma(m+1-A)\Gamma(m)}$$

Furthermore, every path from (A, B) to (m, m) has the same probability  $p((A, B) \to (m, m))$ , and

$$p((A,B) \to (m,m)) = \frac{\prod_{i=A}^{m-1} (\alpha i + \beta) \prod_{j=B}^{m-1} (\alpha j + \beta)}{\prod_{k=A+B}^{2m-1} (\alpha k + 2\beta)}$$
$$= \frac{\alpha^{m-A} \prod_{i=A}^{m-1} (i + \beta/\alpha) \times \alpha^{m-B} \prod_{j=B}^{m-1} (j + \beta/\alpha)}{\alpha^{2m-A-B} \prod_{k=A+B}^{2m-1} (k + 2\beta/\alpha)}$$
$$= \frac{\prod_{i=A}^{m-1} (i + \beta/\alpha) \prod_{j=B}^{m-1} (j + \beta/\alpha)}{\prod_{k=A+B}^{2m-1} (k + 2\beta/\alpha)}.$$

Substituting B = 1 gives

$$p((A,1) \to (m,m)) = \frac{\prod_{i=A}^{m-1} (i+\beta/\alpha) \prod_{j=1}^{m-1} (j+\beta/\alpha)}{\prod_{k=A+1}^{2m-1} (k+2\beta/\alpha)}$$

$$= \frac{\frac{\Gamma(m+\beta/\alpha)}{\Gamma(A+\beta/\alpha)} \frac{\Gamma(m+\beta/\alpha)}{\Gamma(1+\beta/\alpha)}}{\frac{\Gamma(2m+2\beta/\alpha)}{\Gamma(A+1+2\beta/\alpha)}}$$

$$= \frac{\Gamma(A+1+2\beta/\alpha)}{\Gamma(A+\beta/\alpha)\Gamma(1+\beta/\alpha)} \cdot \frac{\Gamma(m+\beta/\alpha)\Gamma(m+\beta/\alpha)}{\Gamma(2m+2\beta/\alpha)}$$

$$\leq \frac{\Gamma(A+1+2\lceil\beta/\alpha\rceil)}{\Gamma(A+\lceil\beta/\alpha\rceil-1)\Gamma(1+\beta/\alpha)} \cdot \frac{\Gamma(m+\beta/\alpha)\Gamma(m+\beta/\alpha)}{\Gamma(2m+2\beta/\alpha)}$$

$$\leq P(A) \cdot \frac{\Gamma^2(m+\beta/\alpha)}{\Gamma(2m+2\beta/\alpha)},$$

where P(A) is a fixed polynomial with degree at most  $\lceil \beta/\alpha \rceil + 2$ . In the rest of this proof, we continue using P(A) to denote a fixed polynomial, although its exact expression may change from line to line. Since

$$f(A,m) = \Theta((A,1) \to (m,m)) \cdot p((A,1) \to (m,m)),$$

we have the bound

$$f(A,m) \le P(A) \frac{\Gamma(2m-A-1)}{\Gamma(m+1-A)\Gamma(m)} \frac{\Gamma^2(m+\beta/\alpha)}{\Gamma(2m+2\beta/\alpha)},$$

for some polynomial P(A). We now use Stirling's bound, which says that for all z > 0, the value of  $\Gamma(z)$  lies with a constant factor of  $\frac{1}{\sqrt{z}} \left(\frac{z}{e}\right)^{z}$ :

$$\Gamma(z) \sim \frac{1}{\sqrt{z}} \left(\frac{z}{e}\right)^z.$$

Substituting and modifying P(A) as convenient, we then have

$$\begin{split} f(A,m) &\leq P(A) \frac{(2m-A-1)^{2m-A-1-1/2}(m+\beta/\alpha)^{2(m+\beta/\alpha-1/2)}}{(m+1-A)^{m+1-A-1/2}m^{m-1/2}(2m+2\beta/\alpha)^{2m+2\beta/\alpha-1/2}} \\ &= P(A) \frac{(2m)^{2m-A-3/2}m^{2m+2\beta/\alpha-1}}{m^{m-A+1/2}m^{m-1/2}(2m)^{2m+2\beta/\alpha-1/2}} \times \frac{(1-\frac{A+1}{2m})^{2m-A-3/2}(1+\beta/m\alpha)^{2m+2\beta/\alpha-1/2}}{(1-\frac{A+1}{2m})^{m-A+1/2}(1+\beta/m\alpha)^{2m+2\beta/\alpha-1/2}} \\ &= \frac{P(A)}{2^{A}m^{2}} \times \frac{(1-\frac{A+1}{2m})^{2m-A-3/2}}{(1-\frac{A-1}{m})^{m-A+1/2}(1+\beta/m\alpha)^{1/2}} \\ &\leq \frac{P(A)}{2^{A}m^{2}} \times \frac{(1-\frac{A+1}{2m})^{2m-A-3/2}}{(1-\frac{A-1}{m})^{m-A+1/2}} \\ &= \frac{P(A)}{2^{A}m^{2}} \times \left(\frac{(1-\frac{A+1}{2m})^{2m-A-3/2}}{(1-\frac{A-1}{m})^{m-A+1/2}} \left(1-\frac{A+1}{2m}\right)^{m-2} \\ &= \frac{P(A)}{2^{A}m^{2}} \times \left(1+\frac{A-3}{2m-2A+2}\right)^{m-A+1/2} \left(1-\frac{A+1}{2m}\right)^{m-2} \\ &\leq \frac{P(A)}{2^{A}m^{2}} \times \exp\left(\frac{A-3}{(2m-2A+2)}(m-A+1/2)-\frac{(A+1)}{2m}(m-2)\right) \\ &\leq \frac{P(A)}{2^{A}m^{2}} \times \exp\left(\frac{A-3}{2}-\frac{A+1}{2}+\frac{A+1}{m}\right) \\ &\leq \frac{P(A)}{2^{A}m^{2}} \times \exp\left(-2+2\right) \\ &= \frac{P(A)}{2^{A}m^{2}}, \end{split}$$

where in (a), we have used the fact that for all  $x \in \mathbb{R}$ , we have  $1 + x \leq e^x$ . Finally, noting that  $\sum_{m=1}^{\infty} \frac{1}{m^2} < \infty$ , we conclude that

$$f(A) = \sum_{m=A}^{\infty} f(A,m) \le \frac{P(A)}{2^A},$$

for a fixed polynomial P. Without loss of generality, we may choose P(A) to be a monomial with a positive coefficient, so  $\frac{P(A)}{2^A}$  is clearly monotonically decreasing for large enough A.

# References

- N. Bailey. The Mathematical Theory of Infectious Diseases and Its Applications. Griffin, London, 1975.
- [2] A.-L. Barabási and R. Albert. Emergence of scaling in random networks. Science, 286(5439):509-512, 1999.
- [3] S. P. Borgatti. Centrality and network flow. Social Networks, 27(1):55–71, 2005.
- [4] S. Bubeck, L. Devroye, and G. Lugosi. Finding Adam in random growing trees. *Random Structures and Algorithms*, page to appear, 2015.
- [5] M. Drmota. Random trees: An interplay between combinatorics and probability. Springer Science & Business Media, 2009.
- [6] P. Galashin. Existence of a persistent hub in the convex preferential attachment model. arXiv preprint arXiv:1310.7513, 2013.
- [7] J. Haslegrave and J. Jordan. Preferential attachment with choice. arXiv preprint arXiv:1407.8421, 2014.
- [8] M. O. Jackson. Social and Economic Networks, volume 3. Princeton University Press, 2008.
- [9] M. O. Jackson and A. Wolinsky. A strategic model of social and economic networks. *Journal of Economic Theory*, 71(1):44 74, 1996.
- [10] C. Jordan. Sur les assemblages de lignes. J. Reine Angew. Math, 70(185):81, 1869.
- [11] A. N. C. Kang and D. A. Ault. Some properties of a centroid of a free tree. Information Processing Letters, 4(1):18–20, 1975.
- [12] J. Khim and P. Loh. Confidence sets for the source of a diffusion in regular trees. ArXiv e-prints, October 2015.
- [13] D. E. Knuth. The Art of Computer Programming, Volume 1 (3rd Ed.): Fundamental Algorithms. Addison Wesley Longman Publishing Co., Inc., Redwood City, CA, USA, 1997.
- [14] S. Laruelle and G. Pagès. Nonlinear randomized urn models: A stochastic approximation viewpoint. arXiv preprint arXiv:1311.7367, 2013.
- [15] Y. Malyshkin and E. Paquette. The power of choice combined with preferential attachment. Electronic Communications in Probability, 19(44):1–13, 2014.
- [16] S. L. Mitchell. Another characterization of the centroid of a tree. Discrete Mathematics, 24(3):277–280, 1978.
- [17] R. Pemantle. A survey of random processes with reinforcement. *Probab. Surveys*, 4:1–79, 2007.
- [18] D. Shah and T. Zaman. Rumors in a network: Who's the culprit? Information Theory, IEEE Transactions on, 57(8):5163-5181, 2011.
- [19] P. J. Slater. Maximin facility location. Journal of National Bureau of Standards B, 79:107–115, 1975.

- [20] P. J. Slater. Accretion centers: A generalization of branch weight centroids. Discrete Applied Mathematics, 3(3):187–192, 1981.
- [21] B. Zelinka. Medians and peripherians of trees. Archivum Mathematicum, 4(2):87–95, 1968.