

# Sequential Estimation of Network Cascades

Anirudh Sridhar

Department of Electrical Engineering  
Princeton University  
Princeton, NJ  
anirudhs@princeton.edu

H. Vincent Poor

Department of Electrical Engineering  
Princeton University  
Princeton, NJ  
poor@princeton.edu

**Abstract**—We consider the problem of locating the source of a network cascade, given a noisy time-series of network data. Initially, the cascade starts with one unknown, affected vertex and spreads deterministically at each time step. The goal is to find an adaptive procedure that outputs an estimate for the source as fast as possible, subject to a bound on the estimation error. For a general class of graphs, we describe a family of matrix sequential probability ratio tests (MSPRTs) that are first-order asymptotically optimal up to a constant factor as the estimation error tends to zero. We apply our results to lattices and regular trees, and show that MSPRTs are asymptotically optimal for regular trees. We support our theoretical results with simulations.

**Index Terms**—Network cascade, sequential estimation, asymptotic optimality, hypothesis testing

## I. INTRODUCTION

Network cascades occur when the behavior of an individual or a small group of individuals diffuses rapidly through a network. Examples include the spread of epidemics in physical or geographical networks [1]–[3], fake news in social networks [4]–[6], and the propagation of viruses in computer networks [7], [8]. In each of these cases, the network cascade compromises the functionality of the network and it is of paramount importance to locate the source of the cascade as fast as possible.

This problem poses several interesting challenges. On one hand, network cascades are typically not directly observable even if one can monitor the network in real time. In the example of an epidemic spreading through a network, an individual’s sickness could be caused by the epidemic or by exogenous factors (e.g., allergies). As the network is monitored over time, one may be able to distinguish between these possibilities at the cost of allowing the cascade to propagate further. Thus there is a fundamental tradeoff between the accuracy of the estimated cascade source and the amount of vertices affected by the cascade. How can we design algorithms that achieve the best possible tradeoff?

In this paper, we take the first steps towards formalizing and solving the challenges addressed above. We begin by reviewing a model for network cascades with real-time noisy observations. We then study the problem of minimizing the expected run time of a sequential estimation algorithm for the cascade source subject to the estimation error being at most  $\alpha$ ,

for some  $\alpha \in (0, 1)$ . We show that simple algorithms based on cumulative log likelihood ratios are first-order asymptotically optimal up to constant factors as we send the estimation error to zero for a large class of networks. In certain cases we can say more: the estimator we construct is first-order optimal in regular trees.

### A. A model of network cascades with noisy observations

Let  $G$  be a graph with vertex set  $V$  and let time be indexed by an integer  $t \geq 0$ . We assume that the cascade starts from some vertex  $v$  such that at  $t = 0$ ,  $v$  is the unique vertex affected. For any  $t \geq 1$ , a vertex  $u$  is affected if and only if  $u \in \mathcal{N}_v(t)$ , where  $\mathcal{N}_v(t)$  denotes the set of all vertices within distance  $t$  from  $v$  in the graph, with respect to the shortest path distance, denoted by  $d(\cdot, \cdot)$ . The cascade is not directly observable, but the system instead monitors *public states*  $\{y_u(t)\}_{u \in V, t \geq 0}$ . Conditioned on the source being  $v$ , the public states are independent over all  $u$  and  $t$ , with distributions given by

$$y_u(t) \sim \begin{cases} Q_0 & u \notin \mathcal{N}_v(t); \\ Q_1 & u \in \mathcal{N}_v(t), \end{cases}$$

where  $Q_0, Q_1$  are two distinct mutually absolutely continuous probability measures over  $\mathbb{R}$ . We can think of  $y_u(t) \sim Q_0$  as typical behavior and  $y_u(t) \sim Q_1$  as anomalous behavior caused by the cascade. This is a standard model which has been previously used, for instance, in studying quickest detection problems on networks [9]–[13].

### B. Formulation as a sequential hypothesis testing problem

Let  $(\Omega, \mathcal{F}, \mathbb{P})$  be a common probability space for all random objects. For each vertex  $u$ , let  $H_u$  be the hypothesis that  $u$  is the cascade source and let  $\mathbb{P}_u := \mathbb{P}(\cdot \mid H_u)$  be the associated measure. Any sequential estimator for the cascade source can be represented by a pair  $(D, T)$ , where  $T$  is a (data dependent) stopping time and  $D = \{D(t)\}_{t=0}^\infty$  is a sequence of estimators such that  $D(t) \in V$  depends on the observations  $y(0), \dots, y(t)$ . The output of the sequential estimation procedure is  $D(T)$ . Given a positive integer  $R$  (which we will call the *confidence radius*), we say that the sequential estimator  $(D, T)$  *succeeds* if  $d(D(T), v) \leq R$ ; else it fails. We can then formalize the tradeoff between estimator accuracy and number of infections as follows. Given  $\alpha \in (0, 1)$ , we want to find the estimator which minimizes

the worst-case expected runtime  $\max_{v \in V} \mathbb{E}_v[T]$ , subject to the probability of failure being at most  $\alpha$ .

A typical assumption in source estimation problems is that the graph  $G$  has infinitely many vertices, is connected, and is locally finite<sup>1</sup> (see [14]). However, the infinite graph setting corresponds to testing infinitely many hypotheses, and it is unclear whether there exists a sequential estimator with small error that will terminate in finite time. To remedy this situation, we will consider the behavior of sequential tests on finite restrictions of the graph. Formally, we fix a sequence  $\{V_n\}_n$  where  $V_n \subset V$  and  $|V_n| = n$ . In our analysis we will specify  $n$  and assume that  $D(t) \in V_n$ ; that is, the true source is an element of  $V_n$ . We will assume that  $V_n$  is a neighborhood of a given vertex  $v_0$  without loss of generality as the problem is only harder when all vertices are close to each other.

We may now put these ideas together more formally to define a class of sequential estimators  $\Delta_G(V_n, R, \alpha)$  for which the probability of failure is at most  $\alpha$ . Equations (1) and (2) below present two natural formulations for this set.

$$\Delta'_G := \{(D, T) : \forall v \in V_n, \mathbb{P}_v(d(D(T), v) > R) \leq \alpha\} \quad (1)$$

$$\Delta_G := \left\{ (D, T) : \max_{\substack{u, v \in V_n: \\ d(u, v) > R}} \mathbb{P}_v(D(T) = u) \leq \frac{\alpha}{n} \right\}. \quad (2)$$

Formulation (1) directly bounds the probability of failure, while formulation (2) provides finer information on the distribution of the estimator outside of the confidence radius, and is thus more mathematically convenient. It is easy to see that  $\Delta_G \subset \Delta'_G$ . For any  $(D, T) \in \Delta_G$  and for any  $v \in V_n$ ,

$$\mathbb{P}_v(d(D(T), v) > R) = \sum_{u \in V_n \setminus \mathcal{N}_v(R)} \mathbb{P}_v(D(T) = u) \leq \alpha. \quad (3)$$

In other words, (3) shows that Formulation (2) is stronger than Formulation (1). We believe that in certain cases, the two formulations are equivalent if  $R$  is sufficiently large and  $G$  satisfies some symmetry properties (e.g., vertex-transitivity). While we use (2) in this paper, we will study the relationship between (1) and (2) in future work.

Putting everything together, our goal will be to characterize

$$T^*(V_n, R_n, \alpha) := \min_{(D, T) \in \Delta_G(V_n, R_n, \alpha)} \max_{v \in V_n} \mathbb{E}_v[T]. \quad (4)$$

For a fixed sequential estimator  $(D, T)$ , the inner maximum is the worst-case expected runtime of the estimator over all possible sources. The outer minimum is over all  $(D, T)$  which meet the requirements of  $\Delta_G(V_n, R_n, \alpha)$ . In general sequential multi-hypothesis testing problems, characterizing the optimal test is intractable. We therefore study the asymptotics of (4) first as  $n \rightarrow \infty$  then as  $\alpha \rightarrow 0$ . We consider confidence radii  $R_n$  that may be fixed with respect to all other parameters, or may grow with  $n$ .

<sup>1</sup>A graph is locally finite if the degree of each vertex is finite.

## C. Related work

Although we are, to the best of our knowledge, the first to study this variant of the cascade source estimation problem, our work has close connections to several bodies of work.

Shah and Zaman gave the first systematic study of estimating the source of a network cascade [14], [15], which spawned several follow-up works, see for example [16]–[21]. In their setup, they assume that the network cascade evolves according to a probabilistic model, and that at some future time a snapshot of the cascade is perfectly observed. The problems we address surrounding network cascades are complementary to this approach, and are more appropriate for the setting where one may monitor the state of the network in real time.

Our work falls under the growing body of literature on sequential detection and estimation in networks. Recently Zou, Veeravalli, Li, Towsley and Rovatsos studied the problem of quickest detection of a network cascade [9]–[12], which is similar in nature to our work. The objective of their work is to detect with minimum delay when a cascade has started propagating in a network. They derive tests based on cumulative log-likelihood ratios that are shown to be asymptotically optimal when the growth rate of the cascade becomes very small. The detection problem with other cascade models were studied by Zhang, Yao, Xie and Qiu [13]. Our work, on the other hand, studies problems of estimation. We assume simple cascade dynamics for ease of exposition, and extensions to other cascade dynamics is an important future direction.

## II. ASYMPTOTIC BEHAVIOR OF OPTIMAL TESTS

At the core of the analysis is a characterization of the rate of convergence of the log-likelihood ratios. It is well known that tests based on log-likelihood ratios are optimal for distinguishing between two hypotheses [22] and are asymptotically optimal as the Type I error tends to 0 in the general multi-hypothesis testing problem [23]–[25]. Thus, to motivate our results, we begin by studying some basic properties of the log-likelihood ratios that arise from our problem structure.

For a positive integer  $s$ , define the shorthand  $y(s) := \{y_u(s)\}_{u \in V}$  to be the collection of public states at time  $s$ . We are interested in the cumulative log-likelihood ratio

$$Z_{vu}(t) := \sum_{s=0}^t \log \frac{d\mathbb{P}_v}{d\mathbb{P}_u}(y(s)).$$

From the cascade dynamics defined in Section I-A, we can write the log-likelihood ratio  $\log \frac{d\mathbb{P}_v}{d\mathbb{P}_u}(y(s))$  as

$$\begin{aligned} & \log \frac{\prod_{w \in \mathcal{N}_v(s)} dQ_1(y_w(s)) \cdot \prod_{w \notin \mathcal{N}_v(s)} dQ_0(y_w(s))}{\prod_{w \in \mathcal{N}_u(s)} dQ_1(y_w(s)) \cdot \prod_{w \notin \mathcal{N}_u(s)} dQ_0(y_w(s))} \\ &= \sum_{w \in \mathcal{N}_v(s) \setminus \mathcal{N}_u(s)} \log \frac{dQ_1}{dQ_0}(y_w(s)) + \sum_{w \in \mathcal{N}_u(s) \setminus \mathcal{N}_v(s)} \log \frac{dQ_0}{dQ_1}(y_w(s)). \end{aligned} \quad (5)$$

Under  $\mathbb{P}_v$ , all observed variables  $y_w(s)$  are independent, with distributions given by

$$y_w(s) \sim \begin{cases} Q_0 & w \in \mathcal{N}_u(s) \setminus \mathcal{N}_v(s); \\ Q_1 & w \in \mathcal{N}_v(s) \setminus \mathcal{N}_u(s). \end{cases} \quad (6)$$

To simplify the analysis we assume that for each  $u, v \in V$  and  $t \geq 0$ ,  $\mathcal{N}_v(t) \setminus \mathcal{N}_u(t)$  is nonempty, and  $|\mathcal{N}_v(t) \setminus \mathcal{N}_u(t)| = |\mathcal{N}_u(t) \setminus \mathcal{N}_v(t)|$ . This assumption holds for a large class of graphs (e.g., vertex-transitive graphs such as regular trees and lattices). For  $u, v \in V$  define

$$f_{vu}(t) := \sum_{s=0}^t |\mathcal{N}_v(s) \setminus \mathcal{N}_u(s)|.$$

Equations (5) and (6) imply  $\mathbb{E}_v[Z_{vu}(t)] = \tilde{D}(Q_0, Q_1)f_{vu}(t)$ , where  $\tilde{D}(Q_0, Q_1)$  is the symmetrized Kullback-Leibler divergence between  $Q_0$  and  $Q_1$ , given explicitly by

$$\tilde{D}(Q_0, Q_1) := \int \log \left( \frac{dQ_1}{dQ_0} \right) dQ_1 + \int \log \left( \frac{dQ_0}{dQ_1} \right) dQ_0.$$

Before stating our first main result on lower bounds for  $T^*$ , we will introduce some useful asymptotic notation. For two functions  $g(n, \alpha)$  and  $h(n, \alpha)$ , we write, whenever it is well-defined,

$$g(n, \alpha) \gtrsim_{n, \alpha} h(n, \alpha) \iff \liminf_{\alpha \rightarrow 0} \liminf_{n \rightarrow \infty} \frac{g(n, \alpha)}{h(n, \alpha)} \geq 1.$$

The orderwise comparison operator  $\lesssim_{n, \alpha}$  is analogously defined, but with limsups instead of liminfs. We also write  $g(n, \alpha) \approx_{n, \alpha} h(n, \alpha)$  if and only if  $g(n, \alpha) \gtrsim_{n, \alpha} h(n, \alpha)$  and  $g(n, \alpha) \lesssim_{n, \alpha} h(n, \alpha)$ .

**Theorem 1.** *Let  $F_{vu}$  be the inverse function of  $f_{vu}$  and suppose  $\log |\mathcal{N}_v(R_n)| \ll \log n$ . Then*

$$T^*(V_n, R_n, \alpha) \gtrsim_{n, \alpha} \max_{u, v \in V_n: d(u, v) > 2R_n} F_{vu} \left( \frac{\log n / \alpha}{\tilde{D}(Q_0, Q_1)} \right).$$

*Proof of Theorem 1.* We briefly discuss the high-level proof strategy. The term on the right hand side involving  $F_{vu}$  is the time it takes for  $\mathbb{E}_v[Z_{vu}(t)]$  to cross the threshold  $\log n / \alpha$ , since  $F_{vu}$  is the inverse function of the log-likelihood growth rate  $f_{vu}$ . Using a change of measure argument, we show  $\mathbb{E}_v[Z_{vu}(t)]$  must cross this threshold to achieve the guarantees of  $\Delta_G(V_n, R_n, \alpha)$ .

To show this formally, fix any  $(D, T) \in \Delta_G(V_n, R_n, \alpha)$  as well as a vertex  $v \in V_n$ . Define the event  $\Omega_{v, L} := \{D(T) \in \mathcal{N}_v(R_n)\} \cap \{T \leq L\}$ , where  $L$  is a positive integer to be chosen later. By a change of measure,

$$\mathbb{P}_u(D(T) \in \mathcal{N}_v(R_n)) = \mathbb{E}_v \left[ \mathbb{1}_{\{D(T) \in \mathcal{N}_v(R_n)\}} e^{-Z_{vu}(T)} \right].$$

For any positive integer  $B$ ,

$$\begin{aligned} \mathbb{P}_u(D(T) \in \mathcal{N}_v(R_n)) &\geq \mathbb{E}_v \left[ \mathbb{1}_{\Omega_{v, L} \cap \{Z_{vu}(T) < B\}} e^{-Z_{vu}(T)} \right] \\ &\geq e^{-B} \mathbb{P}_v \left( \Omega_{v, L} \cap \left\{ \max_{t \leq L} Z_{vu}(t) < B \right\} \right) \\ &\geq e^{-B} \left( \mathbb{P}_v(\Omega_{v, L}) - \mathbb{P}_v \left( \max_{t \leq L} Z_{vu}(t) \geq B \right) \right) \end{aligned}$$

Noting that  $\mathbb{P}_v(\Omega_{v, L}) \geq \mathbb{P}_v(D(T) \in \mathcal{N}_v(R_n)) - \mathbb{P}_v(T > L)$  and substituting this in place of  $\mathbb{P}_v(\Omega_{v, L})$  gives

$$\begin{aligned} \mathbb{P}_v(T > L) &\geq \mathbb{P}_v(D(T) \in \mathcal{N}_v(R_n)) \\ &\quad - e^B \mathbb{P}_u(D(T) \in \mathcal{N}_v(R_n)) - \mathbb{P}_v \left( \max_{t \leq L} Z_{vu}(t) \geq B \right) \end{aligned}$$

From (2),  $\mathbb{P}_v(D(T) \in \mathcal{N}_v(R_n)) \geq 1 - \alpha$  and  $\mathbb{P}_u(D(T) \in \mathcal{N}_v(R_n)) \leq \frac{\alpha |\mathcal{N}_v(R_n)|}{n}$  for  $u \in V_n \setminus \mathcal{N}_v(2R_n)$ . It follows that

$$\begin{aligned} \mathbb{P}_v(T > L) &\geq 1 - \alpha \\ &\quad - e^B \cdot \frac{\alpha |\mathcal{N}_v(R_n)|}{n} - \mathbb{P}_v \left( \max_{t \leq L} Z_{vu}(t) \geq B \right). \end{aligned} \quad (7)$$

Let  $\epsilon > 0$ , and set

$$\begin{aligned} L &:= F_{vu} \left( \frac{1 - \epsilon}{\tilde{D}(Q_0, Q_1) + \epsilon} \log \frac{n}{\alpha |\mathcal{N}_v(R_n)|} \right) \\ B &:= (1 - \epsilon) \log \frac{n}{\alpha |\mathcal{N}_v(R_n)|}. \end{aligned}$$

Then by the strong law of large numbers (see Lemma 2.1 in [23]),  $\mathbb{P}_v(\max_{t \leq L} Z_{vu}(t) \geq B)$  goes to 0 as  $n \rightarrow \infty$ . Plugging in these values to (7) gives the following lower bound on  $\mathbb{P}_v(T > L)$ :

$$1 - \alpha - \left( \frac{\alpha |\mathcal{N}_v(R_n)|}{n} \right)^\epsilon - \mathbb{P}_v \left( \max_{t \leq L} Z_{vu}(t) \geq B \right).$$

Take  $n \rightarrow \infty$  and then  $\alpha \rightarrow 0$  to obtain

$$\liminf_{\alpha \rightarrow 0} \liminf_{n \rightarrow \infty} \mathbb{P}_v \left( T > F_{vu} \left( \frac{(1 - \epsilon) \log \frac{n}{\alpha |\mathcal{N}_v(R_n)|}}{\tilde{D}(Q_0, Q_1) + \epsilon} \right) \right) \geq 1.$$

We conclude from Markov's inequality and by considering only the first-order terms (see Lemma 2.1 in [23]).  $\square$

Next we establish an upper bound on  $T^*$  by considering families of matrix sequential probability ratio tests (MSPRTs) [23]–[25]. Define the stopping time

$$T_v := \min \left\{ t \geq 0 : \min_{u \in V_n \setminus \mathcal{N}_v(R_n)} Z_{vu}(t) \geq \log \frac{n}{\alpha} \right\},$$

and define the pair  $(D_n, T_n)$  via

$$T_n := \min_{v \in V_n} T_v, \quad D_n(t) := \arg \max_{v \in V_n} \min_{u \in V_n \setminus \mathcal{N}_v(R_n)} Z_{vu}(t)$$

so that  $D_n(T_n) := \arg \min_{v \in V_n} T_v$ . It is simple to verify  $(D_n, T_n) \in \Delta_G(V_n, R_n, \alpha)$ . For  $u \in V_n \setminus \mathcal{N}_v(R_n)$ ,

$$\begin{aligned} \mathbb{P}_v(D_n(T_n) = u) &\leq \mathbb{P}_v(T_u < \infty) = \mathbb{E}_u \left[ \mathbb{1}_{\{T_u < \infty\}} e^{-Z_{uv}(T_u)} \right] \\ &= e^{-\log n / \alpha} \mathbb{E}_u \left[ \mathbb{1}_{\{T_u < \infty\}} e^{-(Z_{uv}(T_u) - \log n / \alpha)} \right]. \end{aligned}$$

Since  $Z_{vu}(T_u) \geq \log n / \alpha$  by the definition of  $T_u$ , we have an upper bound of  $\alpha / n$  as desired. The following theorem gives us an upper bound for the expectation of  $T_n$ .

**Theorem 2.** *Let  $\alpha \in (0, 1)$  be fixed. There exists a constant  $C(Q_0, Q_1) \in (0, \tilde{D}(Q_0, Q_1))$  such that for every  $v \in V_n$ ,*

$$\mathbb{E}_v[T_n] \leq \max_{u \in V_n: d(u, v) > R_n} F_{vu} \left( \frac{\log n}{C(Q_0, Q_1)} \right) (1 + o_n(1)),$$

where  $o_n(1) \rightarrow 0$  as  $n \rightarrow \infty$ .

*Proof.* We begin by upper bounding  $\mathbb{P}_v(T_v > t)$ . We can write

$$\begin{aligned} \mathbb{P}_v(T_v > t) &\leq \mathbb{P}_v\left(\min_{u \in V_n \setminus \mathcal{N}_v(R_n)} Z_{vu}(t) < \log \frac{n}{\alpha}\right) \\ &\leq \sum_{u \in V_n \setminus \mathcal{N}_v(R_n)} \mathbb{P}_v\left(Z_{vu}(t) < \log \frac{n}{\alpha}\right). \end{aligned} \quad (8)$$

Fix  $\epsilon > 0$  and suppose that, for all  $u \in V_n \setminus \mathcal{N}_v(R_n)$ ,

$$\log n / \alpha \leq (\tilde{D}(Q_0, Q_1) - \epsilon) f_{vu}(t).$$

To deal with the terms in the summation, we will use Chernoff bounds for  $Z_{vu}(t)$ . Suppose that  $X, Y$  are independent random variables with  $X \sim Q_1$  and  $Y \sim Q_0$ , and define

$$I(x) := \sup_{\lambda \geq 0} \left\{ \lambda x + \log \mathbb{E} \left[ \left( \frac{dQ_0}{dQ_1}(X) \right)^\lambda \left( \frac{dQ_1}{dQ_0}(Y) \right)^\lambda \right] \right\}.$$

Furthermore, we have the strict inequality  $I(x) > 0$  for  $x \leq \tilde{D}(Q_0, Q_1)$  since the  $\lambda$ -moments of  $\frac{dQ_0}{dQ_1}(X)$  and  $\frac{dQ_1}{dQ_0}(Y)$  exist for  $\lambda \in [0, 1]$ . Since  $Z_{vu}(t)$  can be written as a sum of i.i.d. random variables under  $\mathbb{P}_v$ , we have the bound

$$\mathbb{P}_v(Z_{vu}(t) \leq x f_{vu}(t)) \leq e^{-f_{vu}(t) I(x)} \text{ for } x \leq \tilde{D}(Q_0, Q_1).$$

Hence we can bound the summation (8) by

$$\exp \left( \log n - \min_{u \in V_n \setminus \mathcal{N}_v(R_n)} f_{vu}(t) I(\tilde{D}(Q_0, Q_1) - \epsilon) \right). \quad (9)$$

Set  $C(Q_0, Q_1) := \min\{\tilde{D}(Q_0, Q_1) - \epsilon, I(\tilde{D}(Q_0, Q_1) - \epsilon)\}$  and define

$$t_{n,\epsilon} := \max_{u \in V_n: d(u,v) > R_n} F_{vu} \left( \frac{\log n / \alpha}{C(Q_0, Q_1)} \right).$$

Next, write  $\mathbb{E}_v[T_v] = \sum_{t=0}^{\infty} \mathbb{P}_v(T_v > t)$  and apply (9) to bound  $\mathbb{E}_v[T_v]$  by

$$t_{n,\epsilon} + \sum_{t=t_{n,\epsilon}+1}^{\infty} \exp \left( \log n - \min_{\substack{u \in V_n: \\ d(u,v) > R_n}} f_{vu}(t) I(\tilde{D}(Q_0, Q_1) - \epsilon) \right)$$

which is  $t_{n,\epsilon}(1 + o_n(1))$  where  $o_n(1) \rightarrow 0$  as  $n \rightarrow \infty$ . The desired result follows from considering only first-order terms.  $\square$

### III. APPLICATIONS TO REGULAR TREES AND LATTICES

To apply Theorems 1 and 2, it suffices to compute  $F_{vu}$  for the graphs of interest. The following result shows that  $(D_n, T_n)$  is asymptotically optimal as  $n \rightarrow \infty$  and  $\alpha \rightarrow 0$  in regular trees.

**Corollary 1.** *Let  $G$  be the infinite  $k$ -regular tree. If  $R_n \ll \log n$ , then the MSPRT is asymptotically optimal, and*

$$T^*(V_n, R_n, \alpha) \approx_{n,\alpha} \frac{\log \log n}{\log(k-1)}.$$

The idea behind the proof is that for  $k \geq 3$ ,  $f_{vu}(t) \asymp (k-1)^t$  and  $F_{vu}(z) \sim \frac{\log z}{\log(k-1)}$ . This follows from simple counting arguments so we do not include it here. In particular,

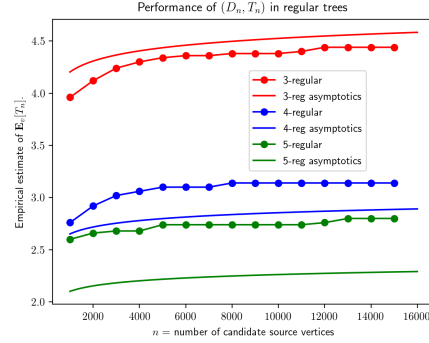


Fig. 1. Numerical results for the MSPRT in regular trees for  $R = 0$ .

the first-order asymptotics of  $F_{vu}$  do not depend on  $d(u, v)$ ; hence the first-order behavior of  $T^*$  does not depend on  $R$ .

To generate the numerical results in Figure 1, we let  $G$  be a balanced,  $k$ -regular tree with height  $h_k$  and root vertex  $v_0$ , where we set  $h_3 = 15$  (32,767 vertices),  $h_4 = 11$  (29,524 vertices) and  $h_5 = 9$  (87,381 vertices). In each case we enumerated the vertices by positive integers so that if vertex  $u$  is assigned a larger number than vertex  $v$ , then  $d(u, v_0) \geq d(v, v_0)$ . We set  $V_n$  to be the set of vertices with label at most  $n$ , and allowed  $n$  to range between 1,000 and 16,000. We set the root vertex  $v_0$  to be the cascade source and also set  $\alpha = 0.1$ ,  $R = 0$ ,  $Q_0 \equiv N(0, 1)$  and  $Q_1 \equiv N(2, 1)$ . An estimate of  $\mathbb{E}_{v_0}[T_n]$  was obtained by averaging over 50 simulations per data point.

Next we consider optimal source estimation in lattices. We establish rigorous results in the one-dimensional case and discuss how things change in higher dimensions. Interestingly, the performance of the optimal estimator depends strongly on the confidence radius.

**Corollary 2.** *Let  $G$  be the infinite line graph. If  $R_n \ll n$ , then  $(D_n, T_n)$  is asymptotically optimal up to a constant factor depending on  $Q_0$  and  $Q_1$ . Let  $C(Q_0, Q_1)$  be the constant from Theorem 2. If  $R_n \ll \sqrt{\log n}$ ,*

$$\frac{\log n}{2R_n \tilde{D}(Q_0, Q_1)} \lesssim_{n,\alpha} T^*(V_n, R_n, \alpha) \lesssim_{n,\alpha} \frac{\log n}{R_n C(Q_0, Q_1)}.$$

*If  $R_n \gg \sqrt{\log n}$  and  $\log R_n \ll \log n$ ,*

$$\sqrt{\frac{\log n}{\tilde{D}(Q_0, Q_1)}} \lesssim_{n,\alpha} T^*(V_n, R_n, \alpha) \lesssim_{n,\alpha} \sqrt{\frac{\log n}{C(Q_0, Q_1)}}.$$

We sketch the derivation of  $F_{vu}$ . If  $r = d(u, v)$  is even, then simple counting arguments show that

$$f_{vu}(t) = \begin{cases} (t+1)^2 & t < \frac{r}{2}; \\ \frac{r}{2} (2t+2 - \frac{r}{2}) & t \geq \frac{r}{2}. \end{cases} \quad (10)$$

Let  $\lambda > 0$  be any constant. As  $n \rightarrow \infty$ , (10) implies

$$F_{vu}(\lambda \log n) \sim \begin{cases} \frac{\lambda \log n}{R_n} & R_n \ll \sqrt{\log n} \\ \sqrt{\lambda \log n} & R_n \gg \sqrt{\log n}. \end{cases} \quad (11)$$

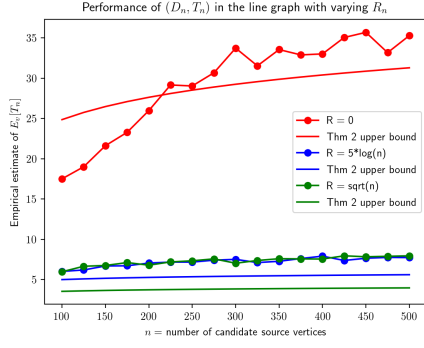


Fig. 2. Numerical results for the MSPRT in the infinite line graph for  $R_n = 0, 5 \log n, \sqrt{n}$ .

Applying Theorems 1 and 2 proves the corollary.

To generate the numerical results in Figure 2 we let  $G$  be a line graph with 1,000 vertices where we enumerate vertices from left to right from 1 to 1,000. We set  $V_n$  to be the set of  $n$  vertices closest to the vertex labelled 500 (where we chose odd values of  $n$ ) and we let  $n$  vary between 25 and 500. We studied the performance of the MSPRT when the cascade source is the vertex labelled 500 and set  $\alpha = 0.2, Q_0 \equiv N(0, 1)$  and  $Q_1 \equiv N(0.5, 1)$ . To estimate  $\mathbb{E}_{500}[T_n]$  we average each data point over 50 simulations. We plot our results for  $R = 0, 5 \log n, \sqrt{n}$ .

Computing  $f_{vu}$  and  $F_{vu}$  in higher-dimensional lattices requires more involved combinatorial arguments, but Corollary 2 gives us a strong intuition of what to expect. The size of  $|\mathcal{N}_v(s)|$  in the  $k$ -dimensional lattice is of order  $s^k$ , so for  $t < \frac{d(u,v)}{2}$ ,  $f_{vu}(t)$  should be on the order of  $t^{k+1}$ . Inverting  $f_{vu}$ , we expect that  $T^*(n, R_n, \alpha)$  will increase as  $(\log n)^{\frac{1}{k+1}}$  when  $R_n \gg (\log n)^{\frac{1}{k+1}}$ .

#### IV. CONCLUSION

In this paper, we studied the problem of estimating the source of a network cascade given noisy time-series data. We found that if  $\min_{v \in V_n} |\mathcal{N}_v(R_n)| \ll n$ , the MSPRT is asymptotically optimal as  $\alpha \rightarrow 0$  in the case of regular trees, and is asymptotically optimal up to a constant factor in general. We have discussed many avenues for future work, including a study of *non-asymptotic* optimality, closing the gap between the upper and lower bounds of Theorems 1 and 2, and investigating the relationship between the formulations (1) and (2) for  $\Delta_G$ .

#### REFERENCES

- [1] N. A. Christakis and J. H. Fowler, "Social network sensors for early detection of contagious outbreaks," *PLOS ONE*, vol. 5, no. 9, pp. 1–8, Sept 2010.
- [2] F. Pervaiz, M. Pervaiz, N. Rehman, and U. Saif, "Flubreaks: Early epidemic detection from google flu trends," *Journal of Medical Internet Research*, vol. 14, p. 125, Oct 2012.
- [3] N. Antulov-Fantulin, A. Lančić, T. Šmuc, H. Štefančić, and M. Šikić, "Identification of patient zero in static and temporal networks: Robustness and limitations," *Phys. Rev. Lett.*, vol. 114, p. 248701, Jun 2015.
- [4] E. Tacchini, G. Ballarin, M. L. Della Vedova, S. Moret, and L. de Alfaro, "Some Like it Hoax: Automated Fake News Detection in Social Networks," in *2nd Workshop on Data Science for Social Good*, 2017, pp. 1–15.
- [5] H. Allcott and M. Gentzkow, "Social Media and Fake News in the 2016 Election," *Journal of Economic Perspectives*, vol. 31, no. 2, pp. 211–236, May 2017.
- [6] A. Fournay, M. Z. Rácz, G. Ranade, M. Mobius, and E. Horvitz, "Geographic and temporal trends in fake news consumption during the 2016 us presidential election," in *Proceedings of the 2017 ACM on Conference on Information and Knowledge Management*, ser. CIKM 17. New York, NY, USA: Association for Computing Machinery, 2017, pp. 2071–2074.
- [7] J. O. Kephart and S. R. White, "Directed-graph epidemiological models of computer viruses," in *Proceedings. 1991 IEEE Computer Society Symposium on Research in Security and Privacy*, May 1991, pp. 343–359.
- [8] G. A. N. Mohamed and N. Ithnin, "Survey on representation techniques for malware detection system," *American Journal of Applied Sciences*, vol. 14, pp. 1049–1069, Nov 2017.
- [9] S. Zou and V. V. Veeravalli, "Quickest detection of dynamic events in sensor networks," in *2018 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, April 2018, pp. 6907–6911.
- [10] S. Zou, V. V. Veeravalli, J. Li, and D. Towsley, "Quickest detection of significant events in structured networks," in *2018 52nd Asilomar Conference on Signals, Systems, and Computers*, Oct 2018, pp. 1307–1311.
- [11] —, "Quickest detection of dynamic events in networks," *IEEE Transactions on Information Theory*, pp. 1–1, 2019.
- [12] G. Rovatsos, V. V. Veeravalli, D. Towsley, and A. Swami, "Quickest Detection of Growing Dynamic Anomalies in Networks," *arXiv e-prints*, p. arXiv:1910.09151, Oct 2019.
- [13] R. Zhang, R. Yao, Y. Xie, and F. Qiu, "Quickest detection of cascading failure," *arXiv e-prints*, p. arXiv:1911.05610, Oct 2019.
- [14] D. Shah and T. Zaman, "Rumors in a Network: Who's the Culprit?" *IEEE Transactions on Information Theory*, vol. 57, no. 8, pp. 5163–5181, 2011.
- [15] —, "Detecting Sources of Computer Viruses in Networks: Theory and Experiment," in *ACM SIGMETRICS*, vol. 38, 2010, pp. 203–214.
- [16] —, "Finding rumor sources on random trees," *Operations Research*, vol. 64, no. 3, pp. 736–755, 2016.
- [17] J. Khim and P.-L. Loh, "Confidence Sets for the Source of a Diffusion in Regular Trees," *IEEE Transactions on Network Science and Engineering*, vol. 4, no. 1, pp. 27–40, 2017.
- [18] A. Y. Lokhov, M. Mézard, H. Ohta, and L. Zdeborová, "Inferring the origin of an epidemic with a dynamic message-passing algorithm," *Phys. Rev. E*, vol. 90, p. 012801, Jul 2014.
- [19] W. Luo, W.-P. Tay, and M. Leng, "Identifying infection sources and regions in large networks," *IEEE Transactions on Signal Processing*, vol. 61, pp. 2850–2865, 2013.
- [20] Z. Wang, W. Dong, W. Zhang, and C. W. Tan, "Rumor Source Detection with Multiple Observations: Fundamental Limits and Algorithms," in *ACM SIGMETRICS Performance Evaluation Review*, vol. 42, 2014, pp. 1–13.
- [21] S. Feizi, K. Duffy, M. Kellis, and M. Mardar, "Network infusion to infer information sources in networks," *IEEE Transactions on Network Science and Engineering*, vol. PP, Jan 2015.
- [22] A. Wald and J. Wolfowitz, "Optimum character of the sequential probability ratio test," *Ann. Math. Statist.*, vol. 19, no. 3, pp. 326–339, Sept 1948.
- [23] A. G. Tartakovsky, "Asymptotic optimality of certain multihypothesis sequential tests: Non i.i.d. case," *Statistical Inference for Stochastic Processes*, vol. 1, no. 3, pp. 265–295, Oct 1998.
- [24] C. W. Baum and V. V. Veeravalli, "A sequential procedure for multihypothesis testing," *IEEE Transactions on Information Theory*, vol. 40, no. 6, pp. 1994–2007, Nov 1994.
- [25] V. V. Veeravalli and C. W. Baum, "Asymptotic efficiency of a sequential multihypothesis test," *IEEE Transactions on Information Theory*, vol. 41, no. 6, pp. 1994–1997, Nov 1995.