

Modeling outcasting behaviour using mean field theory

Utkarsh Upadhyay, Jean-Yves Le Boudec, Rachid Guerraoui

August 18, 2011

Abstract

We model *outcasting* phenomenon in a completely connected network of nodes communicating via a gossip-based protocol. An *outcast* is defined as a node which has been blacklisted by a certain fraction of nodes in the system. We show how to calculate the probability of such a node receiving a packet disseminated in the network. We do it by approximating the evolution of the system with ODEs inspired by the SIR model of epidemic spread. We theoretically bound the probability of the actual system deviating from its approximation using results from mean field theory. We show that as the number of nodes in the network increases, it becomes difficult to outcast misbehaving nodes by distributed blacklisting while maintaining the same quality of service.



ÉCOLE POLYTECHNIQUE
FÉDÉRALE DE LAUSANNE

Page intentionally left blank

1 Introduction

Gossip-based protocols are a viable method of communication in large scale networks with frequently changing topologies. Out of three primary applications of gossip based protocols discussed in [2], the application we are interested in is push based broadcasting of packets on a network with probabilistic guarantees on packet reception [5]. In these algorithms, each node after receiving a packet, forwards it to f other randomly sampled peers [8], where f maybe fixed or dependent on the resources available to the node [6]. It is known that the expected value of f must grow logarithmically with the increasing network size to preserve the quality of service and guarantees provided by the protocol [5]. In the present analysis, we consider f to be a fixed constant.

In such algorithms, it may happen that a fraction of nodes do not communicate with a given node (called the *outcast*) in the network. This could be the result of *blacklisting* due to misbehaviour. We analyse the probability of packet reception by the outcast as a function of the fraction of nodes which stop communicating with it.

This analysis is useful in designing protocols for punishing *freeriding* nodes in a network. For example, if a live stream is being disseminated on the network with 10% Forward Error Correction, then the analysis would allow us to determine the fraction of nodes which should stop communicating with an outcast so that it receives less than 90% of the packets. This would make it impossible for the outcast to reconstruct the stream. Currently, these blacklisting mechanisms are centralized [7]. We investigate the feasibility of doing this blacklisting in a distributed manner, where each node makes the decision of blacklisting an outcast for itself.

We show using our approximate model that the fraction of nodes which must blacklist the outcast becomes prohibitively large as the number of nodes in the system increases, that is, it is difficult to *remove* a node from the network by distributed blacklisting once it has been inducted.

Related works: In [9], the authors introduce and analyse Markov Decision Evolutionary Games and show that as the number of players (nodes) in the system grows to infinity, the random process consisting of one individual player and the remaining population converges weakly to a jump process driven by the solution of a system of differential equations. We instead deal with a situation in which one individual is dealt with in different ways by the two separate fraction of nodes.

In [1], authors investigate an abstraction, which they call mean field method, for performance evaluation of dynamic networks with pairwise communication between nodes. They aim to automate the evaluation by abstracting the topology of the underlying network by defining different *classes* of states. In our work, we concentrate on a similar distinction in topology (with some nodes completely connected and other nodes not connected with the outcast). However, nodes in our network communicate via a gossip based protocol and we are interested in the interactions of a class of nodes (outcast) which is of size 0 in the mean field limit.

Structure: Section 2 describes the exact system and how it is modeled. The model is analysed in Section 3 and the analysis is verified with the help of simulations in Section 4. This is followed by a short discussion and concluding remarks.

2 System description and modeling

The system under consideration is a network of $N + 1$ nodes which communicate with each other using a gossip-based protocol. Initially, only the source node has the packet. It selects f random nodes without replacement in the network and forwards them the packet. This could be done by maintaining a partial view of the network [8]. Then each node which has received the packet in turn chooses f random nodes to forward the packet to and so on. Every node which receives the packet forwards it only once irrespective of how many copies it receives. The evolution of the system happens asynchronously.

The propagation of one packet of data is identical to the spread of an epidemic in the SIR model: the nodes which have not received a packet are the susceptible nodes (set **S**), the nodes which have the packet are the infected nodes (set **I**), packet communication is *infection*, and the nodes which have duly forwarded the packet are the recovered nodes (**R** set). See Figure 1 for two stages during the spread of the packet. This vocabulary would be adopted for the ensuing discussion.

When a node is outcast, the remaining N nodes of the system can be divided into two sets \mathbf{Z}_1 and \mathbf{Z}_2 , with the nodes which can communicate with the complete network in \mathbf{Z}_1 and the nodes

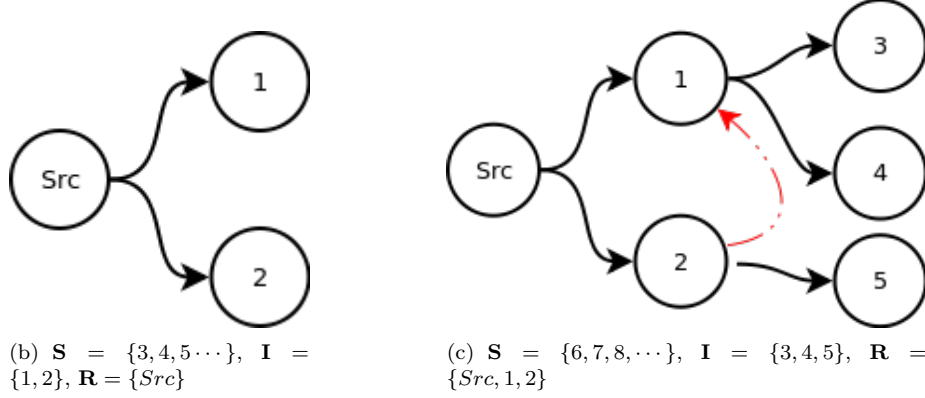


Figure 1: Two rounds of gossiping with random peer selection

which have blacklisted the outcast in \mathbf{Z}_2 . Now we define 5 different subsets of nodes (instead of 3 set of nodes for the SIR model) of \mathbf{Z}_1 and \mathbf{Z}_2 . Since we are not concerned about the individual nodes in the system, we keep track of only the sizes of each set. The description of each set is given in Table 1.

Set descriptor	Description
I_1	Number of nodes which are <i>infected</i> in \mathbf{Z}_1
S_1	Number of nodes which are <i>susceptible</i> in \mathbf{Z}_1
I_2	Number of nodes which are <i>infected</i> in \mathbf{Z}_2
S_2	Number of nodes which are <i>susceptible</i> in \mathbf{Z}_2
R	Number of nodes which have been <i>removed</i>

Table 1: Description of variables in the model

We are interested in finding out what is the fraction of packets which reach the outcast, assuming that propagation of each packet happens independently. By Borel’s law of large numbers, we know that it is equal to the probability of the first packet reaching the outcast.

2.1 System evolution

The actual system evolves asynchronously, with each node taking steps independently. However, the time elapsed between two consecutive states of the system does not effect the actual transition the system makes. Hence, we can model this system as a Discrete time Markov process where at each time step, we choose a random node n from the N nodes in the system (barring the outcast), and do one of the following:

1. If n is an infected node in \mathbf{Z}_1 , it will infect f nodes (chosen without replacement from N nodes, all but the infected node itself) and will be removed from I_1 and added to R . A node can be infected only if it either belongs to S_1 (in which case it is shifted to I_1) or to S_2 (in which case, it is shifted to I_2)
2. If n is an infected node in \mathbf{Z}_2 , then all actions happen as in case 1, but the choice of f nodes is limited to $N - 1$ nodes in the network, excluding the outcast and n itself.
3. Otherwise, no action is taken.

Define $\{I_1[k]\}_{k \geq 0}$, $\{S_1[k]\}_{k \geq 0}$, $\{I_2[k]\}_{k \geq 0}$, $\{S_2[k]\}_{k \geq 0}$ as the Discrete time Markov process which describe the evolution of the system after each step.

The system comes to a halt when the number of infected nodes falls to zero. We are interested in finding out the probability of the outcast being infected before the system halts.

3 Analysis

To analyse the model, we define the state of the system as a tuple: $\xi \triangleq (\mathbf{A}, Inf)$, where $\mathbf{A} = (I_1, S_1, I_2, S_2)^T$ such that $0 \leq I_1, S_1, I_2, S_2 \leq N$ and Inf is an indicator variable which is 1 if the outcast is infected. As we rely on a result by Darling *et. al.*, we try to align our notation as much as possible to the notation followed by them in [4].

To use the result, however, we have to make the transition to continuous time as suggested in [3]. Essentially, we assume that each node in the network (except the outcast) has a Poisson clock associated with it, which ticks with rate 1. Hence, the system evolves at the rate of N clock ticks per unit time.

Let $\{I_1(t)\}_{t \geq 0}, \{S_1(t)\}_{t \geq 0}, \{I_2(t)\}_{t \geq 0}, \{S_2(t)\}_{t \geq 0}$ denote the associated stochastic processes which describe the evolution of the system in continuous time. Then, if the k^{th} tick of the clock happens at time τ_k , the following equalities holds $I_1(\tau_k) = I_1[k], S_1(\tau_k) = S_1[k], I_2(\tau_k) = I_2[k]$ and $S_2(\tau_k) = S_2[k]$.

Define $X_t = (I_1(t), S_1(t), I_2(t), S_2(t), Inf(t))^T$ such that $(X_t)_{t \geq 0}$ is a continuous-time Markov chain with a countable and finite state-space denoted by \mathbf{S} . Let $q(\xi, \xi')$ denote the rate of transitions of this Markov chain. Set $\mathbf{I} = \{0, 1\}$ and define the function $y : \mathbf{S} \rightarrow \mathbf{I}$ as $Y_t \triangleq y(X_t) = Inf(t)$.

Define the coordinate function $\mathbf{x} : \mathbf{S} \rightarrow \mathbb{R}^4$, such that $\mathbf{X}_t \triangleq \mathbf{x}(X_t) = \left(\frac{I_1(t)}{N}, \frac{S_1(t)}{N}, \frac{I_2(t)}{N}, \frac{S_2(t)}{N} \right)^T$. The drift vector for states of this system is defined as:

$$\beta(\xi) = \sum_{\xi' \neq \xi} (\mathbf{x}(\xi') - \mathbf{x}(\xi)) q(\xi, \xi') \quad (1)$$

After some combinatorial arguments (see Appendix), the drift vector can be written as:

$$\beta(\xi) = N \cdot \begin{pmatrix} \frac{I_1}{N} \cdot (f \cdot \frac{S_1}{N} - 1) + \frac{I_2}{N} \cdot f \cdot \frac{S_1}{N-1} \\ -\frac{I_1}{N} \cdot f \cdot \frac{S_1}{N} - \frac{I_2}{N} \cdot f \cdot \frac{S_1}{N-1} \\ \frac{I_1}{N} \cdot f \cdot \frac{S_2}{N} + \frac{I_2}{N} \cdot (f \cdot \frac{S_2}{N-1} - 1) \\ -\frac{I_1}{N} \cdot f \cdot \frac{S_2}{N} - \frac{I_2}{N} \cdot f \cdot \frac{S_2}{N-1} \end{pmatrix}$$

Now consider the associated process Y_t . It is also a continuous time Markov process on the state space \mathbf{I} . After the outcast receives the packet, it remains infected with that packet forever. Hence, 1 is an absorbing state for the chain. The rate of transition from 0 (not-infected) to 1 (infected) at time t can be written as: $N \cdot \frac{I_1(t)}{N} \cdot \binom{N-1}{f-1} / \binom{N}{f}$, which can be understood as the thinning of the Poisson process of the entire system (whose rate was N) such that it only considers the events when the node whose clock ticked is an infected node in \mathbf{Z}_1 (with probability $= \frac{I_1(t)}{N}$), and that one of the nodes out of the f chosen by this infected node is the outcast (with probability $= \binom{N-1}{f-1} / \binom{N}{f}$). On simplification, we have the transition rates of the Markov Chain as:

$$\gamma(y, y') = \begin{cases} N \cdot \frac{I_1(t)}{N} \cdot \frac{f}{N} & \text{if } y = 0, y' = 1 \\ 0 & \text{if } y = 1, y' = 0 \end{cases} \quad (2)$$

Let $U = [0, 1] \times [0, 1] \times [0, 1] \times [0, 1] \subset \mathbb{R}^4$, and $x_0 \in U$. Consider the Lipschitz vector field $b : \mathbb{R}^4 \rightarrow \mathbb{R}^4$:

$$b \left([i_1, s_1, i_2, s_2]^T \right) = \begin{pmatrix} i_1(f s_1 - 1) + i_2 f s_1 \\ -i_1 f s_1 - i_2 f s_1 \\ i_1 f s_2 + i_2(f s_2 - 1) \\ -i_1 f s_2 - i_2 f s_2 \end{pmatrix} \quad (3)$$

and define $(x_t)_{t \leq \zeta}$ as the maximal solution of the differential equation $\dot{x}_t = b(x_t)$ starting from x_0 where ζ is chosen such that there is no other solution in U which is defined for a longer time interval. In our case $\zeta = \infty$.

Name the components of (x_t) as $(i_1(t), s_1(t), i_2(t), s_2(t))$, which are real functions of time. Define y_t as a time inhomogeneous continuous time Markov process on the state space $\mathbf{I} = \{0, 1\}$ with the following transition rates:

$$g(x_t, y, y') = \begin{cases} i_1(t) \cdot f & \text{if } y = 0, y' = 1 \\ 0 & \text{if } y = 1, y' = 0 \end{cases} \quad (4)$$

*The coordinate function scales the components of X_t by N and drops the last component $Inf(t)$, for which the mean field limit does not exist

We compare the stochastic process Y_t with y_t and solutions of the following equations:

$$\mathbf{X}_t = \mathbf{X}_0 + M_t + \int_0^t \beta X_s ds, \quad 0 \leq t \leq T_1 \quad (5)$$

$$x_t = x_0 + \int_0^t b(x_s) ds, \quad 0 \leq t \leq \zeta, \quad (6)$$

where $T_1 = \inf \{t \geq 0 : \beta(X_t) = \infty\}$. For a chosen $\varepsilon > 0$, we aim to bound the probability of the event:

$$\tilde{\Omega}_N = \left\{ \sup_{0 \leq t \leq \infty} \|X_t - x_t\| > \varepsilon \text{ or } Y_t \neq y_t \right\} \quad (7)$$

This is the event when the ODE approximation to the system is *not* accurate enough at *some* point in time, i.e., the state of the system differs by more than ε (with respect to the supremum norm) or the state of the outcast are infected in one and not infected in the other.

We can decompose this event into the union of the following two events at a chosen time horizon t_0 :

$$\tilde{\Omega}_N^1(t_0) = \left\{ \sup_{t \leq t_0} \|X_t - x_t\| > \varepsilon \text{ or } Y_t \neq y_t \right\} \quad (8)$$

$$\tilde{\Omega}_N^2(t_0) = \left\{ \sup_{t \geq t_0} \|X_t - x_t\| > \varepsilon \text{ or } Y_t \neq y_t \right\} \quad (9)$$

$$\tilde{\Omega}_N = \tilde{\Omega}_N^1(t_0) \cup \tilde{\Omega}_N^2(t_0) \quad (10)$$

$$\implies \Pr \left\{ \tilde{\Omega}_N \right\} \leq \Pr \left\{ \tilde{\Omega}_N^1(t_0) \right\} + \Pr \left\{ \tilde{\Omega}_N^2(t_0) \right\} \quad (11)$$

Consider the first of these events $\tilde{\Omega}_N^1(t_0)$. We can bound the probability of this event using Theorem 4.4 in [4]. We can show that we can choose $\varepsilon > 0$ and a time horizon t_0 such that, in the time window $[0, t_0]$ the solution of the ODE x_t and the Markov process \mathbf{X}_t differ by at most ε and that Y_t and y_t remain identical with high probability. The theorem is reiterated in the Appendix for clarity.

Using the theorem, we can prove the following result:

Theorem 1. *If ε and t_0 are chosen such that: $\varepsilon < t_0$, and $\frac{t_0}{N} \leq \delta$, then:*

$$\Pr \left(\sup_{t \leq t_0} \|\mathbf{X}_t - x_t\| > \varepsilon \text{ or } Y_t \neq y_t \right) \leq f\varepsilon t_0 + 10e^{-\varepsilon^2 \cdot N/C} \quad (12)$$

where $C = 18(1 + f)ft_0e^{1+2(4f+1)t_0}$.

Proof. We show this result by calculating the values of the various quantities in the expression.

- The Lipschitz constant K for the vector field (3) is $4f + 1$ (Proof: See Appendix).
- $(t_0 \wedge T) \equiv (t_0 \wedge T_0) \equiv t_0$: At each step in the system, the total number of nodes is conserved and the number of nodes never becomes negative. Therefore, all components of \mathbf{X}_t remain in $[0, 1]$ and the process \mathbf{X}_t never leaves U . Hence, $t_0 \wedge T = t_0$. Also, since Y_t is also confined to the state space $I = \{0, 1\}$, $t_0 \wedge T_0 = t_0$.
- Showing $\Omega_0 = \Omega$: It is ensured because we always start both the ODE and the Stochastic process with the same initial conditions: $\mathbf{X}_0 \equiv x_0$.
- Showing $\Omega_1 = \Omega$: We know:

$$\|\beta(X_t) - b(\mathbf{x}(X_t))\| \leq \frac{I_1}{N} \cdot \sup \left\{ \left(\frac{S_1}{N-1} - \frac{S_1}{N} \right), \left(\frac{S_2}{N-1} - \frac{S_2}{N} \right) \right\} \quad (13)$$

$$\leq \frac{1}{N-1} \cdot \frac{I_1}{N} \sup \left\{ \frac{S_1}{N}, \frac{S_2}{N} \right\} \quad (14)$$

$$\leq \frac{1}{N} \quad (15)$$

Hence, we have:

$$\int_0^{t_0} \|\beta(X_t) - b(\mathbf{x}(X_t))\| dt \leq \int_0^{t_0} \frac{1}{N} dt \quad (16)$$

$$= \frac{t_0}{N} \quad (17)$$

Since $\frac{t_0}{N} \leq \delta$, $\Omega_1 = \Omega$.

- Showing $\Omega_2 = \Omega$: As per Footnote 4 in [4], we know that any choice of A which satisfies:

$$A \geq QJ^2 \exp\{\delta J/(At_0)\} \quad (18)$$

where Q is the maximum jump rate, and J is the upper bound for the supremum norm of the jumps, is sufficient for $\Omega_2 = \Omega$. The maximum jump rate for our process is N , and the maximum supremum norm of the jumps is f/N (as the number of susceptible nodes can fall by a maximum of f in one jump). Hence, set A as:

$$A = \frac{(1+f)fe}{N} \quad (19)$$

Then,

$$\int_0^{t_0} \phi(X_t, \theta) dt = \int_0^{t_0} \max_i \phi^i(X_t, \theta) dt \quad (20)$$

$$= \int_0^{t_0} \max_i \sum_{\xi' \neq X_t} \sigma_\theta(x_t^i(\xi') - x_t^i(X_t)) q(X_t, \xi') dt \quad (21)$$

We know that $\sum_{\xi' \neq X_t} q(X_t, \xi') \leq N$. Since in each step a maximum of f nodes can be changed (from the susceptible sets) $\max_i |x_t^i(\xi') - x_t^i(X_t)| \leq \frac{f}{N}$. It is easy to verify that $\sigma_\theta(x) \leq \frac{1}{2}(\theta|x|)^2 e^{\theta|x|}$. Hence, we have:

$$\int_0^{t_0} \phi(X_t, \theta) dt \leq \int_0^{t_0} \sigma_\theta\left(\frac{f}{N}\right) \cdot N dt \quad (22)$$

$$\leq \frac{1}{2} \left(\frac{\theta f}{N}\right)^2 e^{\frac{\theta f}{N}} N t_0 \quad (23)$$

Since we chose $\varepsilon < t_0 \implies \theta \leq f \cdot N$:

$$\int_0^{t_0} \phi(X_t, \theta) dt \leq \frac{1}{2} \theta^2 \frac{f^2}{N} e^{\frac{\theta f}{N}} t_0 \quad (24)$$

$$\leq \frac{1}{2} \theta^2 \frac{f(1+f)e}{N} t_0 \quad (25)$$

$$= \frac{1}{2} \theta^2 A t_0 \quad (26)$$

Hence, $\Omega_2 = \Omega$ if we set A as $\frac{f(1+f)e}{N}$.

- Showing $\Omega_3 = \Omega$: The transition rates of the processes Y_t are given in equation (2) and the rates of the process y_t in equation (4). Then, we have:

$$\sum_{y \neq y'} |\gamma(X_t, y) - g(\mathbf{x}(X_t), y(X_t), y)| = |\gamma(X_t, 0) - g(\mathbf{x}(X_t), y(X_t), y)| \quad (27)$$

$$= |N \cdot \frac{I_1(t)}{N} \cdot \frac{f}{N} - \frac{I_1(t)}{N} \cdot f| \quad (28)$$

$$= 0 \quad (29)$$

Hence, set $G = 0$ and the condition $\int_0^{T_0 \wedge t_0} \sum_{y \neq y(X_t)} |\gamma(X_t, y) - g(\mathbf{x}(X_t), y(X_t), y)| dt \leq G t_0$ is always satisfied. Hence, $\Omega_3 = \Omega$.

- Calculating κ :

$$\kappa = \sup_{t \leq t_0} \sup_{\|x - x_t\| \leq \varepsilon, y \in I} \sum_{y' \neq y} |g(x, y, y') - g(x_t, y, y')| \quad (30)$$

$$= \sup_{t \leq t_0} \sup_{\|x - x_t\| \leq \varepsilon, y \in I} |g(x, 0, 1) - g(x_t, 0, 1)| \quad (31)$$

$$= \sup_{t \leq t_0} \sup_{\|x - x_t\| \leq \varepsilon, y \in I} |i_1 \cdot f - i_1(t) \cdot f| \quad (32)$$

$$= f\varepsilon \quad (33)$$

Using these values for the said theorem stated in the Appendix, we can show the result holds. \square

Theorem 1 ensures that the probability of the event $\tilde{\Omega}_N^1(t_0)$ is upper bounded and that the bound can be made arbitrarily small as N increases. To bound the probability of $\tilde{\Omega}_N^2(t_0)$, we have the following conjecture.

Conjecture 1. For any $\varepsilon \geq 0$, then $\exists N_0, t_0$ such that $\frac{t_0}{N_0} \leq \delta$ and $\varepsilon < t_0$ and:

$$\Pr \left(\sup_{t \geq t_0} \|\mathbf{X}_t - x_t\| > \varepsilon \text{ or } Y_t \neq y_t \right) \leq \varepsilon + \Pr \left\{ \tilde{\Omega}_N^1(t_0) \right\} \quad (34)$$

We believe that this conjecture is true because it can be shown that $i_1(t)$ and $i_2(t)$ for the ODE approximation fall exponentially with t for any initial state. Hence, we know $\exists t_0$ such that:

1. $\sup_{t \geq t_0} \|x_t - x_{t_0}\| \leq \varepsilon/4$
2. $\Pr \{ \exists t \geq t_0, y_t \neq y_{t_0} \} \leq \int_{t_0}^{\infty} f \cdot i_1(t) dt \leq \varepsilon/4$

A similar result is expected to hold for the processes Y_t and X_t , after fixing a minimum value of N_0 . Using this conjecture, we would be able to show:

Conjecture 2. $\lim_{N \rightarrow \infty} \lim_{t \rightarrow \infty} \Pr \{ Y_t = 1 \mid Y_0 = 0 \} = \lim_{t \rightarrow \infty} \Pr \{ y_t \mid y_0 = 0 \}$

And bound the probability of the event $\tilde{\Omega}_N$ as an explicit function of N . This conjecture is empirically found to be true as suggested in the following section. However, such a verification can only be accepted with partial faith, and work is ongoing for proving the result rigorously.

4 Simulation results

To verify the model, we run experiments for varying number of nodes in the system. Two of these cases are presented here, for $N = 100$ nodes in Figure 2(a) and for $N = 1000$ nodes in Figure 2(b). These show a close fit of the simulation data by the model predictions.

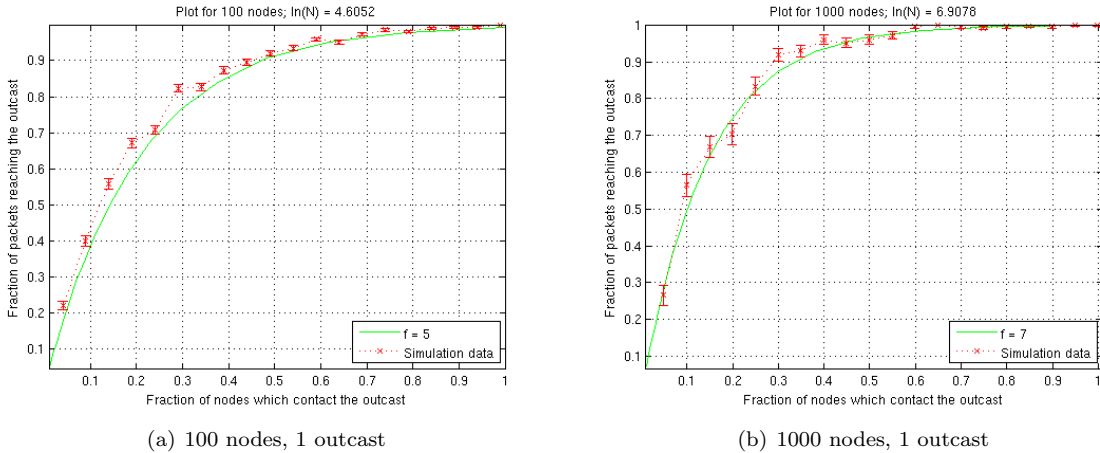


Figure 2: Verification of model predictions

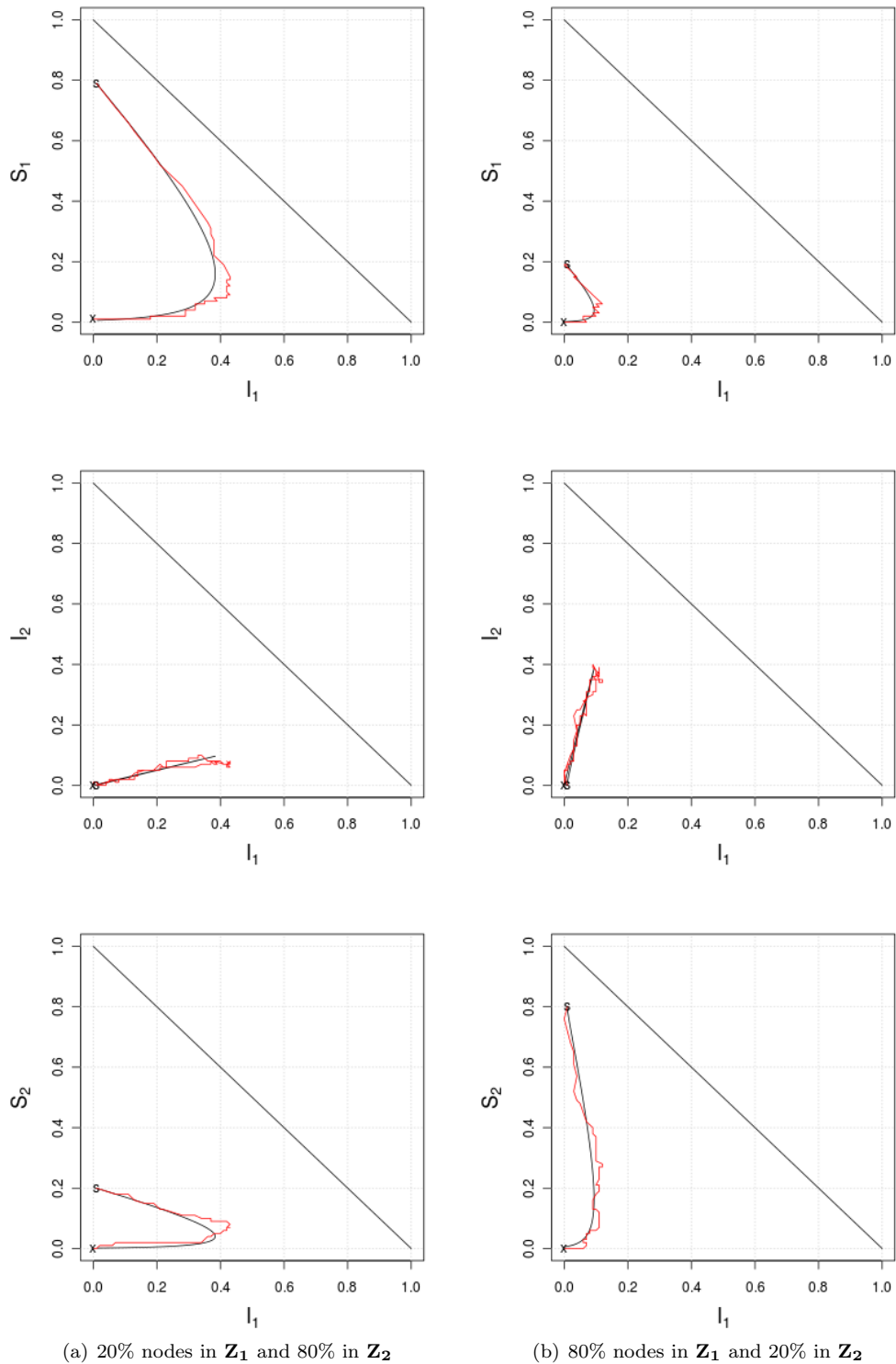


Figure 3: Phase plots of the evolution of the system with two different fractions of nodes blacklisting the outcast as predicted by the ODE (continuous line) and as obtained by an actual simulation (jagged line). The simulation was done for 100 nodes with $f = 5$. The starting state is denoted by \mathbf{S} and point the system halts is denoted by \mathbf{X} . The system starts with one infected (source) node in \mathbf{Z}_1 and all other nodes susceptible.

The phase plot is shown in Figure 3 with the details of the stochastic simulation as well as the path of the deterministic system. It can be seen that the ODE approximation is fairly representative of the system with only 100 nodes. Figure 4 shows the probability of the outcast getting infected as predicted by the deterministic system as well as obtained by 1000 simulation runs.

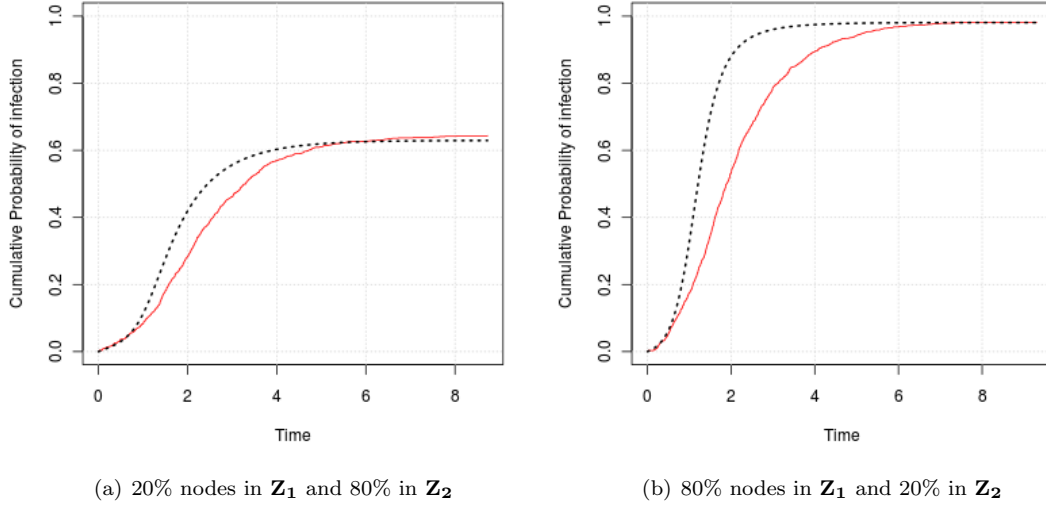


Figure 4: The probability that the outcast node would get infected before time t as predicted by the deterministic system (dotted line) and as obtained by 1000 repeated runs of a simulation (continuous line) with 100 nodes. The difference between 1 and the final value of cumulative probability of infection is the probability that the outcast does not get infected.

After verifying the model for small number of nodes via simulations, we predict using the model what would be the behaviour of the system as the number of nodes increases. The results of simulations with different fractions of nodes blacklisting the outcast with different number of nodes is shown in Figure 5. To ensure that an outcast receives than 90% of the data, the fraction of nodes

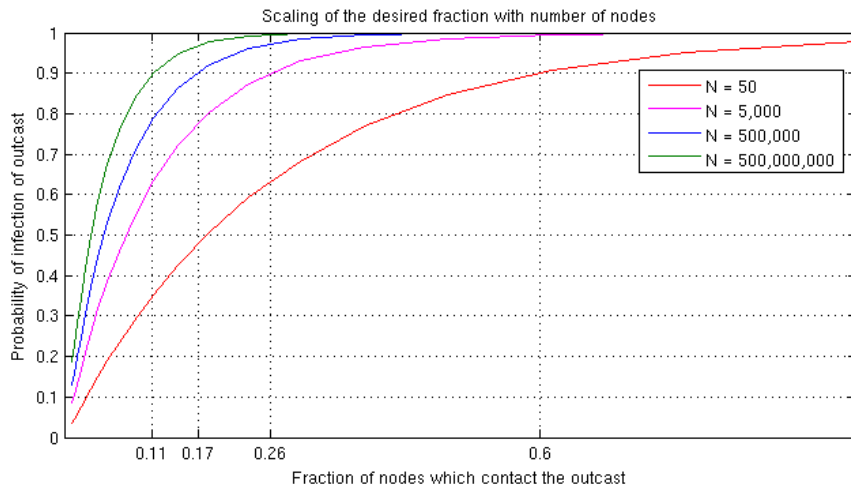


Figure 5: Fraction of nodes which need to contact the outcast to deliver 90% packets to it. Different lines denote different sizes of the network and the x -axis marks the points denoting the fraction of nodes which should blacklist the outcast corresponding to 90% reception.

which need to blacklist it increases with increasing N . This is owing to the increase of the fanout f to maintain the probabilistic guarantees of the system; in particular, the guarantee of atomic broadcasts. The necessary fraction are tabulated in Table 2.

Number of nodes in network	Fanout size	Required Fraction
50	4	40%
5,000	9	74%
500,000	14	82%
500,000,000	21	89%

Table 2: Fraction of nodes which needs to stop communicating with the outcast to get 90% of the data, as shown in Figure 5 the number of nodes in the system were chosen so as to make the increase in fanout size balanced.

5 Discussion

In this paper, we show that it is possible to model the interaction between a system evolving as a discrete state space Markov process and an individual using mean field theory. We take help existing results in arriving at probabilistic guarantees on accuracy of the model up to a certain time horizon and suggest conjectures to prove the probabilistic bounds on the asymptotic behaviour of the model. We also empirically verify the model using simulations and show that the model is reasonably accurate even for a system with only 100 nodes.

Using the model, we show that distributed blacklisting of nodes cannot be scaled up while preserving the quality of service to other nodes. Hence, it is difficult to *outcast* a node in a decentralized manner after it has been inducted in the network. This result helps us develop a better understanding of distributed blacklisting and reiterates the robustness of gossip based protocols.

6 Acknowledgement

The authors would like to thank Rohit Nagpal for his helpful comments and aid in analysis of the system.

Appendix A: Derivation of the drift vector

The drift vector for states of this system is defined as:

$$\beta(\xi) = \sum_{\xi' \neq \xi} (\mathbf{x}(\xi') - \mathbf{x}(\xi)) q(\xi, \xi')$$

To arrive at an expression of the vector, we deal with each component individually. We know that the total rate of the system is N and that if $\xi \neq \xi'$, then $q(\xi, \xi') = N \cdot p_{\xi \rightarrow \xi'}^{(1)}$, where the events are ticks of the Poisson clocks attached to each node. The first component of the drift vector (corresponding to I_1), is:

$$\sum_{\xi' \neq \xi} \left(\frac{I'_1}{N} - \frac{I_1}{N} \right) \cdot N \cdot p_{\xi \rightarrow \xi'}^{(1)} = N \cdot \mathbf{E} \left[\frac{I'_1}{N} - \frac{I_1}{N} \right] \quad (35)$$

To calculate $\mathbf{E} [I'_1 - I_1]$, we condition on the node whose clock ticked. We know that this is equivalent to choosing a node at random from the set of N nodes. Call the node whose clock ticked n . Let E_1 be the event that n is an infected node in \mathbf{Z}_1 and E_2 be the event that n is an infected node in \mathbf{Z}_2 . We know $\Pr\{E_1\} = \frac{I_1}{N}$ and $\Pr\{E_2\} = \frac{I_2}{N}$.

$$\mathbf{E} \left[\frac{I'_1}{N} - \frac{I_1}{N} \right] = \mathbf{E} \left[\frac{I'_1}{N} - \frac{I_1}{N} \mid E_1 \right] \cdot \frac{I_1}{N} + \mathbf{E} \left[\frac{I'_1}{N} - \frac{I_1}{N} \mid E_2 \right] \cdot \frac{I_2}{N} + 0 \cdot \left(1 - \frac{I_1 + I_2}{N} \right) \quad (36)$$

Consider $\mathbf{E} \left[\frac{I'_1}{N} - \frac{I_1}{N} \mid E_1 \right]$. This can be calculated by noting that probability that $I'_1 = I_1 + j$ given E_1 happened means that out of f nodes which were infected, $j + 1$ nodes were chosen from the susceptible nodes in \mathbf{Z}_1 and rest were chosen from the other $N - S_1$ nodes (excluding n itself). It should be noted here that if the node n is an infected node in \mathbf{Z}_1 , it can infect f nodes from N other nodes (including the outcast) while if it is in \mathbf{Z}_2 , it can infect only $N - 1$ nodes.

$$\text{Now, } \mathbf{E} \left[\frac{I'_1}{N} - \frac{I_1}{N} \mid E_1 \right] = \sum_{k=-1}^{f-1} k \frac{\binom{S_1}{k+1} \binom{N-S_1}{f-(k+1)}}{\binom{N}{f}} \quad (37)$$

$$= \sum_{k=0}^f k \frac{\binom{S_1}{k} \binom{N-S_1}{f-k}}{\binom{N}{f}} - \sum_{k=0}^f \frac{\binom{S_1}{k} \binom{N-S_1}{f-k}}{\binom{N}{f}} \quad (38)$$

$$\text{And, } \mathbf{E} \left[\frac{I'_1}{N} - \frac{I_1}{N} \mid E_2 \right] = \sum_{k=0}^f k \cdot \frac{\binom{S_1}{k} \binom{N-1-S_1}{f-k}}{\binom{N-1}{f}} \quad (39)$$

Now consider the binomial expansion of the polynomial $x \cdot \left[\frac{d}{dx} (1+x)^{S_1} \right] \cdot (1+x)^{N-S_1} \equiv S_1 \cdot x \cdot (1+x)^{N-1}$. By comparing the coefficient of x^f on both sides of the equivalent expression, we find:

$$\sum_{k=0}^f k \binom{S_1}{k} \binom{N-S_1}{f-k} \equiv S_1 \binom{N-1}{f-1}$$

It is easy to verify then that:

$$\sum_{k=0}^f k \frac{\binom{S_1}{k} \binom{N-S_1}{f-k}}{\binom{N}{f}} = \frac{S_1 \cdot f}{N}$$

Using similar arguments, for the rest of the expressions and components of the drift vector, we arrive at the final expression:

$$\beta(\xi) = N \cdot \begin{pmatrix} \frac{I_1}{N} \cdot \left(f \cdot \frac{S_1}{N} - 1 \right) + \frac{I_2}{N} \cdot f \cdot \frac{S_1}{N-1} \\ - \frac{I_1}{N} \cdot f \cdot \frac{S_1}{N} - \frac{I_2}{N} \cdot f \cdot \frac{S_1}{N-1} \\ \frac{I_1}{N} \cdot f \cdot \frac{S_2}{N} + \frac{I_2}{N} \cdot \left(f \cdot \frac{S_2}{N-1} - 1 \right) \\ - \frac{I_1}{N} \cdot f \cdot \frac{S_2}{N} - \frac{I_2}{N} \cdot f \cdot \frac{S_2}{N-1} \end{pmatrix}$$

It should be noted here that this drift vector is obtained by using exact calculations which required combinatorial arguments. However, if we relax the model by allowing for selection with replacement of nodes, the choice of nodes could be made independently and it would be possible to use linearity of expectation to readily arrive at the expression for the drift vector. This relaxation would be especially useful in the case when f is not fixed but is drawn randomly from some distribution, allowing the use of Wald's equation.

Appendix B: Differential equation approximations for Markov chains

Theorem 2. Choose $\varepsilon > 0$ and t_0 such that:

$$\forall \xi \in S \text{ and } t \leq t_0, \|\mathbf{x}(\xi) - x_t\| \leq \varepsilon \implies \mathbf{x}(\xi) \in U. \quad (40)$$

Let K be the Lipschitz norm of the vector field b with respect to the supremum norm $\|\cdot\|$ and let d be the dimension of the vector X_t . Let $\delta = \varepsilon e^{-Kt_0}/3$. Fix $A > 0$ and set $\theta = \delta/(At_0)$. Define:

$$\sigma_\theta = e^{\theta|x|} - 1 - \theta|x|, x \in \mathfrak{R} \quad (41)$$

and set :

$$\phi^i(\xi, \theta) = \sum_{\xi' \neq \xi} \sigma_\theta \left(x^i(\xi') - x^i(\xi) \right) q(\xi, \xi') \quad (42)$$

$$\phi(\xi, \theta) = \max_i \phi^i(\xi, \theta) \quad (43)$$

Let $T = \inf \{t \geq 0 : \mathbf{X}_t \notin U\}$ and $T_0 = \inf \{t \geq 0 : \mathbf{X}_t \notin U \text{ or } Y_t \notin I\}$.

Set

$$\kappa = \sup_{t \leq t_0} \sup_{\|x - x_t\| \leq \varepsilon, y \in I} \sum_{y' \neq y} |g(x, y, y') - g(x_t, y, y')| \quad (44)$$

Fix $G > 0$ and define the following events:

$$\Omega_0 = \{\|\mathbf{X}_0 - x_0\| \leq \delta\} \quad (45)$$

$$\Omega_1 = \left\{ \int_0^{T \wedge t_0} \|\beta(X_t) - b(\mathbf{x}(X_t))\| dt \leq \delta \right\} \quad (46)$$

$$\Omega_2 = \left\{ \int_0^{T \wedge t_0} \phi(X_t, \theta) dt \leq \frac{1}{2} \theta^2 A t_0 \right\} \quad (47)$$

$$\Omega_3 = \left\{ \int_0^{T_0 \wedge t_0} \sum_{y \neq y(X_t)} |\gamma(X_t, y) - g(\mathbf{x}(X_t), y(X_t), y)| dt \leq G t_0 \right\} \quad (48)$$

Then, we have:

$$\Pr \left(\sup_{t \leq t_0} \|\mathbf{X}_t - x_t\| > \varepsilon \text{ or } Y_t \neq y_t \right) \leq (G + \kappa) t_0 + 2de^{-\delta^2/2At_0} + \Pr(\Omega_0^c \cup \Omega_1^c \cup \Omega_2^c \cup \Omega_3^c) \quad (49)$$

Proof. See [4]. □

Appendix C: Finding the Lipschitz constant of the vector field (3)

To show that the vector field given in equation (3) has the Lipschitz constant $K = (4f + 1)$, we make use of the following Lemma.

Lemma 1. Let $z_1 = [x_1, y_1]^T$ and $z_2 = [x_2, y_2]^T$. Consider the function $f : [0, 1] \times [0, 1] \rightarrow \mathfrak{R}$:

$$f([x, y]) = c \cdot x \cdot y, \quad 0 \leq x, y \leq 1 \quad (50)$$

Function f has Lipschitz constant $2 \cdot |c|$ with respect to supremum norm $\|\cdot\|$, or, equivalently:

$$|f(z_2) - f(z_1)| \leq 2 \cdot |c| \cdot \|z_2 - z_1\| \quad (51)$$

Proof. Assume the two points on which the $|f(z_2) - f(z_1)|/\|z_2 - z_1\|$ is maximum lie along a line with slope m . By symmetry, we know that there will exist such a line such that $|m| \leq 1$.

Now we consider the three cases where this line cuts two opposite sides or two adjacent sides of the square $[0, 1] \times [0, 1]$ and find the maximum value of $|f(z_2) - f(z_1)|/\|z_2 - z_1\|$ in each case. Note that since we have assumed $m \leq 1$, so that $\|z_2 - z_1\| \equiv |x_2 - x_1|$.

1. **Opposite sides are intercepted:** Say that the sides $x = 0$ and $x = 1$ are intersected at heights a and b , respectively. Then the equation of the line is give by: $y = \frac{(b-a)}{1} \times x + a$ and substituting the value of y_1 and y_2 , we have:

$$\frac{|f(z_1) - f(z_2)|}{|x_2 - x_1|} = |c| \cdot |(b-a) \times (x_2 + x_1) + a| \quad (52)$$

$$\leq |c| \cdot |2(b-a) + a| \quad (53)$$

$$\text{(since, } 0 \leq a, b \leq 1) \leq 2|c| \quad (54)$$

2. **Adjacent $x = 1$ and $y = 0$ sides are intersected:** Say that the intercepts on $x = 1$ and $y = 0$ are at height a and distance b from origin, respectively. Then, the equation of the line is given by $y = \frac{a}{(1-b)} \times x + \frac{ab}{(1-b)}$, where $\frac{a}{(1-b)} \leq 1$ and substituting the values of y_1 and y_2 , we have:

$$\frac{|f(z_1) - f(z_2)|}{|x_2 - x_1|} = |c| \cdot \left| \frac{a}{(1-b)} \times (x_2 + x_1) - \frac{ab}{(1-b)} \right| \quad (55)$$

$$\leq |c| \cdot \left| \frac{a}{(1-b)} \right| \cdot |(x_2 + x_1) - b| \quad (56)$$

$$\leq |c| \cdot |(x_2 + x_1) - b| \quad (57)$$

$$\leq 2 \cdot |c| \quad (58)$$

3. **Adjacent $x = 0$ and $y = 0$ sides are intersected:** Say that the intercepts on $x = 0$ and $y = 0$ are at distance a and b from origin, respectively. Then the equation of the line would be $y = -\frac{b}{a}x + b$. Also, we will have $0 \leq x_1 + x_2 \leq 2a$ since $0 \leq x_1, x_2 \leq a$. Again, substituting the values of y_1 and y_2 , we have:

$$\frac{|f(z_1) - f(z_2)|}{|x_2 - x_1|} = |c| \cdot \left| b - \frac{b}{a}(x_1 + x_2) \right| \quad (59)$$

$$= |c| \cdot |b| \cdot \left| 1 - \frac{(x_1 + x_2)}{a} \right| \quad (60)$$

$$\leq |c| \cdot |b| \quad (61)$$

$$\leq |c| \quad (62)$$

Hence, the Lipschitz constant in for this function is $\sup \{2|c|, |c|\} = 2|c|$. \square

Lemma 2. If function $f(\cdot)$ is uniformly continuous with Lipschitz constant k_1 and function $g(\cdot)$ is uniformly continuous with Lipschitz constant k_2 , then:

1. $f(\cdot) \pm g(\cdot)$ is uniformly continuous with Lipschitz constant $K \leq k_1 + k_2$.

Proof.

$$|(f(x_2) \pm g(x_2)) - (f(x_1) \pm g(x_1))| \leq |f(x_2) - f(x_1)| + |g(x_2) - g(x_1)| \quad (63)$$

$$\leq (k_1 + k_2) \cdot \|x_2 - x_1\| \quad (64)$$

□

2. If $F(\cdot) \triangleq [f(\cdot), g(\cdot)]^T$, then under the supremum norm, that is, $\|F(x_2) - F(x_1)\| \triangleq \sup \{f(x_2) - f(x_1), g(x_2) - g(x_1)\}$, $F(\cdot)$ is Uniformly continuous with Lipschitz constant $K = \sup \{k_1, k_2\}$.

Proof.

$$\sup \{f(x_2) - f(x_1), g(x_2) - g(x_1)\} \leq \sup \{k_1 \cdot \|x_2 - x_1\|, k_2 \cdot \|x_2 - x_1\|\} \quad (65)$$

$$\leq \sup \{k_1, k_2\} \cdot \|x_2 - x_1\| \quad (66)$$

□

For the vector field defined by equation (3), the Lipschitz constant can be calculated for each component:

$$b \left([i_1, s_1, i_2, s_2]^T \right) = \begin{pmatrix} i_1(f s_1 - 1) + i_2 f s_1 \\ -i_1 f s_1 - i_2 f s_1 \\ i_1 f s_2 + i_2(f s_2 - 1) \\ -i_1 f s_2 - i_2 f s_2 \end{pmatrix} \xleftarrow{\text{L-const.}} \begin{pmatrix} 2f + 1 + 2f \\ 2f + 2f \\ 2f + 2f + 1 \\ 2f + 2f \end{pmatrix} \quad (67)$$

Now taking the maximum of these constants, we can show that $K = 4f + 1$.

References

- [1] Rena Bakhshi, J. Endrullis, Stefan Endrullis, Wan Fokkink, and Boudewijn Haverkort. Automating the mean-field method for large dynamic gossip networks. In *Proceedings of the 7th International Conference on the Quantitative Evaluation of Systems, QEST 2010*, pages 241–250. IEEE Computer Society, September 2010.
- [2] Ken Birman. The promise, and limitations, of gossip protocols. *Operating Systems Review*, pages 8–13, 2007.
- [3] R. W. R. Darling. Fluid limits of pure jump markov processes: a practical guide. *Construction*, (July):16, 2002.
- [4] R. W. R. Darling and J. R. Norris. Differential equation approximations for Markov chains. *Probability Surveys*, 5:37–39, 2008.
- [5] P. T. Eugster, Rachid Guerraoui, Anne-Marie Kermarrec, and L. Massoulié. From epidemics to distributed computing. *IEEE Computer*, 37:60–67, 2004.
- [6] Davide Frey, Rachid Guerraoui, Anne-Marie Kermarrec, Boris Koldehofe, Martin Mogensen, Maxime Monod, and Vivien Quéma. Heterogeneous Gossip. In *Proceedings of the 10th ACM/IFIP/USENIX International Middleware Conference (Middleware)*, 2009.
- [7] Rachid Guerraoui, Kévin Huguenin, Anne-Marie Kermarrec, Maxime Monod, and Swagatika Prusty. Lifting: Lightweight freerider-tracking protocol in gossip. 2010.
- [8] Márk Jelasity, Spyros Voulgaris, Rachid Guerraoui, Anne-Marie Kermarrec, and Maarten van Steen. Gossip-based peer sampling. *ACM Transactions on Computer Systems*, 25(3):8, 2007.
- [9] Hamidou Tembine, Jean-Yves Le Boudec, Rachid El-Azouzi, and Eitan Altman. Mean Field Asymptotic of Markov Decision Evolutionary Games and Teams. In *Gamenets*, 2009. Invited Paper.