

# Finding Influential Neighbors to Maximize Information Diffusion in Twitter

Hyounghshick Kim  
Department of Computer  
Science and Engineering  
Sungkyunkwan University  
Suwon, Korea  
hyoung@skku.edu

Konstantin Beznosov  
Department of Electrical and  
Computer Engineering  
University of British Columbia  
Vancouver, Canada  
beznosov@ece.ubc.ca

Eiko Yoneki  
Computer Laboratory  
University of Cambridge  
Cambridge, United Kingdom  
eiko.yoneki@cl.cam.ac.uk

## ABSTRACT

The problem of spreading information is a topic of considerable recent interest, but the traditional influence maximization problem is inadequate for a typical viral marketer who cannot access the entire network topology. To fix this flawed assumption that the marketer can control any arbitrary  $k$  nodes in a network, we have developed a decentralized version of the influential maximization problem by influencing  $k$  neighbors rather than arbitrary users in the entire network. We present several reasonable neighbor selection schemes and evaluate their performance with a real dataset collected from Twitter. Unlike previous studies using network topology alone or synthetic parameters, we use real propagation rate for each node calculated from the Twitter messages during the 2010 UK election campaign. Our experimental results show that information can be efficiently propagated in online social networks using neighbors with a high propagation rate rather than those with a high number of neighbors.

## Categories and Subject Descriptors

H.3.5 [Information Storage and Retrieval]: Online Information Services; J.4 [Computer Applications]: Social and Behavioral Sciences

## General Terms

Human Factors, Measurement

## Keywords

Information diffusion; Information dissemination; Online social networks; Viral marketing

## 1. INTRODUCTION

In the field of social network analysis, a fundamental problem is to develop an epidemiological model for finding an

efficient way to spread information through the model. It seems natural that many people are often influenced by their friends' opinions or recommendations. This is called the "word of mouth" effect and has for long been recognized as a powerful force affecting product recommendation [2].

Recent advances in the network theory have provided us with the mathematical and computational tools to understand them better. For example, in the *Independent Cascade* (IC) model proposed by Goldenberg et al. [8], (1) some non-empty set of nodes are initially *activated* (or influenced); (2) at each successive step, the influence is propagated by activated nodes, independently activating their inactive neighbors based on the *propagation probabilities* of the adjacent edges. Here, activated nodes mean the nodes that have adopted the information or have been infected. This models how a piece of information will likely be spread through a network over time. It enables us to investigate what sort of information diffusion scheme might be the most effective one under certain conditions.

This model is also highly relevant to security. For example, cyberstalkers might be interested in spreading rumors, gossips, news or pictures through social networks to damage their victims' (e.g., celebrity, political party, company or country) reputation. The same model works in social media campaign where spammers and propagandists want to share their advertisements on online social networks; fake accounts with automated bots are often used to amplify advertising campaigns in social media [17, 3, 19].

Thus far, however, the models and analytic tools used to analyze epidemics have been somewhat limited. Most previous studies [7, 11] aimed to analyze the process of information diffusion by choosing a set of arbitrary  $k$  nodes in a network as the initially activated nodes from a bird's eye perspective based on the full control of the entire network, which may indeed be unacceptable in many real life applications since there is no such central entity (except the online social network service provider itself).

From the point of view of an individual user (e.g., viral marketer) who wants to efficiently spread a piece of information (or a rumor) through a network, a more reasonable epidemiological model would not assume the knowledge of the entire network topology. Kim and Yoneki [14] recently introduced the problem called *Influential Neighbor Selection* (INS) where a spreader  $s$  spreads a piece of information through carefully chosen  $k$  neighbors of hers instead of a set of any arbitrary  $k$  nodes in a network. Under this model, each user can only communicate with the user's immediate

Copyright is held by the International World Wide Web Conference Committee (IW3C2). IW3C2 reserves the right to provide a hyperlink to the author's site if the Material is used in electronic media.

WWW'14 Companion, April 7–11, 2014, Seoul, Korea.

ACM 978-1-4503-2745-9/14/04.

<http://dx.doi.org/10.1145/2567948.2579358>.

neighbors and has no knowledge of the global network topology except for her own connections. However, their work has two limitations: (1) it was simply assumed to use a constant propagation rate, despite variations in user propagation rates in practice. For example, in real-world online social network services such as Twitter or Facebook, each user has a distinct propagation rate for her neighbors on spreading information according to the user’s reputation and/or role, such as opinion formers, leaders or followers [4]; (2) their experimental results were limited to undirected graphs with parameter values chosen in a somewhat ad-hoc manner.

More recently, Kim [12] extended this model by introducing several parameters (*user propagation weight*, *content interestingness*, and *decay factor*) to provide a more general and practical information diffusion model. This gives much finer granularity than the previous model [14]. However, their experiments still depended on synthetic parameters that might significantly affect the information diffusion process.

With a real dataset (Twitter users and messages related to the 2010 UK election campaign), we revisited the INS problem and evaluated the performance of four spreading schemes from the simple random neighbor selection to a sophisticated neighbor selection scheme using both the “number of friends” and “user propagation rate” each neighbor has. To measure the performance of these schemes, we used the conventional *Independent Cascade* (IC) model [8], which is widely used for the analysis of information diffusion [8, 11, 9].

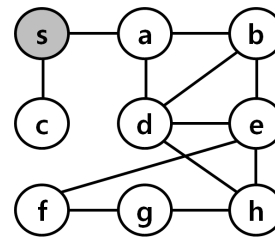
We performed simulation with various parameters. Our experimental results suggest that the scheme to select neighbors who wrote popular posts produced the best overall results, even without consideration of the “number of friends”. That is, we can efficiently spread information without knowing the “number of friends” each user has. Moreover, we found that the information diffusion speed of some schemes (e.g., random neighbor selection) in the previous study [12] was quite exaggerated and thus contributed to the reduction of the performance gap between information diffusion schemes. For example, we observed that the **Random** selection scheme is not practically effective even with a high number  $k$  initially activated nodes; this is quite different from previous studies [14, 12], which showed that the **Random** selection scheme achieved reasonable performance when  $k \geq 3$ .

The rest of this paper is organized as follows. In Section 2 we formally define the INS problem and notations. Then, we present the four reasonable neighbor selection schemes in Section 3. In Section 4, we evaluate their performance through simulation with a real dataset collected from Twitter, and recommend the best neighbor selection scheme with various conditions. Related work is discussed in Section 5. We conclude in Section 6.

## 2. INFLUENTIAL NEIGHBOR SELECTION PROBLEM

We begin with the definition of the *Independent Cascade* (IC) model [8], and then introduce the *Influential Neighbor Selection* (INS) problem, which will be used in the rest of the paper.

We model an *influence network* as a directed graph  $G = (V, E)$  consisting of a set of nodes  $V$  and a set of ordered



**Figure 1: An example of the INS problem. With the spreader node  $s$ , when  $k = 1$ , we should choose  $a$  as an initially activated node to maximize  $|S_t|$ ; however, in the traditional IM problem, the optimal choice might be either  $d$  or  $e$  rather than  $a$ .**

pairs of nodes  $E$  called the edge set, representing the communication channels between node pairs. A directed edge  $(u, v)$  from node  $u$  to node  $v$  of  $G$  is associated with a *propagation probability*  $\lambda_{u,v}$ , which is the probability that  $v$  is activated by  $u$  through the edge in the next time step if  $u$  is activated. Here,  $v$  is said to be a *neighbor* (or successor) of node  $u$ . For node  $u \in V$ , we use  $N(u)$  to denote the set of  $u$ ’s neighbors. The *out-degree* of node  $u$  is denoted as  $d(u) = |N(u)|$ , which could be used simply in estimating the node  $u$ ’s influence on information propagation.

In the IC model [8], we assume that the time during which a network is observed is finite; without loss of generality, the time period is divided into fixed discrete steps  $\{1, \dots, t\}$ . Let  $S_i \subseteq V$  be the set of nodes that are activated at the time step  $i$ . We consider the dynamic process of information diffusion starting from the set of nodes  $S_0 \subseteq V$  that are initially activated until the time step  $t$  as follows: At each time step  $i$  where  $1 \leq i \leq t$ , every node  $u \in S_{i-1}$  activates its inactivated neighbors  $v \in V \setminus S_{i-1}$  with a propagation probability  $\lambda_{u,v}$ . The process ends after the time step  $t$  with  $S_t$ . A conventional *Influential Maximization* (IM) problem is to find a set  $S_0$  consisting of  $k$  nodes to maximize  $|S_t|$ .

The *Influential Neighbor Selection* (INS) problem [14] is a variant of the IM problem: Given a spreader  $s \in V$  and a budget constraint  $k$ , we aim to maximize the number of activated nodes in a network after the time step  $t$  by selecting  $s$ ’s  $\min(k, d(s))$  neighbors only (rather than any subset of  $k$  nodes), as the set of nodes  $S_0 \subseteq V$  that are initially activated. Compared to the conventional IM problem, the INS problem has three additional requirements: (1) each node only communicates with its immediate neighbors; (2) each node has no knowledge about the entire network topology except for its own connections; and (3) each message size is bounded to  $O(\log |V|)$  bits (more intuitively, each message can only contain the node identity and some constant values of the node properties). Figure 1 shows an example of the INS problem. Given this graph with the spreader node  $s$ , when  $k = 1$ , we should choose  $a$  as an initially activated node to maximize the number of activated nodes in future; however, in the traditional IM problem, the optimal choice might be either  $d$  or  $e$  rather than  $a$ .

However, the initial IM problem in [14] – every edge has the same propagation probability – is too simple to correctly reflect the characteristics of the information diffusion process in real-world situations. Clearly, in the most popular online social network services such as Twitter or Facebook, each user has a different propagation rate for her neighbors on

spreading information in a network according to the user’s reputation or role such as opinion formers, leaders or followers [4]. Kim [12] extended this epidemiological model by introducing the three important parameters (user propagation weight  $\omega$ , content interestingness  $\phi$ , and decay factor  $\gamma$ ) to establish a more general and practical information propagation model. The details are as follows:

The user propagation weight  $\omega$  represents each user’s average propagation rate to her neighbors. Given a user  $u$ ,  $\omega(u)$  is defined as  $\tau(u)/(\rho(u)/d(u))$  where  $\tau(u)$  and  $\rho(u)$  are the number of  $u$ ’s posts shared by  $u$ ’s neighbors and the number of  $u$ ’s all posts, respectively. For example, if a user  $u$  with 1,000 neighbors wrote 10 posts and gets 100 shares,  $\omega(u)$  is  $100/(10 \cdot 1000) = 0.01$ .

The content interestingness  $\phi(r)$  of information  $r$  represents a measure to determine how much users want to share the information  $r$  with their neighbors. Naturally, higher content interestingness  $\phi$  of a piece of information may facilitate higher propagation for the information through a network. Previous studies [22, 10] showed that propagation probability  $\lambda$  can be greatly changed with the content of information (i.e., content interestingness  $\phi$ ).

The decay factor  $\gamma$  at hop  $N$  represents the ratio between the propagation probability at hop  $N$  and the propagation probability at hop  $N - 1$ . In practice, the propagation probability might decay exponentially as the cascades spreads away from the information source. Here, one possible explanation would be that the freshness of the information would drop as the time goes on.

With these parameters, given an edge  $(u, v) \in E$ , a spreader  $s \in V$  and a piece of information  $r$ ,  $\lambda(u, v, s, r)$  is finally defined as follows [12]:

$$\lambda(u, v, s, r) = \min\{\omega(u) \cdot \phi(r) \cdot \gamma^{\delta(u, s, r)-1}, 1\} \quad (1)$$

where  $\delta(u, s, r)$  is the number of times the information  $r$  is to be relayed from  $s$  to  $u$ .

For example, when  $\phi(r) = 0.0136$ ,  $\delta(u, s, r) = 3$  and  $\gamma = 0.2$ , a user  $u$  with  $\omega(u) = 1$  would activate his (or her) neighbor  $v$  with the probability of about  $0.0005 (\approx 1 \cdot 0.0136 \cdot (0.2)^2)$ .

In this paper, we also use these parameters and the propagation probability equation to provide more realistic simulation results.

### 3. NEIGHBOR SELECTION SCHEMES

For the INS problem described in Section 2, we basically use a greedy strategy to select the influential neighbors.

Assume that a spreader  $s \in V$  wants to spread a piece of information  $r$  through the network  $G = (V, E)$  by sharing  $r$  with its  $\min(k, d(s))$  neighbors at the initial step. Node  $s$  first tries to assess the influence of information diffusion for each neighbor  $v \in N(s)$ , respectively, by collecting the information about  $v$ . We note that neighbors’ influence should be estimated based on  $s$ ’s local information only, rather than the whole network. Since online social networks, such as Facebook, typically provide APIs to obtain the neighborhood information about user,  $s$  can automatically collect the information about her own neighbors. After estimating the neighbors’ influences,  $s$  selects the top  $\min(k, d(s))$  nodes with the highest influence values from  $N(s)$ ; that is, for the IC model in Section 2, these nodes are selected as the set of initially activated nodes  $S_0 \subseteq V$ .

For the purpose of influence estimation, we test the following four selection schemes based on the “number of friends” and “user propagation weight” each user has:

- **Random selection:** Pick  $\min(k, d(s))$  nodes randomly from  $N(s)$ . This scheme is very simple and easy to implement – the spreader  $s$  does not need any knowledge of the network topology.
- **Degree selection:** Pick the  $\min(k, d(s))$  highest-degree nodes from  $N(s)$ . This scheme requires the degree knowledge of neighbors.
- **Propagation-weight selection:** Pick the  $\min(k, d(s))$  highest user propagation weight nodes from  $N(s)$ . This scheme requires the user propagation weight knowledge of the nodes. To calculate  $\omega(v)$  for  $s$ ’s neighbor  $v \in N(s)$ , the information about  $\tau(v)$ ,  $\rho(v)$  and  $d(v)$  is required where  $\tau(v)$  and  $\rho(v)$  are the number of  $v$ ’s posts shared by  $v$ ’s neighbors and the number of  $v$ ’s all posts, respectively.
- **Hybrid selection:** Pick the  $\min(k, d(s))$  nodes  $v \in V$  with the highest *weighted* node degree  $\omega d(v)$  which is defined as  $\omega d(v) = \omega(v) \cdot d(v)$ . At the first glance, this scheme requires the knowledge of both the degree and the user propagation weight of neighbors. In fact, however, this scheme can be simply implemented without the knowledge about node degree since  $\omega(v) \cdot d(v)$  is calculated as  $\tau(v)/\rho(v)$ ;  $d(v)$  is automatically canceled in the calculation.

We note that these schemes seem the most reasonable, since we cannot calculate network centrality metrics, such as closeness and betweenness [13], which require the knowledge of the entire network topology. Here, we do not consider the other metrics (e.g., [21]) to estimate node centrality based on localized information alone since previous work [14] already showed that these metrics are ineffective for the INS problem compared with node degree.

The communication costs of all these schemes are  $O(d(s))$  since the spreader  $s$  can obtain  $d(v)$ ,  $\omega(v)$  or  $\omega d(v)$  through only direct communications with each neighbor  $v \in N(s)$ .

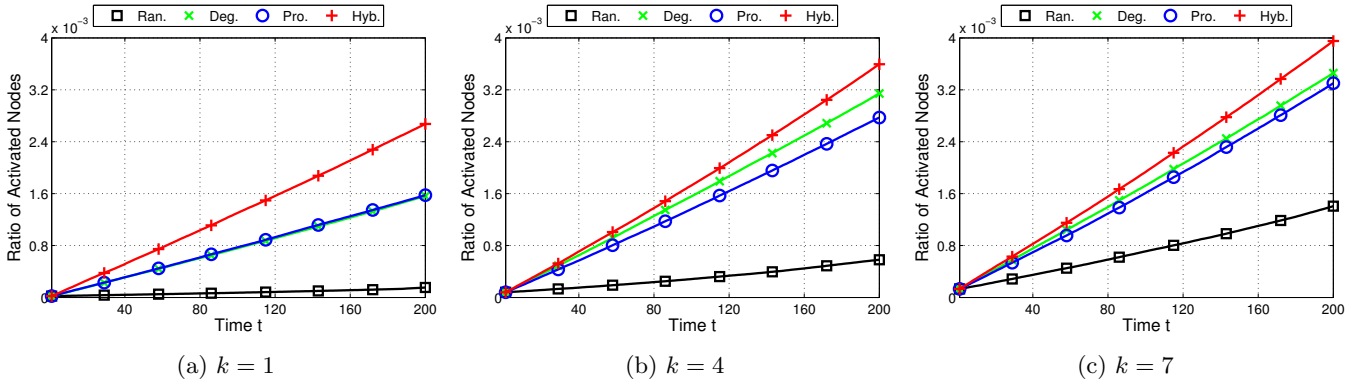
## 4. EXPERIMENTAL RESULTS

In this section, we analyze the performance of the selection schemes presented in Section 3.

For experiments, we used the Twitter dataset [1] related to the 2010 UK general election between the 5th and 12th of May since this dataset reflects typical behavior of information diffusion in a political campaign.

To remove insignificant test cases, we filtered out users who did not either write any posts or had any followers. The used graph consists of 45,179 nodes and 1,938,734 edges representing a sub-network of Twitter. Also, this graph has the following properties: (1) its average degree is 42.91; (2) its number of strongly connected components is 4559; and (3) the number of weakly connected components is 18 (i.e., this graph is divided into 18 disconnected components).

In these experiments, our goal was to find the best neighbor selection scheme to maximize information diffusion in Twitter. Unlike previous studies that used network topology alone [14] or synthetic datasets [12], we used real propagation rate  $\omega(u)$  for each node  $u$  calculated from the Twitter messages.



**Figure 2: Changes in the ratio of the average number of activated nodes to the total number of nodes in the network over time  $t$ .**

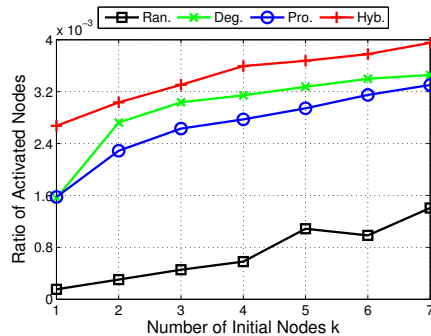
We used the IC model described in Section 2 to evaluate the performance of the schemes presented in Section 3, with varying the number of initially activated neighbors  $k$ . The propagation probability  $\lambda(u, v, s, r)$  on an edge  $(u, v) \in E$  was defined with the spreader  $s \in V$  and a piece of information  $r$  described in Section 2.

In each simulation run, we randomly picked a spreader with a piece of information  $r$  and then selected its  $k$  neighbors according to a selection criterion presented in Section 3, where the content interestingness  $\phi(r)$  was randomly drawn from the normal distribution with the mean of 0.0136 and standard deviation of 0.0501, according to real data [22]. We also set the decay factor  $\gamma$  from 0.2 to 1.0.

For evaluation, we observed the changes in the number of activated nodes during the 200th time steps. With a fixed  $k$ , we repeated this 500 times to minimize the bias of the test samples (randomly selected spreaders); we measured the ratio of the average number of activated nodes per test sample to the total number of nodes in the network. To establish a fair comparison, the parameter values were the same for all selection schemes in the  $i$ th run. Figure 2 shows how these values are changed over time  $t$  with  $k = 1, 4$  or  $7$  and  $\gamma = 0.2$  under the IC model.

From this figure, we can see that the Hybrid selection scheme outperformed the other selection schemes: When  $k = 1$ , in the Hybrid selection scheme, the ratio of the average number of activated nodes to the total number of nodes were over 0.0024 while the ratios were around 0.0016 in Degree and Propagation-weight selection schemes. As  $k$  increased to 7, the gap between Hybrid and other selection schemes was rather reduced, but still seemed significant. This shows that we can effectively spread information using the Hybrid scheme, even without consideration of the “number of friends” information since the node degree is not needed to use the Hybrid scheme. Interestingly, the Degree selection was slightly better than the Propagation-weight selection when  $k = 4$  or  $7$ , while these schemes produced almost the same results when  $k = 1$ .

Paired one-tailed t-tests with  $\alpha = 0.01$  were used to compare the performance of the neighbor selection schemes in a statistically significant manner. We tested whether the distributions of the numbers of the activated nodes between schemes after the final time step (i.e., the 200th) were statistically different. The test results show that the performance



**Figure 3: Changes in the ratio of the average number of activated nodes to the total number of nodes in the network with the number of initially activated neighbors  $k$ .**

of all the schemes appeared to be significantly different, except for the comparison of Propagation-weight and Degree when  $k = 1$  ( $p$ -value = 0.5202). These test results are quite different from those in the previous study [12] using synthetic parameters where Hybrid and Propagation-weight selection schemes achieved almost the same performance results.

We now discuss how the performance of the different neighbor selection schemes may change with the number of initially activated nodes  $k$ . To accelerate the speed of information diffusion, a possible straightforward approach is to increase the number of initially activated neighbors  $k$ . Probably, we can imagine that even the naive Random selection scheme can also be used to efficiently disseminate a piece of information if  $k$  increases sufficiently.

To demonstrate the effects of  $k$ , we analyzed the ratio of the average number of activated nodes after the 200th time steps with  $k$  ranging from 1 to 7. The experimental results are shown in Figure 3.

Unsurprisingly, the performance of all selection schemes overall improved as  $k$  increased. The ratios of activated nodes, except the Random selection scheme, show almost a similar pattern – the curves commonly had gentle slopes. Although the Random selection scheme was relatively highly affected by  $k$ , the average number of activated nodes in the

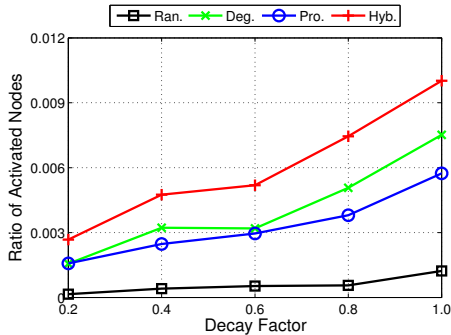


Figure 4: Changes in the ratio of the average number of activated nodes to the total number of nodes in the network with the decay factor  $\gamma$ .

Random selection scheme was still below 0.0016 even for  $k = 7$ . We can also see that the performance gaps between the schemes still existed with  $k$ . The performance of the Hybrid selection scheme, when  $k = 4$ , was better than the other schemes even when  $k = 7$ . We note that these results are quite different from previous work [12] where the effects of  $k$  were rather limited – in all selection strategies, the number of activated nodes were not greatly increased with  $k \geq 3$ .

Finally, we now discuss the effects of the decay factor  $\gamma$  presented in Section 2. We performed simulation with the same parameters, while varying  $\gamma$  from 0.2 to 1. The experimental results are shown in Figure 4. We use a different y-axis scale on this figure since the numbers of activated nodes were greatly increased with  $\gamma$ .

From this figure, we can see that the performance of all selection schemes generally improved and the gaps between the schemes grew with  $\gamma$ . Thus our suggestion is to use the Hybrid selection scheme even with a large decay factor  $\gamma$ . Interestingly, we can also observe two different patterns in Figure 4: one is for the Hybrid and Degree selection schemes, which tends to increase quickly when  $\gamma = 0.6$ , and the other one is for the Propagation-weight and Random selection schemes, which tends to increase relatively slowly.

## 5. RELATED WORK

*Influential Maximization* (IM) problem has recently received increasing attention, given the growing popularity of online social networks, such as Facebook and Twitter, which have provided great opportunities for the diffusion of information, opinions and adoption of new products.

The IM problem was originally introduced for marketing purposes by Domingos and Richardson [7]: The goal is to find a set of  $k$  initially activated nodes with the maximum number of activated nodes after the time step  $t$ . Kempe et al. [11] formulated this problem under two basic stochastic influence cascade models: the *Independent Cascade* (IC) model [8] and the *Linear Threshold* (LT) model [11]. In the IC model, each edge has a propagation probability and influence is propagated by activated nodes independently activating their inactive neighbors based on the edge propagation probabilities. In the LT model, each edge has a weight, each node has a threshold chosen uniformly at random, and a node becomes activated, if the weighted sum of its activated neighbors exceeds its threshold. Kempe et al. [11]

showed that the optimization problem of selecting the most influential nodes in a graph is NP-hard for both models, and also proposed a greedy algorithm that provides a good approximation ratio of 63% of the optimal solution. However, their greedy algorithm relies on the Monte-Carlo simulation on influence cascade to estimate the influence spread, which makes the algorithm slow and not scalable.

A number of papers in recent years have tried to overcome the inefficiency of this greedy algorithm by improving the original algorithm [16, 6] or proposing new algorithms [15, 6, 5]. Leskovec et al. [16] proposed the *Cost-Effective Lazy Forward* (CELFF) scheme in selecting new seeds to reduce the number of influence spread evaluations, but it is still slow and not scalable to large graphs, as demonstrated in [5]. Kimura and Saito [15] proposed shortest-path based heuristic algorithms to evaluate the influence spread. Chen et al. [6] proposed two faster greedy algorithms called *Mixed-Greedy* and *DegreeDiscount* for the IC model where the propagation probabilities on all edges are the same; MixedGreedy removes the edges that have no contribution for propagating influence, which can reduce the computation on the unnecessary edges; DegreeDiscount assumes that the influence spread increases with node degree. Chen et al. [5] proposed the *Maximum Influence Arborescence* (MIA) heuristic based on local tree structures to reduce computation costs. Wang et al. [20] proposed a community-based greedy algorithm for identifying the most influential nodes. The main idea is to divide a social network into communities, and estimate the influence spread in each community instead of the whole network topology.

As a variant of the conventional IM problem, Kim and Yoneki [14] introduced the problem called *Influential neighbor selection* (INS) to select the most influential neighbors of a node, rather than the most influential arbitrary nodes in a network. More recently, Kim [12] extended this epidemiological model by introducing several parameters (*user propagation weight*, *content interestingness*, and *decay factor*) to provide a more general and practical information diffusion model. However, they still used synthetic parameters that might significantly affect the information diffusion process. In this paper, we found that the information diffusion speed of some schemes (e.g., Propagation-weight and Random) in the previous study [12] was quite overestimated.

Many studies noted that the levels of information sharing activity varied greatly between users in social networks. Romero et al. [18] argued that a majority of Twitter users might be passive, not engaging in creating and sharing information. Cha et al. [4] found that users with many followers do not necessarily influence in terms of spawning retweets or mentions – the Spearman’s rank correlation coefficient between the “ranking by followers” and “ranking by retweets” for all users was 0.549. Zhou et al. [22] showed that in Twitter, the content of a tweet might be an important factor in determining the “retweet rate” – the mean retweet rate was 0.0136 but standard deviation was as high as 0.0501. Also, they observed that cascades tend to be wide and not too deep suggesting that the retweet rate may decay as the cascades spreads away from the source – the mean of decay factors was about 0.2.

## 6. CONCLUSIONS

Given the increasing popularity of online social networking services, there has been growing interest in investigating

the characteristics of epidemic spreading, in order to accelerate or mitigate it. Kim and Yoneki [14] introduced the optimization problem to find *influential neighbors* for maximizing information diffusion. We have extended their work by introducing several important parameters (*user propagation weight*, *decay factor* and *content interestingness*) to provide a more general and practical information diffusion model.

We presented four neighbor selection schemes (**Random**, **Degree**, **Propagation-weight** and **Hybrid** selection) and explored their feasibility. We compared these selection schemes by computing the ratio of the average number of activated nodes to the total number of nodes in the network. We discussed which selection methods are generally recommended under which conditions.

Our experimental results showed that the **Hybrid** selection scheme produced the best results of maximizing information diffusion through intensive simulation. Even with a small  $k$ , the **Hybrid** selection scheme outperformed the other selection schemes with a relatively large  $k$ . Since the **Hybrid** selection scheme can use the information about users' posts alone, we can efficiently spread information without the information about the "number of friends" each user has.

As an extension to this work, we are considering a theoretical study to formally generalize and verify our results in order to consider a wide range of application environments (e.g., each of which will have different levels of content interestingness). We will also develop a more extended framework for information diffusion. We may consider not only a spreader with the knowledge about user's neighbors but also a spreader with a partial knowledge of network topology (e.g., a subset of users or neighbors of neighbors). For example, we will extend the concept of the INS problem by expanding the set of the initially activated nodes with the distance from the information spreader.

Another interesting problem is to consider a new problem in the opposite direction to prevent (or reduce) the spread of information (e.g., rumor) by carefully monitoring (important) users with a high "user propagation weight" and/or "number of friends".

## 7. ACKNOWLEDGMENTS

The research is funded in part by EU grant FP7-ICT-318398 eCOUSIN project.

## 8. REFERENCES

- [1] A. Boutet, H. Kim, and E. Yoneki. What's in Twitter, I know what parties are popular and who you are supporting now! *Social Network Analysis and Mining*, 3(4):1379–1391, 2013.
- [2] J. J. Brown and P. H. Reingen. Social ties and Word-of-Mouth referral behavior. *The Journal of Consumer Research*, 14(3):350–362, 1987.
- [3] Q. Cao, M. Sirivianos, X. Yang, and T. Pregueiro. Aiding the detection of fake accounts in large scale social online services. In *Proceedings of the 9th USENIX Conference on Networked Systems Design and Implementation (NSDI)*, 2012.
- [4] M. Cha, H. Haddadi, F. Benevenuto, and K. Gummadi. Measuring user influence in twitter: The million follower fallacy. In *Proceedings of the 4th AAAI Conference on Weblogs and Social Media (ICWSM)*, 2010.
- [5] W. Chen, C. Wang, and Y. Wang. Scalable influence maximization for prevalent viral marketing in large-scale social networks. In *Proceedings of the 16th ACM Conference on Knowledge Discovery and Data mining (KDD)*, 2010.
- [6] W. Chen, Y. Wang, and S. Yang. Efficient influence maximization in social networks. In *Proceedings of the 15th ACM Conference on Knowledge Discovery and Data mining (KDD)*, 2009.
- [7] P. Domingos and M. Richardson. Mining the network value of customers. In *Proceedings of the 7th ACM Conference on Knowledge Discovery and Data mining (KDD)*, 2001.
- [8] J. Goldenberg, B. Libai, and E. Muller. Talk of the Network: A Complex Systems Look at the Underlying Process of Word-of-Mouth. *Marketing Letters*, pages 211–223, 2001.
- [9] D. Gruhl, R. Guha, D. Liben-Nowell, and A. Tomkins. Information diffusion through blogspace. In *Proceedings of the 13th international conference on World Wide Web (WWW)*, 2004.
- [10] L. K. Hansen, A. Arvidsson, F. A. Nielsen, E. Colleoni, and M. Etter. Good friends, bad news - affect and virality in twitter. In *Future Information Technology*, volume 185 of *Communications in Computer and Information Science*, pages 34–43. Springer Berlin Heidelberg, 2011.
- [11] D. Kempe, J. Kleinberg, and E. Tardos. Maximizing the spread of influence through a social network. In *Proceedings of the 9th ACM Conference on Knowledge Discovery and Data mining (KDD)*, 2003.
- [12] H. Kim. Don't Count The Number of Friends When You Are Spreading Information in Social Networks. In *Proceedings of the 8th International Conference on Ubiquitous Information Management and Communication (ICUIMC)*, 2014.
- [13] H. Kim, J. Tang, R. Anderson, and C. Mascolo. Centrality prediction in dynamic human contact networks. *Computer Networks*, 56(3):983–996, 2012.
- [14] H. Kim and E. Yoneki. Influential neighbours selection for information diffusion in online social networks. In *Proceedings of the 21th International Conference on Computer Communication Networks (ICCCN)*, 2012.
- [15] M. Kimura and K. Saito. Tractable models for information diffusion in social networks. In *Proceedings of the 10th European conference on Principle and Practice of Knowledge Discovery in Databases (PKDD)*, 2006.
- [16] J. Leskovec, A. Krause, C. Guestrin, C. Faloutsos, J. VanBriesen, and N. Glance. Cost-effective outbreak detection in networks. In *Proceedings of the 13th ACM Conference on Knowledge Discovery and Data mining (KDD)*, 2007.
- [17] P. Metaxas, E. Mustafaraj, and D. Gayo-Avello. How (not) to predict elections. In *Proceedings of the 3rd international conference on Privacy, Security, Risk and Trust (PASSAT) and the 3rd international conference on Social Computing (SocialCom)*, 2011.
- [18] D. M. Romero, W. Galuba, S. Asur, and B. A. Huberman. Influence and passivity in social media. In *Proceedings of the 20th international conference on World Wide Web (WWW)*, 2011.
- [19] K. Thomas, C. Grier, and V. Paxson. Adapting social spam infrastructure for political censorship. In *Proceedings of the 5th USENIX Conference on Large-Scale Exploits and Emergent Threats (LEET)*, 2012.
- [20] Y. Wang, G. Cong, G. Song, and K. Xie. Community-based greedy algorithm for mining top-k influential nodes in mobile social networks. In *Proceedings of the 16th ACM Conference on Knowledge Discovery and Data mining (KDD)*, 2010.
- [21] K. Wehmuth and A. Ziviani. Daccer: Distributed assessment of the closeness centrality ranking in complex networks. *Computer Networks*, 57(13):2536–2548, 2013.
- [22] Z. Zhou, R. Bandari, J. Kong, H. Qian, and V. Roychowdhury. Information resonance on twitter: watching iran. In *Proceedings of the first ACM Workshop on Social Media Analytics (SOMA)*, 2010.