

“

**t**

7

7,

١

, *Member, IEEE,*

, *Senior Member, IEEE*

## Abstract—

**Abstract—**

### Index Terms—

## 1 INTRODUCTION

$$I$$

- J. Zhang and Z. Fang are with the Department of Computer Science and Technology, Tsinghua University, Beijing 100084, China.  
E-mail: {zhangjing12, fzp13}@mails.tsinghua.edu.cn.
- W. Chen is with the Theory Group, Microsoft Research, Beijing 100080, China. E-mail: weic@microsoft.com.
- J. Tang is with the Department of Computer Science and Technology, Tsinghua University, Beijing 100084, China, and Tsinghua National Laboratory for Information Science and Technology (TNList).  
E-mail: jietang@tsinghua.edu.cn.

*Manuscript received 11 Jan. 2014; revised 5 Jan. 2015; accepted 9 Feb. 2015.  
Date of publication 25 Feb. 2015; date of current version 2 July 2015.*

*Recommended for acceptance by G. Das.*

For information on obtaining reprints of this article, please send e-mail to: [reprints@ieee.org](mailto:reprints@ieee.org), and reference the Digital Object Identifier below.

Digital Object Identifier no. 10.1109/TKDE.2015.2407351

## 2 “FOLLOWING” LINK CASCADE MODEL

Let  $G = (V, E, t)$  be a directed graph with vertex set  $V$ , edge set  $E$ , and a time function  $t: E \rightarrow \mathbb{N} \cup \{\perp\}$ . For each edge  $e_{uv} \in E$ , let  $t(e_{uv}) = n$  denote the time when the edge was discovered. If  $t(e_{uv}) = \perp$ , then the edge has not been discovered yet.

Let  $A$  and  $B$  be two sets of vertices. Let  $C$  be a set of vertices. Let  $t(e_{uv}) = n$  denote the time when the edge was discovered. If  $t(e_{uv}) = \perp$ , then the edge has not been discovered yet.

Let  $t(e_{uv}) = n$  denote the time when the edge was discovered. If  $t(e_{uv}) = \perp$ , then the edge has not been discovered yet.

Let  $t(e_{uv}) = n$  denote the time when the edge was discovered. If  $t(e_{uv}) = \perp$ , then the edge has not been discovered yet.

Let  $t(e_{uv}) = n$  denote the time when the edge was discovered. If  $t(e_{uv}) = \perp$ , then the edge has not been discovered yet.

**A.1. Diffusion effect between links decays over time.**

Let  $t(e_{uv}) = n$  denote the time when the edge was discovered. If  $t(e_{uv}) = \perp$ , then the edge has not been discovered yet.

Let  $t(e_{uv}) = n$  denote the time when the edge was discovered. If  $t(e_{uv}) = \perp$ , then the edge has not been discovered yet.

Let  $t(e_{uv}) = n$  denote the time when the edge was discovered. If  $t(e_{uv}) = \perp$ , then the edge has not been discovered yet.

Let  $t(e_{uv}) = n$  denote the time when the edge was discovered. If  $t(e_{uv}) = \perp$ , then the edge has not been discovered yet.

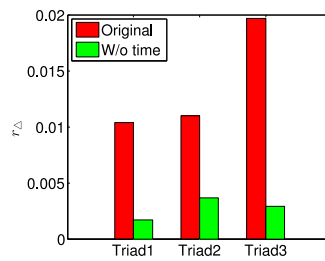
Organization.



### 3.1 Data Collection

### 3.2 Observations

$$r_{\Delta} = \frac{|C_{\Delta}^{+}|}{|C_{\Delta}|}. \quad (1)$$



(

[illegible]

$$\begin{pmatrix} + & - \\ \xi & \xi \end{pmatrix} \begin{matrix} A & C \\ C & B \end{matrix} \begin{matrix} \xi \\ \xi \end{matrix} \quad C \begin{pmatrix} + & - \\ 0 & 0 \end{pmatrix}.$$

## 4 MODEL LEARNING

Let  $\mathcal{D}$  be a dataset of  $n$  i.i.d. samples  $(e', e)$  from a joint distribution  $p(e', e)$ .

*Likelihood function.* For a parameter vector  $\theta$ , the likelihood function is defined as

$$L(\theta) = \prod_{(e', e) \in \mathcal{D}} p(e', e; \theta) = \prod_{(e', e) \in \mathcal{D}} \exp\left(\sum_{\xi, \eta} \theta_{\xi\eta} \xi_{e'} \eta_e\right).$$

Let  $\theta = \{h_{e'e}, g_{e'e}\}$  be the parameters of the model. The log-likelihood function is then

$$\ell(\theta) = \log L(\theta) = \sum_{(e', e) \in \mathcal{D}} \log p(e', e; \theta) = \sum_{(e', e) \in \mathcal{D}} \sum_{\xi, \eta} \theta_{\xi\eta} \xi_{e'} \eta_e.$$

The maximum likelihood estimate (MLE) of  $\theta$  is the parameter vector that maximizes the log-likelihood function:

$$\hat{\theta} = \arg \max_{\theta} \ell(\theta).$$

$$e \qquad \qquad \qquad y_{e'e} \qquad \qquad \qquad \mathfrak{f} \ e' \\ [t_{e'}, t_e], \qquad \qquad \qquad \mathfrak{f} \qquad \qquad \qquad t_{e'} \\ t_e, e' \qquad \qquad \qquad e \qquad \qquad \qquad \mathfrak{f} \ . \mathbb{I} \qquad \qquad \qquad , y_{e'e} \\ e' \qquad \qquad \qquad e \ \mathfrak{f} \qquad \qquad \qquad t_e$$

$$\begin{aligned} y_{e'e} &= 1 - h_{\Delta} g_{\Delta} \sum_{t=t_{e'}}^{t_e} (1 - g_{\Delta})^{t-t_{e'}} \\ &= h_{\Delta} (1 - g_{\Delta})^{t_e-t_{e'}+1} + (1 - h_{\Delta}). \end{aligned} \tag{6}$$

$$\begin{array}{c} \mathfrak{f} \qquad \qquad \qquad \mathfrak{f} \qquad \qquad \qquad \mathfrak{f}\mathfrak{f} \\ \mathfrak{f} \qquad \qquad \qquad \cdot \qquad \qquad \qquad e \in \mathcal{E} \\ \qquad \qquad \qquad \mathfrak{f} \\ \delta \qquad \qquad \qquad \mathbb{I} \ \mathfrak{f} \qquad \qquad \qquad - \\ e' \qquad \qquad \qquad y_{ee'} \ \mathfrak{f} \ e' \in R_{e'} \qquad \qquad \qquad R_e \\ \mathfrak{f} \qquad \qquad \qquad \mathfrak{f} \ e \qquad \qquad \qquad \mathfrak{f} \qquad \qquad \qquad t_e + \delta. \qquad \qquad \qquad - \\ y_{ee'} \qquad \qquad \qquad \cdot \ (\ ) , \qquad \qquad \qquad t_{e'} \\ t_e \qquad \qquad \qquad t_e \qquad \qquad \qquad t_e + \delta \qquad \qquad \qquad \cdot \\ \qquad \qquad \qquad - \qquad \qquad \qquad \mathfrak{f} \end{array}$$

$$\log \mathcal{L} = \sum_{e \in \mathcal{E}} \left\{ \log \sum_{\vec{\alpha}_{S_e}} \prod_{e' \in S_e} x_{e'e}^{\alpha_{e'}} y_{e'e}^{1-\alpha_{e'}} + \sum_{e' \in R_e} \log y_{ee'} \right\}.$$

EM algorithm.

$$\begin{array}{c} \cdot \\ q(e|\vec{\alpha}_{S_e}) = \frac{p(e|\vec{\alpha}_{S_e})}{\mathfrak{f}} \\ , \\ - \qquad \qquad \qquad \mathfrak{f} \end{array}$$

$$\begin{aligned} \log \mathcal{L} &= \sum_{e \in \mathcal{E}} \left\{ \log \sum_{\vec{\alpha}_{S_e}} \hat{q}(e|\vec{\alpha}_{S_e}) \frac{p(e|\vec{\alpha}_{S_e})}{\hat{q}(e|\vec{\alpha}_{S_e})} + \sum_{e' \in R_e} \log y_{ee'} \right\} \\ &\geq \sum_{e \in \mathcal{E}} \left\{ \sum_{\vec{\alpha}_{S_e}} \hat{q}(e|\vec{\alpha}_{S_e}) \log \frac{p(e|\vec{\alpha}_{S_e})}{\hat{q}(e|\vec{\alpha}_{S_e})} + \sum_{e' \in R_e} \log y_{ee'} \right\}, \\ &\qquad \qquad \qquad \hat{\phantom{q}} \qquad \qquad \qquad \mathfrak{f} \qquad \qquad \qquad - \\ &\qquad \qquad \qquad \hat{q}(e|\vec{\alpha}_{S_e}) \log \hat{q}(e|\vec{\alpha}_{S_e}) \\ &\qquad \qquad \qquad \mathfrak{f} \qquad \qquad \qquad , \end{aligned}$$

$$Q(\theta, \hat{\theta})$$

$$Q(\theta, \hat{\theta}) = \sum_{e \in \mathcal{E}} \left\{ \sum_{\vec{\alpha}_S}$$

$$h_{\Delta} = \frac{\sum_{(e',e) \in C_{\Delta}^{+}} \hat{D}_{e'e} + \sum_{(e',e) \in C_{\Delta}^{-}} \hat{B}_{e'e}}{|C_{\Delta}|}, \quad (12)$$

$$g_{\Delta} = \frac{\sum_{(e',e) \in C_{\Delta}^{+}} \hat{A}_{e'e}}{\sum_{(e',e) \in C_{\Delta}^{-}} \hat{B}_{e'e}(\delta + 1) + \sum_{(e',e) \in C_{\Delta}^{+}} \hat{D}_{e'e}(t_e - t_{e'} + 1)}. \quad (13)$$

$$D_{e'e} = B_{e'e} + A_{e'e} - A_{e'e}B_{e'e}. \quad (14)$$

Algorithm 1.	
$G = (V, E, t)$	
$\theta = \{h_{\Delta}, g_{\Delta}\}$	$(0, )$
$E$	$e \in \mathcal{E}$
$e' \in S_e$	$x_{e'e} \leftarrow 0$
$e' \in S_e$	$y_{e'e} \leftarrow 0$
$e' \in S_e$	$A_{e'e} \leftarrow 0$
$e' \in S_e$	$B_{e'e} \leftarrow 0$
$e' \in S_e$	$D_{e'e} \leftarrow 0$
$e' \in R_e$	$B_{ee'} \leftarrow 0$
$\Delta = 1$	$24$
$h_{\Delta}$	$( )$
$g_{\Delta}$	$( )$
Convergence	

## 5 APPLICATIONS

Followee maximization.

Let  $S$  be a set of nodes in  $G$ . The goal is to find a set of nodes  $S$  that maximizes the number of edges from  $S$  to  $V \setminus S$ . This is a classic problem in graph theory, and it is known to be NP-hard. However, there are several approximation algorithms that can be used to find a set of nodes  $S$  that is close to optimal.

One such algorithm is the greedy algorithm, which starts with an empty set  $S$  and iteratively adds nodes to  $S$  until no more nodes can be added without decreasing the number of edges from  $S$  to  $V \setminus S$ . This algorithm is guaranteed to find a set of nodes  $S$  that is at least half as large as the optimal set.

Another algorithm is the local search algorithm, which starts with a set of nodes  $S$  and iteratively replaces nodes in  $S$  with nodes outside of  $S$  until no more improvements can be made. This algorithm is also guaranteed to find a set of nodes  $S$  that is at least half as large as the optimal set.

Followee maximization.

Let  $S$  be a set of nodes in  $G$ . The goal is to find a set of nodes  $S$  that maximizes the number of edges from  $S$  to  $V \setminus S$ . This is a classic problem in graph theory, and it is known to be NP-hard. However, there are several approximation algorithms that can be used to find a set of nodes  $S$  that is close to optimal.

One such algorithm is the greedy algorithm, which starts with an empty set  $S$  and iteratively adds nodes to  $S$  until no more nodes can be added without decreasing the number of edges from  $S$  to  $V \setminus S$ . This algorithm is guaranteed to find a set of nodes  $S$  that is at least half as large as the optimal set.

Another algorithm is the local search algorithm, which starts with a set of nodes  $S$  and iteratively replaces nodes in  $S$  with nodes outside of  $S$  until no more improvements can be made. This algorithm is also guaranteed to find a set of nodes  $S$  that is at least half as large as the optimal set.

Algorithm 2.	
$G = (V, E), v, k$	
$S = \emptyset$	$R = 0,000$
$i = 1$ to $k$	$u \in V \setminus S$
$s_u = 0$	$r = 1$ to $R$
$s_u + =  FCM(S \cup \{u\}) $	$s_u = s_u / R$
$S = S \cup \{argmax_{u \in V \setminus S} s_u\}$	

## 6 EXPERIMENTS

### 6.1 Experimental Setup

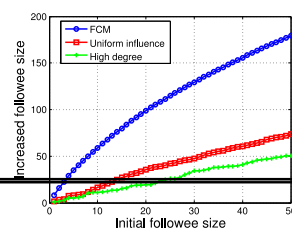
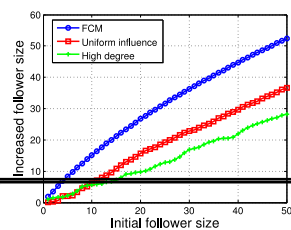
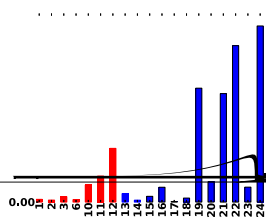
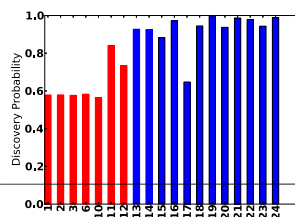
The experiments were conducted on a set of graphs generated using a random graph model. The graphs were generated with a fixed number of nodes  $n$  and a fixed number of edges  $m$ . The number of nodes  $n$  was varied from 100 to 1000, and the number of edges  $m$  was varied from 100 to 1000. The graphs were generated using a random graph model, and the edges were generated independently of each other.

The experiments were conducted on a set of graphs generated using a random graph model. The graphs were generated with a fixed number of nodes  $n$  and a fixed number of edges  $m$ . The number of nodes  $n$  was varied from 100 to 1000, and the number of edges  $m$  was varied from 100 to 1000. The graphs were generated using a random graph model, and the edges were generated independently of each other.





$$\begin{array}{ccccccc}
\xi & & u, & & & & \\
& & & 0 \cdot V & & \xi \xi & \\
\xi & & u, & & \xi & & \xi \\
& & u. & & & & \xi - \\
\xi & & & \xi \alpha & \beta \xi & 0 & 0. \cdot \\
& & & & & & -
\end{array}$$



*Per-triad analysis.*

*Delay analysis.*

*Convergence analysis.*

### Model parameter analysis.

### 6.3 Application Improvement



## REFERENCES

- Proc. 14th ACM SIGKDD Int. Conf. Knowl. Discovery Data Mining, 00, .
- Proc. 4th ACM Int. Conf. Web Search Data Mining, 0, .
- Science, .
- Proc. IEEE 12th Int. Conf. Data Mining, 0, .
- Nature, .
- Proc. 11th SIAM Int. Conf. Data Mining, 0, .
- Proc. 16th ACM SIGKDD Int. Conf. Knowl. Discovery Data Mining, 00, .
- Proc. 15th ACM SIGKDD Int. Conf. Knowl. Discovery Data Mining, 00, .
- Phys. Rev. E, . 0, .
- Proc. 16th Int. Conf. World Wide Web, 00, .
- Proc. 7th ACM SIGKDD Int. Conf. Knowl. Discovery Data Mining, 00, .
- Networks, Crowds, and Markets: Reasoning about a Highly Connected World, .
- Tsinghua Sci. Technol., .
- Permutation, Parametric and Bootstrap Tests of Hypotheses, .
- Proc. 3rd ACM Int. Conf. Web Search Data Mining, 00, .
- Amer. J. Sociol., .
- Proc. 13th Int. Conf. World Wide Web, 00, .
- Proc. 11th SIAM Int. Conf. Data Mining Workshop, 00, .
- Proc. 8th ACM SIGKDD Int. Conf. Knowl. Discovery Data Mining, 00, .
- Psychometrika, .
- Proc. 9th ACM SIGKDD Int. Conf. Knowl. Discovery Data Mining, 00, .
- Proc. Roy. Soc. A, . 00
- Proc. ACM Conf. Inf. Knowl. Manage., 00, .
- Intell. Data Anal., .
- Proc. 19th ACM SIGKDD Int. Conf. Knowl. Discovery Data Mining, 0, .
- Proc. 12th ACM SIGKDD Int. Conf. Knowl. Discovery Data Mining, 00, .
- Proc. 21st Int. Conf. World Wide Web, 00, .
- Proc. 11th ACM SIGKDD Int. Conf. Knowl. Discovery Data Mining, 00, .
- Proc. 13th ACM SIGKDD Int. Conf. Knowl. Discovery Data Mining, 00, .
- Math. Biosci., .
- Proc. ACM Int. Conf. Web Search Data Mining, 0, .
- J. Amer. Soc. Inf. Sci. Technol., .
- Proc. 12th ACM SIGKDD Int. Conf. Knowl. Discovery Data Mining, 00, .
- Proc. 16th ACM SIGKDD Int. Conf. Knowl. Discovery Data Mining, 00, .
- Data Mining Knowl. Discovery, .
- Sci. Technol., .
- ACM Trans. Knowl. Discovery Data, .
- Proc. 18th ACM SIGKDD Int. Conf. Knowl. Discovery Data Mining, 0, .
- Phys. Rev. E, . 00, .
- Computer Intensive Methods for Testing Hypotheses, .
- Proc. 4th Int. AAAI Conf. Weblogs Social Media, 00, .
- Int. J. Semantic Web Inf. Syst., .
- Proc. 17th ACM SIGKDD Int. Conf. Knowl. Discovery Data Mining, 0, .
- Proc. 15th ACM SIGKDD Int. Conf. Knowl. Discovery Data Mining, 00, .
- Proc. 19th ACM SIGKDD Int. Conf. Knowl. Discovery Data Mining, 0, .
- Proc. 6th Int. Conf. Data Mining, 00, .
- Proc. 7th IEEE Int. Conf. Data Mining, 00, .

