

香港中文大學 The Chinese University of Hong Kong

Analyzing Competitive Influence Maximization Problems with Partial Information

Yishi Lin, John C.S. Lui The Chinese University of Hong Kong

Background: Word-of-Mouth



In social interactions, we **influence** each other.

Background: Viral Marketing

- Assumption: the *word-of-mouth* effect
- Idea: exploiting the *social influence* for marketing
- Targeting "influencers" who are likely to produce the word-of-mouth diffusion



Background: Classical Independent Cascade model

- Single source Independent Cascade (IC) model (Kempe et al. KDD'03)
 - Initially, a set of "seed nodes" *S* are activated.
 - Influenced node u influences its neighbor v with probability p_{uv} .
 - Influence spread $\sigma(S)$: the expected number of influenced nodes



Kempe, David, Jon Kleinberg, and Éva Tardos. "Maximizing the spread of influence through a social network." *Proceedings* of the ninth ACM SIGKDD international conference on Knowledge discovery and data mining. ACM, 2003.

Background: Classical Influence Max. Problem

Input: *G* and *k*



Problem (Influence Maximization) Select k seed nodes so to maximize the expected spread of influence.

Output:

Seed set of size k



Under the IC model:

- The IM problem is NP hard. oxtimes
- Even computing $\sigma(S)$ is #P hard. \otimes

Motivation

Competition among products



• Partial information: It is not always possible to have full information about viral marketing strategies of the competitor.

Main Contributions

General Competitive Independent Cascade Model

- Many specific models proposed previously are its special cases
 - Distance-based model (Carnes et al. ICEC'2007)
 - Wave propagation model (Carnes et al. ICEC'2007)
 - Campaign-Oblivious Independent Cascade model (Budak et al. WWW'11)

General Competitive Influence Maximization Problem

• Assuming only partial knowledge about competitor's seeding strategy

General algorithmic framework

- It solves the general problem.
- It works for any specific instances of the general model.

Model General Competitive Independent Cascade Model

- **Network** G = (V, E):
 - Every edge e_{uv} is associated with a probability p_{uv} .
- **Sources:** two competing sources *A* and *B*.
- State of a node: *Susceptible*, *Inf_A* or *Inf_B*
 - "Influenced" cannot change its state.
- Seeds / initial adopters: $S_A \subseteq V$, $S_B \subseteq V$
 - We assume $S_A \cap S_B = \emptyset$.



Model

General Competitive Independent Cascade Model

• **Given seeds:** S_A and S_B

- Determine propagation results
 - Active edges E_a : edge e_{uv} is "active" w.p. p_{uv} .
 - Node u will be in the same state as that of *one of* its *nearest seeds* in $G = (V, E_a)$.
 - A specific model should specify how the influence propagates in detail.
- The expected influence
 - $\sigma(S_B|S_A) = \mathbb{E}_{E_a}[\# \text{ of nodes in state } Inf_B]$
- Assumption
 - \circ monotonicity and submodularity of $\sigma(S_B|S_A)$



 $S_A = \{3\} S_B = \{5\}$

Problem Definition

Competitive Influence Maximization problem with Partial information (CIMP)

Input:

- *G*, *k*, propagation model
- Competitor's seed distribution \mathcal{D}_A



Problem

Select a set S_B^* of k nodes so that the expected spread of influence of source B under the presence of competitors, $\sigma(S_B^*|\mathcal{D}_A)$, is maximized. = $\mathbb{E}_{S_A \sim \mathcal{D}_A} [\sigma(S_B^*|S_A)]$

Monotone & Submodular

Output: Seed set S_B of size k



Problem Definition

Competitive Influence Maximization problem with Partial information (CIMP)

The CIMP problem is **NP hard**. $\ensuremath{\mathfrak{S}}$ Even computing $\sigma(S_B | \mathcal{D}_A)$ is **#P hard**. $\ensuremath{\mathfrak{S}}$

Solution Two-phase Competitive Influence Maximization (TCIM)

TCIM: Estimating the Expected Influence

Random <u>Reverse</u> <u>A</u>ccessible <u>P</u>ointed <u>G</u>raph (RAPG)

Input:

- Random root r
- Random seeds $S_A \sim D_A$
- Random active subgraph *g*

Output: $R = (V_R, E_R, S_{R,A})$

- *V_R*: nodes that might influence *r* in *g*
- E_R : all shortest paths from V_R to r in g
- $S_{R,A} = S_A \cap V_R$: seeds of source A in R.



$$g = G \setminus \{e_1, e_6, e_7\}, S_A = \{4, 9\}$$

$$R = (V_R, E_R, S_{R,A})$$

TCIM: Estimating the Expected Influence

- <u>**R**everse</u> <u>**A**</u>ccessible <u>**P**</u>ointed <u>**G**</u>raph *R*
- Seed set S_B
- Specific competitive propagation model

"Score" of S_B in R: Pr[S_B influences the root of R]



TCIM: High Level Ideas



TCIM: Main Results

(Theorem 4) Two-phase Competitive Influence Maximization

Practical performance guarantee

- $\sigma(S_B | \mathcal{D}_A) \ge (1 1/e \epsilon) \cdot \sigma(S_B^{OPT} | \mathcal{D}_A)$, with probability at least $1 n^{-\ell}$
- the best approximate ratio one could obtain in polynomial time

Practical efficiency

- $\circ O((c(\ell+k)(m+n) \log n)/\epsilon^2)$
- $^{\circ}\,$ the value of c is related to the specific GCIC model

Application: A Special Case of the GCIC model

Distance-based Model (Carnes et al. ICEC'2007)

• Given S_A , S_B and a set of active edges E_a .

• Probability that source B influences node *u*:

 $\frac{\# \text{ of } u' \text{ s nearest seeds of source } B}{\# \text{ of } u' \text{ s nearest seeds of both sources}}$

• TCIM Complexity:
$$O((k(\ell + k)(m + n) \log n)/\epsilon^2)$$

 \downarrow
 $c = O(k)$



 $S_A = \{3\}, S_B = \{4, 5\}$

Experiments

Comparison among the TCIM framework and previous methods

• Dataset: a Facebook-like social networks (1,899 nodes and 20,296 directed edges)

• Baselines:

- CELF (Leskovec et al. ICDM'07): a greedy method
- CELF++ (Goyal et al. WWW'11): a greedy method
- DegreeDiscount (Chen et al. KDD'09): a heuristic method
- Settings:
 - For each edge e_{uv} ∈ E, we set $p_{uv} = 1/d_v^-$ (IC-Weighted Cascade model).
 - \circ We select 50 nodes using single source influence maximization method for source A.

Comparison among the TCIM framework and previously methods



The influence spread of S_B returned by TCIM, CELF and CELF++ are comparable.

Figure 3: Results on the *Facebook-like network*: Influence versus k under three propagation models. $(|S_A| = 50, \ell = 1)$



Figure 4: Results on the *Facebook-like network*: Running time versus k under three propagation models. $(|S_A| = 50, \ell = 1)$

Experimental Results

Results on larger datasets

- The *NetHEPT* collaboration network (15,233 nodes and 58,891 undirected edges)
- The *Epinion* social network (508,837 directed relationships among 75,879 users)



Figure 7: Results on large datasets: Running time versus ϵ under three propagation models. ($|S_A| = 50, k = 50, \ell = 1$)

Remarks:

- 1. When $\epsilon = 0.5$, TCIM finishes within 7 seconds for the *NetHEPT* dataset and finishes within 23 seconds for the *Epinion* dataset.
- 2. If we do not require a very tight approximation ratio, we could choose a larger ϵ .

Experimental Results

Results on larger datasets

- The *NetHEPT* collaboration network (15,233 nodes and 58,891 undirected edges)
- The Epinion social network (508,837 directed relationships among 75,879 users)



Figure 6: Results on large datasets: Running time versus k under three propagation models. ($|S_A| = 50, \epsilon = 0.1, \ell = 1$)

Remarks:

1. With the increase of k, the running time of TCIM tends to drop first, because the number of RAPG instances needed decreases.

2. TCIM is especially efficient for large k.

Experiments: TCIM with partial information

		Competitor's							
strategy influence given explicit S_A selected by different methods ($ S_A = 50$)									
		COICM				Wave propagation model			
dataset	estimated \mathcal{D}_A/S_A	greedy	degree	centrality	average	greedy	degree	centrality	average
NetHEPT	mixed method √ greedy ★ degree ★ centrality	$599.82 \\ 658.38 \\ 400.18 \\ 233.14$	$\begin{array}{c} 632.23 \\ 515.72 \\ 702.93 \\ 478.74 \end{array}$	$\begin{array}{c} 657.49 \\ 519.50 \\ 622.15 \\ 763.43 \end{array}$	629.85 564.53 575.09 491.77	$586.58 \\ 644.53 \\ 372.58 \\ 201.72$	$\begin{array}{c} 624.41 \\ 525.70 \\ 693.95 \\ 462.66 \end{array}$	650.39 515.37 613.98 752.97	$620.46 \\ 561.87 \\ 560.17 \\ 472.45$
Epinion	mixed method greedy degree centrality	$2781.71 \\ 4440.93 \\ 3130.99 \\ 224.93$	$\begin{array}{r} 4603.63\\ 3958.87\\ 5473.33\\ 2809.74\end{array}$	$10683.26 \\ 6372.13 \\ 7283.28 \\ 12078.70$	6022.87 4923.98 5295.87 5037.79	$\begin{array}{c} 2773.17\\ 4265.87\\ 2983.56\\ 204.01 \end{array}$	$\begin{array}{r} 4494.80\\ 3813.06\\ 5299.18\\ 2721.87\end{array}$	$\begin{array}{c} 10517.00\\ 6377.30\\ 7258.24\\ 12075.78\end{array}$	$5928.32 \\4818.74 \\5180.33 \\5000.55$

Table 1: Expected influence of seeds S_B returned by the TCIM framework given the "*mixed method distribution*" (mixed method) as seed distribution for source A or given the guess of explicit seeds of A. Seeds "greedy" for source A is the set of nodes selected by single source influence maximization algorithm. The set "degree" for source A (resp. "centrality") denotes the top 50 nodes ranked by (out)degree (resp. closeness centrality). (k = 50, $\epsilon = 0.1$, $\ell = 1$)

Conclusion

➤General problem formulation

- General Competitive Independent Cascade (GCIC) model
- Competitive Influence Maximization problem with Partial information (CIMP)

General Two-phase Competitive Influence Maximization (TCIM) framework

- It solves the CIMP problem under the GCIC model.
- With probability at least $1 n^{-\ell}$, it guarantees a $(1 1/e \epsilon)$ -approximate solution.
- It runs in $O((c(\ell + k)(n + m) \log n)/\epsilon^2)$ expected time, where *c* depends on the specific propagation model.

> We conduct extensive experiments using real datasets. For example,

• When S_A is given explicitly, we achieve up to **four orders of magnitude speedup** as compared to previous algorithms with the same quality guarantee.

Thank you!



香港中文大學 The Chinese University of Hong Kong