Influence Maximization at Community Level: A New Challenge with Non-submodularity

Lan N. Nguyen
CISE Department
University of Florida
Gainesville, FL, 32611
Email: {lan.nguyen}@ufl.edu

Kunxiao Zhou Computer School Dongguan University of Technology Dongguan, China Email: zhoukx@dgut.edu.cn My T. Thai CISE Department University of Florida Gainesville, FL, 32611 Email: mythai@cise.ufl.edu

Abstract—Motivated by various settings, we study a new Influence Maximization problem at the Community level (IMC) which aims at finding k users to maximize the benefit of influenced communities where a community is influenced iff the number of influenced users belong to this community exceeds its predefined threshold. In general, IMC objective function is not submodular nor supermodular, thereby making it very challenging to apply existing greedy solutions of the classic influence maximization (IM) where submodular function is required. Furthermore, the major challenge in the traditional methods for any related IM problem is the inefficiency in estimating the influence spread. IMC brings this difficulty to a higher level when considering influenced communities instead of influencing each individual user. In this paper, we propose different approximation algorithms for IMC: (1) Using Sandwich approach with a tight submodular function to bound the IMC objective function, (2) Activating the top-k influencing nodes found from network sampling. Furthermore, when the activated thresholds of communities are bounded by a constant, we propose an algorithm with performance guarantee tight to the inapproximability of IMC assuming the exponential time hypothesis. Each algorithm has its own strengths in a trade-off between effectiveness and running time, which are illustrated both in theory and comprehensive experimental

Index Terms—Influence Maximization; Approximation Algorithms; Online Social Networks; Optimization; Viral Marketing

I. INTRODUCTION

Information propagation in Online Social Networks (OSNs) has become a popular research topic recently. Social networks have been shown to be a useful tool for large scale behavioral manipulation and a social network user can be influenced by their friends to adopt certain behaviors. In this context, the Influence Maximization (IM) problem [1] has been studied extensively: Given a network and an integer k, IM asks for k influential users that maximizes the expected influence spread. IM has been widely adopted in viral marketing: companies target top k users by offering them free products or services with the hope to trigger a chain reaction of product adoption. A lot of research [2], [3], [4], [5] (and references therein) has been done to develop a scalable solution with $1-1/e-\epsilon$ performance guarantee for IM.

In this paper, we study a new problem, Influence Maximization at the Community level (IMC). Given a social network and a collection of disjoint set of users, we call each set of users is a community. An *influenced community* is a community in which the number of influenced users in this community exceeds its predefined threshold. A community, if being influenced, will bring back some benefit. The objective

is to identify k seed nodes that maximize the expected benefit of influenced communities.

The IMC problem is motivated by many different applications, where communities are disjoint. As mentioned in the "Collaborative Based" concept on viral marketing: the product may be useful for a single user but is much more useful when used in a group context [6]. In this type of viral marketing, a certain group of users must be influenced to some degree (threshold) so the product or service can be adopted. Moreover, with an increasing in integrating social network to power grid for energy management, IMC could also be applied in the information attack of a power grid through a social network [7]. In this kind of attack, an adversary tries to influence a certain amount of electric users of each community in order to cause inter-area oscillations, causing cascading failure and devastating impact on large geographical areas. In this context, communities represent geographical neighborhoods, hence they are disjoint. Another readily application of IMC is in the context of election where each community represents a state of population.

As IM is a special case of IMC, we ask what will happen if we adopted IM's solutions to IMC. Indeed, most of IM's solutions [1], [2], [3], [4], [8] exploited a trait that IM exhibits submodular behavior, a critical property on how the $(1-1/e-\epsilon)$ -performance guarantee could be acquired. IMC problem, on the other hand, is shown to be neither submodular nor supermodular, making IMC more challenging to devise an efficient algorithm. Moreover, computing the influence spread of a given seed set is a challenging task and IMC elevates this challenge while considering the impact to communities rather than individual independently. Motivated by these observations, the main contributions of this work are:

- We show that IMC is inapproximable within $O(\mathtt{r}^{1/(2(\log\log\mathtt{r})^c)})$ ratio and does not exhibit submodular behaviors, where \mathtt{r} is the number of communities and c>0 is a universal constant independent of \mathtt{r} .
- We present a new sampling technique, modifying from the Reverse Influence Sampling (RIS) [5] so that we take communities as the main subject and randomly sample deterministic instances of the social network in which a community can be influenced. This sampling technique then helps us acquire an estimation on the benefit of influenced communities given an initial seed set.
- We propose three approximation algorithms to solve IMC: (1) Using sandwich approach to provide upper bounded submodular function for the objective function,

- (2) Classifying communities/users by their impact on influencing the other and activating the most influencing users, and (3) When the activation threshold of each community is bounded, we obtain an $O(\mathbf{r}^{1/2})$ -approximation algorithm, which is tight to the inapproximable ratio result assuming the exponential time hypothesis [9].
- We conduct extensive experiments using real world datasets to show the effectiveness of our proposed algorithms and illustrate the above trade-offs.

Related work. Kempe et al [1] first formulated IM and proved its NP-hardness. In addition, computing the expected influence spread for a given seed set is #P-hard. Series of studies have been done to approximate the influence maximization problem [2], [3], [4], [5], [8]. However, different to our work, all of these studies exploited IM's submodularity to obtain the ratio of $1-1/e-\epsilon$ and focused on improving the time complexity.

Some versions of non-submodular related-influenced maximization have been studied recently. Most notably, Lin et al. [10] introduced a k-boosting problem which aims at finding k users to boost so to trigger a maximized "boosted" influence spread. Lu et al. [11] proposed the Comparative Independent Cascade model where the IM under this model is no longer submodular. Both works utilized the Sandwich Approximation strategy, finding a submodular function that bounds the original objective function; the derived results associated with these relaxed functions are then used as the estimated solution to the original non-submodular problem. Furthermore, Ma et al [12] proposed a seed selection strategy using network graphical properties in another non-submodular diffusion model, called the influence barricade model. However, this method provides no theoretical approximation guarantee. All of these models still focused on maximizing influence spread. Different from the above studies, we study how the spread of influence makes an impact to communities given a seed set. The difference in objective function requires us to devise new approaches.

Organization. Section II formally defines the IMC problem and its challenges, including its inapproximability. In Section III, we propose a new approach to estimate the benefit of influenced communities given an initial seed set. Sections IV and V present our three approximation solutions to IMC. The experimental evaluation of our methods is demonstrated in Section VI. Finally, Section VII concludes the paper.

II. MODELS AND PROBLEM DEFINITION

A. Propagation Model and Problem Definition

We abstract a social network using a weighted graph G=(V,E,w) with |V|=n nodes and |E|=m directed edges. Each edge $e\in E$ is associated with a weight $w_e\in [0,1]$ which indicates the probability that u influences v. By convention, $w_e=0$ if $e\not\in E$. Note that for each edge $e=(u,v)\in E$, the notation w_e and w(u,v) are used interchangeably. Technically, we can view G as a generative model for deterministic graphs. A deterministic graph $G=(V,E_G)$ is generated from G by selecting each edge $(u,v)\in E$, independently, with probability w(u,v). We refer to G as a sample graph of G.

We use the fundamental diffusion model namely Independent Cascade (IC) [1]. Assume that we have a set of seed nodes S, the propagation processes happen in discrete rounds. At round 0, all nodes in S are active (influenced) and the others

are inactive. At round $t \ge 0$, when a node u gets activated, initially or by another node, it has a single chance to activate each inactive neighbor v with the successful probability proportional to the edge weight w(u, v). An activated node remains active till the end of the diffusion process. Note that even we only use IC model, similar to [2], [13] the solution can be easily extended to the Linear Threshold model [1].

Given a collection of disjoint sets $Com = \{C_i | 1 \le i \le r, C_i \subset V\}$, r = |Com|. We call each set C_i a community. Let h_i be activation threshold of community C_i . A community exceeds its activation threshold. Denote b_i as a benefit if the community C_i gets influenced. Let c(S) be the expected benefit of influenced communities under a seed set S and a given propagation model. We formally define IMC as follows: **Definition 1.** (IMC) Given a graph G = (V, E, w) an integer $1 \le k \le |V|$, a propagation model and a set of communities Com. The IMC problem asks for a seed set $S \subset V$ of k nodes that maximizes the expected benefit of influenced communities c(S) under the given propagation model.

B. Challenges of the ICM Problem

We now provide key properties and challenges of the IMC problem. Because IM is a special case of IMC, IMC is also NP-hard and computing c(S) given S is #P-hard. Moreover, we show the following inapproximable result, indicating the complexity of IMC.

Theorem 1. Assuming the exponential time hypothesis [9], there is no polynomial time algorithm that approximates IMC to within $O(\mathbf{r}^{1/2(\log\log \mathbf{r})^c})$ factor of the optimum, where c>0 is a universal constant independent of \mathbf{r} .

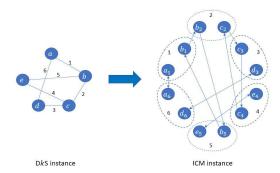


Fig. 1: Conversion from DkS instance to IMC instance. The labels on the edges of DkS graph are corresponding to the labels on communities of IMC instance.

Proof. We consider the following *Densest k-Subgraph* (DkS) problem: given a connected undirected graph G and an integer k, the goal is to find a set S of k-nodes that an induced subgraph of G from S contains maximum number of edges. Let e(S) be the number of edges in induced subgraph of G from S. DkS has been shown to be inapproximable within factor of $O(n^{1/(\log\log n)^c})$ [14] assuming the exponential time hypothesis [9]. To reach the relation between the inapproximability of DkS and IMC, given an instance (G_D, k) of DkS, we construct an instance $(G_I, \operatorname{Com}, k)$ of IMC as follows: For each edge $e = \{a, b\} \in E$, we create a community C_e with activation threshold $\operatorname{h}_e = 2$ and two nodes a_e and b_e in C_e . Node a(b) is called *corresponding node* of $a_e(b_e)$. By that

construction, a multiple nodes in G_I would have the same corresponding node in G_D . Let U_a be a set of nodes having the same corresponding node a. We create edges connecting between nodes in Ua such that Ua is strongly connected and each edge is weighted by 1. Figure 1 shows a simple conversion from DkS instance to ICM instance. Let S_D and S_I are optimal solutions of the DkS and IMC instance. We have following observations:

- We create a solution S_I^\prime for the IMC instance by: for each node $a \in S_D$, pick an arbitrary node in U_a and put into S_I' . Because U_a is strongly connected, all nodes in U_a will be activated. Hence, all communities, corresponding to the edges in the induced subgraph of S_D , are influenced. So we have: $e(S_D) = c(S'_D) \le c(S_I)$
- We create a solution S'_D for the DkS instance by: for each node $u \in S_I$, pick its corresponding node in G_D and put into S'_D . Consider a community C_e which is influenced by S_I , both two node $a_e, b_e \in C_e$ are activated, which means that at least one node in U_a and one node in U_b are in S_I . So a,b will be put into S_D' and edge $e=\{a,b\}$ is in the induced subgraph of S_D' . So we have: $\operatorname{c}(S_I)=$ $e(S_D) \le e(S_D).$

Hence, $e(S_D) = c(S_I)$. Assuming there is a Theree, $e(S_D) = c(S_I)$. Assuming there is a $O(\mathbf{r}^{1/2(\log\log\mathbf{r})^c})$ approximation algorithm K for IMC problem. Denote S_I^{K} as K's solution for IMC instance. We create solution S_D^{K} for DkS instance by picking all corresponding node to S_I^{K} . It is trivial to see that: $e(S_D^{\mathsf{K}}) = \mathsf{c}(S_I^{\mathsf{K}}) \geq O(\mathbf{r}^{-1/2(\log\log\mathbf{r})^c})\mathsf{c}(S_I) \tag{1}$

$$e(S_D^{\mathsf{K}}) = \mathsf{c}(S_I^{\mathsf{K}}) \ge O(\mathsf{r}^{-1/2(\log\log\mathsf{r})^c})\mathsf{c}(S_I) \tag{1}$$

$$\geq O(\mathbf{r}^{-1/2(\log\log\mathbf{r})^c})e(S_D) \tag{2}$$

During the construction of two instances, we have r = $|G_D.E| \leq O(n^2)$ while r = O(n). Replacing into equation (1), we observe that: By using K, we could achieve a solution S_D^{K} for the DkS problem which guarantees:

$$e(S_D^{\mathsf{K}}) \geq O(n^{-1/(\log\log n)^c})e(S_D)$$
 which contradicts to the inapproximability of DkS problem.

Non-submodularity of $c(\cdot)$. In IM, the expected influence of the seed set S is a monotone and submodular function. Thus, a natural greedy algorithm returns a solution with 1-1/eapproximation ratio. However, the objective function c(S) in IMC is neither submodular nor supermodular on the set S of initial activated nodes. Consider each community needs only one node to be influenced, then $c(\cdot)$ exhibits a *submodular* behavior. On the other hand, when several nodes are activated initially, their benefit of influenced communities can be more than the sum of their individual activation effect, which is a supermodular behavior. To illustrate, consider an example in Fig. 2, each edge has weight 0.3 and each community has the activation threshold 2. Therefore, we have $c(\emptyset) =$ $0, c(\{a\}) = 0.327, c(\{b\}) = 0.39, c(\{a,b\}) = 1.09$. Hence, $c(\{b\}) - c(\emptyset) < c(\{a,b\}) - c(\{a\})$, which means $c(\cdot)$ does not exhibit submodular behavior. The non-submodularity of $c(\cdot)$ means that the seed set returned by the greedy algorithm may not have the (1-1/e)-approximation ratio.

Differences to IM. Two notable state-of-the-art influence maximization frameworks are Influence Maximization via Martigale (IMM) [4] and SSA/DSSA [2]. Both of them are based on the sampling method called Reverse Influence Sampling

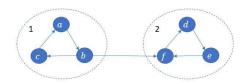


Fig. 2: Non-submodular example

(RIS). Briefly, RIS captures the influence landscape of G by: for a node $v \in V$, a random set R of nodes is generated such that for any seed set S, the probability that $R \cap S \neq \emptyset$ equals the probability that v can be influenced by S. If v is selected uniform randomly from V, the expected influence of any seed set $S \subseteq V$ equals $n \cdot \mathbb{E}[\mathbf{1}_{R \cap S \neq \emptyset}]$ where expectation is taken over the randomness of R. Then, both SSA and IMM used the greedy algorithm for the max-coverage problem to find a seed set that covers the maximum number of sampling sets.

However, the IMC problem is fundamentally different due to the objective function that considers influenced communities instead of individuals separately. RIS seems to be inapplicable in IMC because we find no relationship between the randomly selected node in RIS with the community it belongs to and how it can contribute to estimate the expected benefit of influenced communities. Although such challenge can be solved by subtly modifying RIS, the most challenging task is that IMC's objective function is non-submodular and inapproximable within $O(r^{1/(2(\log\log r)^c)})$. Hence, using greedy solutions for the max-coverage problem given a collection of samples will not provide 1-1/e approximation ratio. Therefore, we need to devise efficient approximation algorithms based on the new sampling technique.

III. BENEFIT ESTIMATION

In this section, we propose a sampling method, called Reverse Influenceable Community (RIC), to estimate the expected benefit of influenced communities given a seed set S.

We denote $\rho: \mathsf{Com} \to \mathbb{R}^{\geq}$ as a probability distribution among Com, where $\rho(C_i) = \frac{b_i}{b}$ and $b = \sum_{j=1}^{r} b_i$. Given a sample graph \mathcal{G} of G, a community C is called being *touched* by $u \in V$ (or u touches C) in \mathcal{G} if there exist a path in \mathcal{G} that connects u to any nodes in C. The RIC sample is defined as follows:

Definition 2. (RIC) Given G = (V, E, w), Com, a random RIC sample g is generated by: 1) selecting a random community C_q given probability distribution ρ ; 2) generating a sample graph \mathcal{G}_q from G and 3) returning a set of nodes that touch C_g in G_g . C_g is called the "source" community of g.

Fig. 3 shows an example of RIC-sample g. The "source" community C_g is marked by the dotted circle and contains node $\{v_1, v_2, v_3, v_4\}$. $h_g = 3$ is an activation threshold of C_g . In this example, we need at least 3 activated users among $\{v_1, v_2, v_3, v_4\}$ to make g influenced. As definition 2, g contains all nodes (even outside C_g) that connect to the nodes in C_q , including v_5, v_6 and v_7 . For simplicity, if C_q is touches by u in \mathcal{G}_g , we also call u touches g or g is touched by u.

Alg. 1 presents in detail how a RIC sample is generated. For each $u \in C_q$, $R_q(u)$ is a reachable set of node u, which is defined as a set of nodes $v \in V$ that can reach u in \mathcal{G}_g (there exists a path from v to u). First, C_g is randomly selected given probability distribution ρ . Next, we mark each edge e in Gwith one of three states: (1) \perp - e has potential to be in \mathcal{G}_q ,

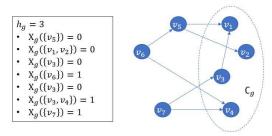


Fig. 3: Example of RIC-sample. Additionally, we also store $R_q(\cdot)$ values of v_1, v_2, v_3, v_4 , i.e $R_q(v_1) = \{v_1, v_3, v_5, v_6, v_7\}$ (2) y - e is in \mathcal{G}_g , and (3) n - e is not in \mathcal{G}_g . Each edge's state is set to be \perp initially and will be changed to y or n when we first process it. This state information is stored by st[·] array. The algorithm then works in backward Breadth-First Search (BFS) manner as follows: First, we put all nodes in C_q into Q - a queue storing nodes we are going to visit. While Q is not empty, we take a node u out of Q and mark it as "visited" (ck[u] = true if u is "visited", false otherwise). Next, we consider all nodes v that u has incoming edge (e = (v, u))with. If $st[e] = \perp$, edge e is never processed before, so we randomly generate edge e = (v, u) with probability w_e . If edge e = (v, u) exists in \mathcal{G}_g and v does not exist in Q as well as was not visited before, we put v into Q. This process is repeated until Q is empty and we will have a deterministic sample graph \mathcal{G}_g that contains paths by which the nodes in C_g could be activated. Finally, for all $u \in C_g$, we determine $R_q(u)$ by using Depth-First Search with graph \mathcal{G}_q .

Estimating c(S). We now present how we use the RIC-sampling to estimate c(·). Let g be a given RIC sample with the source community C_g and the activation threshold h_g . Given a seed set S, we say that g is influenced by S if S could reach at least h_g nodes in C_g . In Fig. 3, g is influenced by $\{v_5, v_6\}$ or $\{v_7\}$ but not by $\{v_1\}$ or $\{v_1, v_4\}$. Define $X_g(S): 2^V \to \{0, 1\}$ as an indicator function whether g is influenced by S.

Algorithm 1 Generate a random RIC graph

Input A social graph G, set of communities Com Output Source community C_g and its activation threshold h_g , reachable set $R_g(u) \quad \forall u \in C_g$

```
1: Select source community C_q randomly
 2: Create a graph \mathcal{G}_g with all nodes in G and no edge
 3: \operatorname{st}[e_{uv}] = \perp \forall (u,v) \in E
 4: Initiate empty queue Q
 5: Put all nodes v \in C_g into Q
 6: while Q is not empty do
         u \leftarrow Q.\text{dequeue}()
 7:
         \mathtt{ck}[u] \leftarrow \mathtt{true}
 8:
         for each in-comming edge e_{vu} of u do
 9:
              if st[e_{vu}] = \perp then
10:
11:
                  st[e_{vu}] = y with prob w_e: n otherwise
                  if st[e_{vu}] = y then create edge e_{vu} in \mathcal{G}_q
12:
              if st[e_{vu}] = y \& !ck[v] \& !Q.contain(v) then
13:
                  Q.enqueue(v)
14:
15: for each node u \in C_q do
         Run DFS to find R_a(u)
Return C_g, h_g, R_g(u) \quad \forall u \in C_g
```

Lemma 1. For any $S \subseteq V$, we have $c(S) = b \cdot \mathbb{E}[X_q(S)]$

where the expectation is taken over the randomness of g.

$$\begin{split} \textit{Proof.} & \text{ Because} \quad \mathbf{X}_g(S) \quad \text{is an indicator variable, then} \\ & \mathbb{E}[\mathbf{X}_g(S)] = \Pr[\mathbf{X}_g(S) = 1]. \text{ We have:} \\ & \Pr[\mathbf{X}_g(S) = 1] = \sum_{\mathbf{C}_i \in \mathsf{Com}} \Pr[\mathbf{C}_i = \mathbf{C}_g] \cdot \Pr[\mathbf{C}_i \text{ is influenced by } S] \\ & = \sum_{\mathbf{C}_i \in \mathsf{Com}} \frac{b_i}{\mathbf{b}} \Pr[\mathbf{C}_i \text{ is influenced by } S] = \frac{1}{\mathbf{b}} \mathbf{c}(S) \end{split}$$

Let \mathcal{R} be a set of independent random RIC samples, define:

$$\hat{\mathbf{c}}_{\mathcal{R}}(S) = \frac{\mathbf{b}}{|\mathcal{R}|} \cdot \sum_{g \in \mathcal{R}} \mathbf{X}_g(S) \quad \forall S \subseteq V \tag{3}$$
 Intuitively, $\hat{\mathbf{c}}_{\mathcal{R}}(S)$ closely estimates $\mathbf{c}(S)$ for any $S \subseteq V$ if

Intuitively, $\hat{c}_{\mathcal{R}}(S)$ closely estimates c(S) for any $S\subseteq V$ if $|\mathcal{R}|$ is sufficiently large. Therefore, IMC could be solved by generating a large collection \mathcal{R} of RIC samples, then finding a seed set S that influences the maximum number of RIC graphs in \mathcal{R} . We define a problem of finding S given a collection \mathcal{R} of RIC samples - the MAXR problem, which is essentially different with the max-coverage problem solved in IM.

Definition 3. (MAXR) Given a set \mathcal{R} of RIC samples, find a seed set $S \subset V$ of k nodes that maximizes the number of influenced RIC samples in \mathcal{R} . The RIC sample g is called influenced by S if S can reach at least \mathbf{h}_g nodes in \mathbf{C}_g .

By this definition, we observe that the objective function of MAXR is equivalent to $\hat{c}_{\mathcal{R}}(\cdot)$. So from now on, we will also call $\hat{c}_{\mathcal{R}}(\cdot)$ the objective function of MAXR.

Clearly, MAXR is NP-hard as max-coverage is a special case. Moreover, the additional challenge of solving MAXR is presented by the following lemma, indicating that a natural greedy algorithm cannot solve MAXR either.

Lemma 2. $\hat{c}_{\mathcal{R}}(\cdot)$ is non-submodular.

which completes the proof.

Proof. We consider a simple instance where \mathcal{R} contains only one RIC graph g. The source community C_g has only 2 nodes u,v and $R_g(u)=\{u\},R_g(v)=\{v\}$. We have $\hat{c}_{\mathcal{R}}(\emptyset)=0,\hat{c}_{\mathcal{R}}(\{u\})=0,\hat{c}_{\mathcal{R}}(\{v\})=0$ and $\hat{c}_{\mathcal{R}}(\{u,v\})=1$. Therefore, $\hat{c}_{\mathcal{R}}(\{u\})-\hat{c}_{\mathcal{R}}(\emptyset)<\hat{c}_{\mathcal{R}}(\{u,v\})-\hat{c}_{\mathcal{R}}(\{v\})$

Overall, two critical pieces needed to solve IMC are:

- Solving the MAXR problem to find a seed set S that influences the maximum number of RIC samples.
- Generating a sufficiently large number of RIC samples such that the returned set S from the MAXR problem guarantees the bounded error.

IV. ALGORITHMS FOR THE MAXR PROBLEM

In this section, we tackle the first piece by devising three approximation algorithms to solve the MAXR problem.

A. Upper Bound Greedy (UBG)

The first algorithm is called Upper Bounded Greedy. Overall, UBG uses a concept of Sandwich Approximation [11] by bounding $\hat{c}_{\mathcal{R}}(\cdot)$ with a submodular function. The most challenging task is to devise a submodular function that is tight to our objective function. In this solution, we propose submodular functions $\nu(\cdot), \nu_{\mathcal{R}}(\cdot)$ to bound $c(\cdot), \hat{c}_{\mathcal{R}}(\cdot)$ respectively and prove the tightness by observing that those functions overlap our objective function in a special case. The Sandwich Approximation is defined as follows:

Algorithm 2 UBG algorithm

Input \mathcal{R}, k Output S

- 1: $S_{\nu} \leftarrow$ greedy selection with objective function $\nu_{\mathcal{R}}(\cdot)$
- 2: $S_c \leftarrow$ greedy selection with objective function $\hat{c}_{\mathcal{R}}(\cdot)$
- 3: $S \leftarrow \arg \max_{S_{\nu}, S_{\mathfrak{c}}} \hat{\mathfrak{c}}_{\mathcal{R}}(S)$

Return S

Theorem 2. (Sandwich Approximation [11]) Let $\sigma: 2^V \to \mathbf{R}$ be a non-submodular. Let μ and ν be submodular and defined on the same ground set V such that $\mu(S) \leq \sigma(S) \leq \nu(S)$ for all $S \subseteq V$. That is $\mu(\cdot)(\nu(\cdot))$ is a lower (resp., upper) bound on σ everywhere. Consider the problem of maximizing σ subject to a cardinality constraint k: run the greedy algorithm on all three functions. It produces an approximate solution for μ and ν . Let $S_{\sigma}, S_{\mu}, S_{\nu}$ be the solution obtained from σ, μ, ν respectively. Then, select the final solution to be

$$S_{\text{sand}} = \underset{S \in \{S_{\sigma}, S_{u}, S_{v}\}}{\arg \max} \sigma(S) \tag{4}$$

Then we have:
$$\sigma(S_{\texttt{sand}}) \geq \max \left\{ \frac{\sigma(S_{\nu})}{\nu(S_{\nu})}, \frac{\mu(S_{\sigma}^*)}{\sigma(S_{\sigma}^*)} \right\} (1 - 1/e) \sigma(S_{\sigma}^*) \quad (5)$$
 Therefore, we have to find at least one of $\mu(\cdot)$ or $\nu(\cdot)$

that should be close to our objective function. In this work, we derive a submodular upper bound $\nu_{\mathcal{R}}(\cdot)$ of $\hat{c}_{\mathcal{R}}(\cdot)$. $\nu_{\mathcal{R}}(\cdot)$ is considerably closer to $\hat{c}_{\mathcal{R}}(\cdot)$ than any other submodular bounded functions we have tested. Define:

$$\nu(S) = \mathbf{b} \cdot \mathbb{E}\left[\min\left(\frac{|\mathbf{I}_g(S)|}{\mathbf{h}_a}, 1\right)\right] \tag{6}$$

$$\nu_{\mathcal{R}}(S) = \frac{\mathbf{b}}{|\mathcal{R}|} \sum_{g \in \mathcal{R}} \min\left(\frac{|\mathbf{I}_g(S)|}{\mathbf{h}_g}, 1\right) \tag{7}$$
 where $\mathbf{I}_g(S)$ is the a set of nodes in \mathbf{C}_g that can be connected

by S in RIC sample g.

 $\nu(S)$ presents the expected fractional benefit from being influenced of all communities while $\nu_{\mathcal{R}}(S)$ provides an estimation on $\nu(S)$ from the sample set \mathcal{R} . Recalling that RIC sample g is called influenced if $|I_g(S)| \ge h_g$. Therefore, to display the fractional value of being influenced, we take min value between $|I_q(S)|/h_q$ and 1.

Lemma 3. For any set $S \subseteq V$, $c(S) \leq \nu(S)$. Also, for a given collection \mathcal{R} of RIC samples, $\hat{c}_{\mathcal{R}}(S) \leq \nu_{\mathcal{R}}(S)$.

Proof. First, let compare $\hat{c}_{\mathcal{R}}(\cdot)$ and $\nu_{\mathcal{R}}(\cdot)$. Given $g \in \mathcal{R}$, if $X_g(S) = 1$ then $|I_g(S)| \ge h_g$ and $\min(|I_g(S)|/h_g, 1) = 1$. Otherwise $X_g(S) = 0$, but $\min (|I_g(S)|/h_g, 1) \ge 0$. Therefore, $\hat{c}_{\mathcal{R}}(S) \leq \nu_{\mathcal{R}}(S)$, $\forall S \subseteq V$.

To prove $\nu(\cdot)$ is an upper bound of $c(\cdot)$, recalling c(S) = $b \cdot \mathbb{E}[X_q(S)]$ where:

 $\mathbb{E}[\mathtt{X}_g(S)] = \mathbb{E}[\mathbf{1}_{|\mathtt{I}_g(S)| \geq \mathtt{h}_g}] = \Pr[|\mathtt{I}_g(S)|/\mathtt{h}_g \geq 1]$ Without loosing generality, we have:

 $\Pr[|\mathtt{I}_q(S)|/\mathtt{h}_q \ge 1] = \Pr[\min(|\mathtt{I}_q(S)|/\mathtt{h}_q, 1) \ge 1]$

Using Markov inequality [15], we have: $\Pr[\min(|\mathtt{I}_g(S)|/\mathtt{h}_g,1)\geq 1]\leq \mathbb{E}[\min(|\mathtt{I}_g(S)|/\mathtt{h}_g,1)]$

So $c(\cdot)$ is upper bounded by $\nu(\cdot)$.

Considering an RIC sample g, it is trivial that $|\mathbf{I}_g(\cdot)|$ is a submodular function, which also means $\nu_{\mathcal{R}}(S) =$ $\frac{\mathbf{b}}{|\mathcal{R}|}\sum_{g\in\mathcal{R}}\min\left(\frac{|\mathbf{I}_g(S)|}{\mathbf{h}_g},1\right)$ is a *submodular* function. Furthermore, $\nu_{\mathcal{R}}(\cdot)$ overlaps $\hat{c}_{\mathcal{R}}(\cdot)$ when the activation thresholds of communities are bounded by 1. We have the following lemma.

Algorithm 3 MAF algorithm

Input \mathcal{R}, k Output S

- 1: Initiate $S_1, S_2 \leftarrow \emptyset$
- 2: $SC \leftarrow sorted \ list \ of \ Com \ in \ order \ of \ their \ appearance \ in \ \mathcal{R}$
- while SC is not empty do
- $C \leftarrow \text{take out } 1^{\text{st}} \text{ community of SC};$
- 5: $X \leftarrow pick h nodes in C$
- if $|S_1 \cup X| \leq k$ then $S_1 = S_1 \cup X$
- 7: $S_2 \leftarrow k$ nodes that appear the most in \mathcal{R}
- 8: $S = \arg\max_{S' \in \{S_1, S_2\}} \hat{\mathsf{c}}_{\mathcal{R}}(S')$

Lemma 4. For any set $S \subseteq V$, $c(S) = \nu(S)$ and $\hat{c}_{\mathcal{R}}(S) =$ $\nu_{\mathcal{R}}(S)$ if $\mathbf{h}_q = 1$ for all $g \in \mathcal{R}$.

Proof. If
$$h_g = 1$$
 for all g , then
$$\min(\frac{|\mathtt{I}_g(S)|}{h_g}, 1) = \min(|\mathtt{I}_g(S)|, 1)$$
$$= 1 \text{ is a report } = \mathtt{X}_g(S)$$

the lemma follows

Lemma 4 also implies that $c(\cdot)$ and $\hat{c}_{\mathcal{R}}(\cdot)$ are submodular if the activation thresholds are bounded by 1.

Our UBG algorithm to MAXR is presented in Alg. 2. Denote S_{ν} as a solution obtained from the greedy algorithm with objective function $\nu_{\mathcal{R}}(\cdot).$ By Theorem 2, UBG guarantees $\frac{\hat{\mathbf{c}}_{\mathcal{R}}(S_{\nu})}{\nu_{\mathcal{R}}(S_{\nu})}(1-1/e)$ approximation ratio to the MAXR problem.

B. Most Appearance First (MAF)

The second algorithm we present is called Most-Appearance-First (MAF). Accordingly, the idea of MAF is that we compute the frequency of appearance of communities or nodes in \mathcal{R} , then try to activate the most influential one. First, MAF considers the communities that appear the most in \mathcal{R} . The algorithm starts by determining frequency each community in Com appears in R then sorts these communities in descending order of their frequency. Let SC be a sorted list of communities. MAF then sequentially takes the community C with highest frequency out of SC. Let h be the activation threshold of C. MAF then randomly picks h nodes in C and puts into S_1 . This process then repeats until $|S_1| = k$. Next, MAF considers nodes that appear the most in \mathcal{R} . A node v is considered to appear in RIC sample g if v is either in C_g or there exists $u \in C_g$ that $v \in R_q(u)$. MAF then picks k nodes that appear the most in \mathcal{R} and put into S_2 . Finally, MAF returns a solution $S \in \{S_1, S_2\}$ that influences the most RIC samples in \mathcal{R} . The full MAF is presented by Alg. 3.

Theorem 3. MAF returns a $\frac{1}{r} \lfloor \frac{k}{h} \rfloor$ approximate result to the MAXR problem, where $h = \max_{i=1}^{r} h_i$

Proof. Denote
$$S_{\text{OPT}}$$
 an optimal solution to MAXR, we have:
$$\hat{\mathbf{c}}_{\mathcal{R}}(S_1) = \frac{\mathbf{b}}{|\mathcal{R}|} \sum_{g \in \mathcal{R}} \mathbf{X}_g(S_k) \quad \geq \frac{\mathbf{b}}{|\mathcal{R}|} \cdot \frac{1}{\mathbf{r}} \cdot \left\lfloor \frac{k}{h} \right\rfloor \cdot |\mathcal{R}| \qquad (9)$$

$$\geq \frac{1}{\mathtt{r}} \left \lfloor \frac{k}{h} \right \rfloor \cdot \frac{\mathtt{b}}{|\mathcal{R}|} \sum_{g \in \mathcal{R}} \mathtt{X}_g(S_{\mathtt{OPT}}) \geq \frac{1}{\mathtt{r}} \left \lfloor \frac{k}{h} \right \rfloor \hat{\mathtt{c}}_{\mathcal{R}}(S_{\mathtt{OPT}})$$

The inequality (9) comes from the observation that MAF chooses to influence communities appearing the most in Rand the "budget" to influence each community is at most $h = \max_{i=1}^{r} \mathbf{h}_i$. Therefore, there are at least $|\mathcal{R}| \frac{1}{r} \lfloor \frac{k}{h} \rfloor$ RIC samples which are influenced by the seed set S_1 .

 S_2 , on the other hand, does not provide any approximation guarantee. We consider an example as follows: Assume \mathcal{R} contains 6 RIC samples $\{g_1,g_2,..g_6\}$ in which C_{g_i} s are disjoint sets, $|C_{g_i}|=3$ and $h_{g_i}=2$ for all $i\in[1,6]$. Given a node $u\not\in C_{g_i}$ for $i\in[1,6]$, there exists an edge in g_i that connects u with a node in C_{g_i} for $i\in[1,6]$, there exists an edge in g_j that connects v with a node in C_{g_i} for $i\in[1,6]$, there exists an edge in g_j that connects v with a node in C_{g_j} for $j\in\{4,5,6\}$. Other than these, there is no other edges in g_i for $i\in[1,6]$. Let k=2. In this scenario, u,v are nodes that appear the most in \mathcal{R} , thus $S_2=\{u,v\}$. However $\hat{c}_{\mathcal{R}}(S_2)=0$ since each community C_{g_i} required at least 2 nodes to be influenced while S_2 only guarantees to influence one node in each community. The optimal solution is $\hat{c}_{\mathcal{R}}(S_{\mathrm{OPT}})=1$. Therefore, S_2 does not provide any approximation guarantee to MAXR. However, S_2 actually performs well in experiments. By selecting the best among S_1, S_2 , we can bound the performance guarantee of MAF and integrate MAF into IMC framework, which will be presented in next section.

C. Algorithm for Bounded Activation Threshold

Having considered two approximation algorithms to MAXR, we now propose solutions with a tight approximation ratio for a special case where the activation threshold of communities is bounded by 2, which means a community needs at least 2 or 1 activated users to be influenced. Then we extend our solution to a case that the input network has bounded activation thresholds by a constant $d \geq 2$.

Denote $I_g(S)$ as a set of nodes in C_g that can be connected by S in the RIC sample g, $I_g(S) = \{u \mid u \in C_g, R_g(u) \cap S \neq \emptyset\}$. Next, let $\mathcal{G}_{\mathcal{R}}(u)$ be a set of RIC samples in \mathcal{R} that u can touch - $\mathcal{G}_{\mathcal{R}}(u) = \{g \mid g \in \mathcal{R} : \exists v \in C_g, u \in R_g(v)\}$; and $\mathcal{D}_{\mathcal{R}}(S,u)$ is a set of RIC samples in $\mathcal{G}_{\mathcal{R}}(u)$ that S can influence, $\mathcal{D}_{\mathcal{R}}(S,u) = \{g \mid g \in \mathcal{G}_{\mathcal{R}}(u), I_g(S) \geq h_g\}$. Before going further, we have the following lemma.

Lemma 5. Given a seed set
$$S \subseteq V$$
, we have
$$\max_{u \in S} |\mathcal{D}_{\mathcal{R}}(S, u)| \leq \sum_{g \in \mathcal{R}} \mathtt{X}_g(S) \leq \sum_{u \in S} |\mathcal{D}_{\mathcal{R}}(S, u)|$$

Proof. The first inequality is trivial because $\mathcal{D}_{\mathcal{R}}(S,u)$ is a subset of the RIC samples set that S can influence. The second inequality comes from observing that: a sample g, if being influenced by S, would appear on at least one set $\mathcal{D}_{\mathcal{R}}(S,u)$ for all $u \in S$.

Lemma 5 gives us an observation that: the number of influenced RIC samples by S is lower bounded by the size of $\mathcal{D}_{\mathcal{R}}(S,u)$ with any $u\in S$ separately but being upper bounded by sum of $|\mathcal{D}_{\mathcal{R}}(S,u)| \ \forall u\in S$. This motivates us to propose the BT algorithm, which is presented in Alg. 4.

The idea of BT is that: for each node $u \in V$, we find a seed set $\mathrm{K}(u)$ that maximizes the number of influenced RIC samples in $\mathcal{G}_{\mathcal{R}}(u)$. To find $\mathrm{K}(u)$, we first add u in $\mathrm{K}(u)$. Because every RIC sample $g \in \mathcal{G}_{\mathcal{R}}(u)$ is touched by u and $\mathrm{h}_g \leq 2$, to make g be influenced by $\mathrm{K}(u)$, we only need at most one more node in C_g to be connected by $\mathrm{K}(u)$ beside the one that is already connected by u. That brings us back to the case when the activation threshold is bounded by 1 and a natural greedy algorithm can return (1-1/e)-approximation result for $\mathrm{K}(u)$.

Algorithm 4 BT algorithm

Input \mathcal{R}, k Output S

```
1: for each node u \in V do
           G \leftarrow \text{Copy of } \mathcal{G}_{\mathcal{R}}(u)
 2:
           for each RIC sample g \in G do
 3:
                for each node v \in C_g do
 4:
                      if u \in R_g(v) then
 5:
                           Remove v out of g
 6:
                           h_a = h_a - 1
 7:
 8:
           \mathcal{T} \leftarrow \text{greedily select } k-1 \text{ nodes with collection } G
          K(u) \leftarrow \{u\} \cup \mathcal{T}
10: S = \arg \max_{\mathbb{K}(u); u \in V} |\mathcal{D}_{\mathcal{R}}(\mathbb{K}(u), u)|
```

Return S

BT returns a seed set S which is K(u) that has the maximum $|\mathcal{D}_{\mathcal{R}}(K(u), u)|$. The approximate guarantee of BT is obtained by the following theorem.

Theorem 4. BT provides a $\frac{1}{k}(1-\frac{1}{e})$ approximation guarantee to the MAXR problem if the activation threshold of each community is at most 2.

Proof. Denote $S_{\mathtt{OPT}}$ as an optimal solution for the MAXR problem. For any $u \in S_{\mathtt{OPT}}$, the natural greedy algorithm guarantees

$$|\mathcal{D}_{\mathcal{R}}(\mathtt{K}(u), u)| \ge (1 - 1/e) \max_{T \subset V; |T| = k - 1} |\mathcal{D}_{\mathcal{R}}(T \cup \{u\}, u)|$$

$$> (1 - 1/e)|\mathcal{D}_{\mathcal{R}}(S_{\mathtt{DRT}}, u)|$$

 $\geq (1-1/e)|\mathcal{D}_{\mathcal{R}}(S_{\mathsf{OPT}},u)|$ Denote $\mathtt{v} = \arg\max_{u \in S_{\mathsf{OPT}}} |\mathcal{D}_{\mathcal{R}}(S_{\mathsf{OPT}},u)|$ and $\mathtt{u} = \arg\max_{u \in V} |\mathcal{D}_{\mathcal{R}}(\mathtt{K}(u),u)|.$ According to BT, $S = \mathtt{K}(\mathtt{u}).$ Therefore, we have:

$$\begin{split} \hat{\mathbf{c}}_{\mathcal{R}}(S_{\mathsf{OPT}}) &= \frac{\mathbf{r}}{|\mathcal{R}|} \sum_{g \in \mathcal{R}} \mathbf{X}_g(S_{\mathsf{OPT}}) \leq \frac{\mathbf{r}}{|\mathcal{R}|} \sum_{u \in S_{\mathsf{OPT}}} |\mathcal{D}_{\mathcal{R}}(S_{\mathsf{OPT}}, u)| \\ &\leq \frac{\mathbf{r}}{|\mathcal{R}|} k \cdot |\mathcal{D}_{\mathcal{R}}(S_{\mathsf{OPT}}, \mathbf{v})| \leq \frac{\mathbf{r}}{|\mathcal{R}|} \frac{k}{1 - 1/e} |\mathcal{D}_{\mathcal{R}}(\mathbf{K}(\mathbf{v}), \mathbf{v})| \\ &\leq \frac{k}{1 - 1/e} \frac{\mathbf{r}}{|\mathcal{R}|} |\mathcal{D}_{\mathcal{R}}(\mathbf{K}(\mathbf{u}), \mathbf{u})| \leq \frac{k}{1 - 1/e} \hat{\mathbf{c}}_{\mathcal{R}}(S) \\ \text{which completes the proof.} \end{split}$$

In Alg. 4, BT starts by iterating through every node $u \in V$. For each node, BT creates a sample set G which contains all RIC samples that u can touch. For each sample $g \in G$, BT removes nodes in C_g that u can connect to and decreases h_g . After this step, each sample $g \in G$ will have an activation threshold at most 1. Then the greedy algorithm is applied to find k-1 nodes that maximizes the number of influenced samples in G. K(u) will contain all these nodes and u.

Extending to d-bounded thresholds: An interesting observation is that we can extend BT to solve IMC with activation thresholds are bounded by a constant d. Denote $\mathrm{BT}^{(d)}$ as an algorithm to solve IMC with bounded threshold d.

Assume we have $\mathrm{BT}^{(d-1)}$ in hand with α -approximation

Assume we have $\mathrm{BT}^{(d-1)}$ in hand with α -approximation guarantee. We devise $\mathrm{BT}^{(d)}$ by modifying from Alg. 4 that the set \mathcal{T} (line 8) is selected by running $\mathrm{BT}^{(d-1)}$ to find k-1 nodes that maximize influenced RIC samples in $\mathcal{G}_{\mathcal{R}}(u)$. Using similar proof as Theorem 4, $\mathrm{BT}^{(d)}$ provides $\frac{\alpha}{k}$ approximation guarantee to IMC with d-bounded threshold. As the approximation ratio of $\mathrm{BT}^{(2)}$ is $\frac{1-1/e}{k}$, using induction, we bound the approximation guarantee of $\mathrm{BT}^{(d)}$ by $\frac{1-1/e}{k^{d-1}}$.

Combining with MAF: In subsection IV-B, we proposed the MAF algorithm which guarantees $\frac{rh}{k}$ approximate ratio. Interestingly, if we combine MAF and BT by running both these two algorithms to get two different results - called S_{MAF} , S_{BT} - then return $S = \arg\max_{S' \in \{S_{\mathtt{MAF}}, S_{\mathtt{BT}}\}} \hat{\mathsf{c}}_{\mathcal{R}}(S')$, we obtain a solution with a new approximation guarantee, which is tight to the inapproximability result in Theorem 1. Let us call this solution MB (MAF and BT).

Theorem 5. MB provides $\Theta(\sqrt{\frac{1-1/e}{r}})$ approximation guarantee to the MAXR problem if the activation threshold of each community is at most 2.

Proof. By definition, we have:

$$\begin{split} \hat{\mathbf{c}}_{\mathcal{R}}^2(S) & \geq \hat{\mathbf{c}}_{\mathcal{R}}(S_{\mathtt{MAF}}) \cdot \hat{\mathbf{c}}_{\mathcal{R}}(S_{\mathtt{BT}}) \\ & \geq \frac{1 - 1/e}{k} \cdot \frac{1}{\mathtt{r}} \cdot \left\lfloor \frac{k}{2} \right\rfloor \cdot \hat{\mathbf{c}}_{\mathcal{R}}^2(S_{\mathtt{OPT}}) \end{split}$$

$$\begin{array}{l} \frac{\lfloor k/2 \rfloor}{k} = \Theta(1) \text{ with sufficient large } k. \text{ Therefore, } \hat{\mathsf{c}}_{\mathcal{R}}(S) \geq \\ \Theta(\sqrt{\frac{1-1/e}{\mathsf{r}}}) \hat{\mathsf{c}}_{\mathcal{R}}(S_{\mathsf{OPT}}), \text{ which completes the proof.} \end{array} \ \Box$$

V. ALGORITHMIC FRAMEWORK TO IMC

In this section, we solve the second required piece: finding the minimum number of RIC samples to provide a boundederror guarantee ϵ , where ϵ is an arbitrary small number. Then, putting together with any algorithm for MAXR described in section IV, we provide an algorithmic framework IMCAF for the IMC problem such that: given an α -approximation algorithm to MAXR, the IMCAF returns an $\alpha(1-\epsilon)$ -approximation guarantee to IMC with probability at least $1 - \delta$, where δ is an arbitrary small number. Accordingly, we have three approximation algorithms to IMC with the following ratios: $\frac{1}{r}\lfloor\frac{k}{h}\rfloor(1-\epsilon) \text{ using MAF, } \Theta(\!\sqrt{\frac{1-1/e}{2r}})(1-\epsilon) \text{ using MB, and } \frac{\operatorname{c}(S_{\nu})}{\nu(S_{\nu})}(1-1/e-\epsilon) \text{ using UBG.}$

A. Minimum size of R

To bound the error between $\hat{c}_{\mathcal{R}}(\cdot)$ and $c(\cdot)$ within an ϵ value, we utilize the following lemma, which is trivially derived from the martingale theory in [16].

Lemma 6. [16] Given a collection R of RIC samples, a seed set S and $\epsilon > 0$, the following inequalities hold,

$$\Pr[\hat{\mathsf{c}}(S) > (1+\epsilon)\mathsf{c}(S)] \le \exp(\frac{-|\mathcal{R}|\epsilon^2}{3\mathsf{b}}\mathsf{c}(S)) \tag{11}$$

$$\Pr[\hat{\mathsf{c}}(S) < (1 - \epsilon)\mathsf{c}(S)] \leq \exp(\frac{-|\mathcal{R}|\epsilon^2}{2\mathsf{b}}\mathsf{c}(S)) \tag{12}$$
 Denote S^* as an optimal solution to IMC. We have the

following observations:

Corollary 1. Given $0 < \epsilon_1, \delta_1 < 1$, with the number of RIC samples satisfies

$$|\mathcal{R}| \geq rac{2\mathsf{b} \ln(1/\delta_1)}{\epsilon_1^2 \mathsf{c}(S^*)}$$

the following condition is guaranteed

$$\Pr[\hat{c}(S^*) \ge (1 - \epsilon_1)c(S^*)] \ge 1 - \delta_1$$
 (13)

This corollary is trivially derived from equation (12).

Corollary 2. Given $0 < \alpha, \epsilon_2, \delta_2 < 1$, with the number of RIC samples satisfies

$$|\mathcal{R}| \ge \frac{3b \ln(\binom{n}{k}/\delta_2)}{\alpha^2 \epsilon_2^2 \mathbf{c}(S^*)} \tag{14}$$

for any seed set S_k of k nodes, the following condition is guaranteed:

$$\Pr[\hat{\mathsf{c}}(S) \ge \mathsf{c}(S) + \alpha \epsilon_2 \mathsf{c}(S^*)] \le \delta_2 \tag{15}$$

Proof. Consider an arbitrary seed set of k nodes S_k , we have: $\Pr[\hat{\mathsf{c}}(S) \ge \mathsf{c}(S) + \alpha \epsilon_2 \mathsf{c}(S^*)]$

$$=\Pr[\hat{\mathsf{c}}(S) \geq (1 + \alpha \epsilon_2 \frac{\mathsf{c}(S^*)}{\mathsf{c}(S)})\mathsf{c}(S)]$$

$$\leq \exp(-\frac{|\mathcal{R}|}{3\mathsf{b}}\frac{\alpha^2\epsilon_2^2\mathsf{c}^2(S^*)}{\mathsf{c}^2(S)}\mathsf{c}(S)) \leq \exp(-\frac{|\mathcal{R}|\alpha^2\epsilon_2^2}{3\mathsf{b}}\mathsf{c}(S^*))$$
 Using the union bound theory, to let condition (15) satisfy for

any seed set S of k nodes, we have:

any seed set
$$S$$
 of k nodes, we have.
$$\Pr[\hat{\mathbf{c}}(S) \geq \mathbf{c}(S) + \alpha \epsilon_2 \mathbf{c}(S^*)] \leq \binom{n}{k} \exp(-\frac{|\mathcal{R}|\alpha^2 \epsilon_2^2}{3\mathbf{b}} \mathbf{c}(S^*))$$
 The equation (15) follows by replacing $\binom{n}{k} \exp(-\frac{|\mathcal{R}|\alpha^2 \epsilon_2^2}{3\mathbf{b}} \mathbf{c}(S^*)) \leq \delta_2$

Given an α -approximation algorithm to MAXR, we can now formally define the minimum number of RIC samples to get the bounded error guarantee ϵ as follows.

Theorem 6. Given an α -approximation algorithm to the MAXR problem and $0 \le \epsilon_1, \epsilon_2, \delta_1, \delta_2 \le 1$, let $\epsilon \ge \epsilon_1 + \epsilon_2$ and $\delta \ge 1$

$$\begin{aligned}
& \delta_1 + \delta_2. & \text{ If the number of RIC graphs satisfies:} \\
& |\mathcal{R}| \geq \frac{\mathsf{b}}{\mathsf{c}(S^*)} \max \left(\frac{2\ln(1/\delta_1)}{\epsilon_1^2}, \frac{3\ln(\binom{n}{k}/\delta_2)}{\alpha^2 \epsilon_2^2} \right) \\
& \text{ the returned seed set S will guarantee:}
\end{aligned} \tag{16}$$

$$\Pr[\mathsf{c}(S) > \alpha(1 - \epsilon)\mathsf{c}(S^*)] > 1 - \delta \tag{17}$$

 $\Pr[\mathsf{c}(S) \geq \alpha(1-\epsilon)\mathsf{c}(S^*)] \geq 1-\delta \tag{17}$ which means S is an $\alpha(1-\epsilon)$ -approximate solution to the *IMC* instance with probability at least $1 - \delta$.

Proof. Let $S_{\mathtt{OPT}}$ be an optimal solution to MAXR problem.

$$c(S) > \hat{c}(S) - \alpha \epsilon_2 c(S^*) \tag{18}$$

$$\geq \alpha \hat{\mathsf{c}}(S_{\mathtt{OPT}}) - \alpha \epsilon_2 \mathsf{c}(S^*) \geq \alpha \hat{\mathsf{c}}(S^*) - \alpha \epsilon_2 \mathsf{c}(S^*) \tag{19}$$

$$\geq \alpha(1 - \epsilon_1)\mathsf{c}(S^*) - \alpha\epsilon_2\mathsf{c}(S^*) \tag{20}$$

$$\geq \alpha (1 - \epsilon) c(S^*) \tag{21}$$

Inequality (18) happens with probability $1-\delta_2$ while inequality (20) happens with probability $1 - \delta_1$. Overall, $c(S) \ge \alpha(1 - \delta_1)$ ϵ)c(S*) with probability at least $(1-\delta_1)(1-\delta_2) \geq 1-\delta$

A drawback of this threshold is that it depends on $c(S^*)$, which is intractable. However, we can replace $c(S^*)$ by its lower bound while still making our conditions be satisfied. As long as the input k is greater than at least one threshold h_i of any community, the optimal solution always guarantees that at least one community will be influenced. We then tighten this lower bound by observation that $c(S^*) \geq \frac{\beta k}{\hbar}$ - where $\beta = \min_{i=1}^{r} b_i$ and $h = \max_{i=1}^{r} h_i$. Replacing this bound into (16), we achieve an official bound of $|\mathcal{R}|$ (denoted Ψ) that we will use in our entire algorithm.

$$\Psi = \frac{bh}{\beta k} \max\left(\frac{2\ln(1/\delta_1)}{\epsilon_1^2}, \frac{3\ln(\binom{n}{k}/\delta_2)}{\alpha^2 \epsilon_2^2}\right)$$
(22)

B. IMC Algorithm Framework

Since the bound Ψ in equation (22) can be a large number, generating Ψ samples then applying any MAXR algorithms should not be a good strategy. Therefore, we utilize the SSA method [2], [17], [18] in which we could reduce the number of generated samples but still keep $\alpha(1-\epsilon)$ approximation guarantee with probability at least $1 - \delta$. In short, the SSA

Algorithm 5 IMC Algorithmic Framework (IMCAF)

Input G, Com, k, ϵ , δ and an α -approx alg κ to MAXR **Output** An $\alpha(1-\epsilon)$ -approximation solution S with probability at least $(1 - \delta)$

```
1: Initiate \epsilon_1, \epsilon_2, \delta_1, \delta_2 such that \epsilon_1 + \epsilon_2 \leq \epsilon and \delta_1 + \delta_2 \leq \delta

2: \Psi = \frac{bh}{\beta k} \max \left( \frac{2 \ln(1/\delta_1)}{\epsilon_1^2}, \frac{3 \ln(\binom{n}{k}/\delta_2)}{\alpha^2 \epsilon_2^2} \right)

3: Initiate \epsilon_1, \epsilon_2, \epsilon_3 such that \epsilon \geq \epsilon_1 + \epsilon_2 + \epsilon_3 + \epsilon_1 \epsilon_2

4: \Lambda = (1 + \epsilon_1)(1 + \epsilon_2)\frac{3}{\epsilon_2^3} \ln \frac{3}{2\delta}
   5: \mathcal{R} \leftarrow generate \Lambda RIC samples using Alg. 1.
 6: do // stop stage
7: S \leftarrow \kappa(\mathcal{R}, k)
8: if \frac{|\mathcal{R}|}{r} \hat{c}_{\mathcal{R}}(S) \geq \Lambda then
9: \mathbf{c}^* \leftarrow \mathrm{Estimate}(G, S, \varepsilon_2, \frac{\delta}{3 \log \Psi / \Lambda}, |\mathcal{R}| \frac{1 + \varepsilon_2}{1 - \varepsilon_2} \frac{\varepsilon_3^2}{\varepsilon_2^2})
10: if \hat{c}_{\mathcal{R}}(S) \leq (1 + \varepsilon_1) \mathbf{c}^* then return S
                                    Double size of R with new RIC samples
11:
```

method keeps generating samples and stop at exponential check points to verify if there is adequate statistical evidence on the solution quality for termination. The SSA method has been successfully applied on IM. But because IMC is different to IM, some modification should be done. Therefore, we provide an IMC algorithmic framework (IMCAF) presented in Alg. 5. Comparing to SSA, IMCAF has the following modifications:

12: **while** $|\mathcal{R}| \leq \Psi$

- We setup the maximum number of samples by equation (22) (line 2)
- We use Alg. 1 to generate RIC samples (line 5 and 11)
- We find a candidate seed set S in stop stage by solving the MAXR problem, which we already proposed several approximation algorithms in section IV. (line 7)
- We use an Estimate procedure (Alg. 6) to verify whether the candidate set S is within $\alpha(1-\epsilon)$ -approximate ratio to optimal solution. (line 9)

The first and second modification are to adapt the IMC objective which is to maximize the benefit of influenced communities. The third modification comes from the fact that finding a candidate set S given a collection of samples is defined as MAXR problem, which is different to the max coverage problem in SSA. (In SSA, S could be found by solving max coverage, which is submodular). The final modification is to utilize Dagum's estimation [19] to estimate c(S). An estimation is presented by Alg. 6. The key part of Estimate procedure is that we iteratively generate RIC samples until the number of influenced samples reaches Λ' value. Line 6-9 is used to verify whether a newly generated sample is influenced by S. Based on Stopping Rule Algorithm (Section 2.1 [19]), Estimate procedure guarantees an estimation c^* of c(S) such that $c^* \geq (1 - \epsilon')c(S_k)$ with probability at least $1 - \delta'$.

Theorem 7. Given $0 < \epsilon, \delta < 1$ and an α -approximation algorithm to MAXR, IMCAF provides an $\alpha(1-\epsilon)$ -approximate guarantee to IMC with probability at least $1 - \delta$. Specifically, $\Pr[\mathsf{c}(S) \ge \alpha(1 - \epsilon)\mathsf{c}(S^*)] \ge 1 - \delta$ (23)

Proof. According to IMCAF, the seed set S would be returned when one of following conditions is satisfied: (1) $|\mathcal{R}| \geq \Psi$ or (2) $\frac{|\mathcal{R}|}{r} \hat{c}_{\mathcal{R}}(S) \ge \Lambda$ and $\hat{c}_{\mathcal{R}}(S) \le (1 + \varepsilon_1) c^*$.

By Theorem 6 and equation (22), we already proved that: if $|\mathcal{R}| \geq \Psi$, the seed set S guarantees an $\alpha(1-\epsilon)$ -approximate ratio with probability at least $1 - \delta$. So let us consider the

Algorithm 6 Estimate procedure

Input $G, S, \epsilon', \delta', T_{max}$ **Output** Estimation of c(S) with error ϵ' 1: $\Lambda' = 1 + 4(e - 2) \ln \frac{2}{\delta'} \frac{1}{(\epsilon')^2} (1 + \epsilon')$ 3: for T from 1 to T_{max} do $C, h, R(\cdot) \leftarrow$ generate a RIC sample using Alg. 1 tmp = 05: 6: **for** each node $v \in C$ **do** 7: if $R(v) \cap S \neq \emptyset$ then tmp = tmp + 1if $tmp \ge h$ then 8: Inf = Inf + 19: 10: if $Inf \geq \Lambda'$ then return $r\Lambda'/T$ 11: **Return** -1

second condition: Given $\frac{|\mathcal{R}|}{r}\hat{c}_{\mathcal{R}}(S) \geq \Lambda$, Lemma 6 in [2] proved that:

$$\Pr[\hat{c}_{\mathcal{R}}(S^*) \le (1 - \varepsilon_3)c(S^*)] \le 2\delta/3 \tag{24}$$

 $\Pr[\hat{\mathbf{c}}_{\mathcal{R}}(S^*) \leq (1-\varepsilon_3)\mathbf{c}(S^*)] \leq 2\delta/3 \tag{24}$ Denote $S_{\mathtt{OPT}}$ as an optimal solution of MAXR given the collection \mathcal{R} . We have:

$$\hat{c}_{\mathcal{R}}(S) \ge \alpha \hat{c}_{\mathcal{R}}(S_{\text{OPT}}) \ge \alpha \hat{c}_{\mathcal{R}}(S^*) \tag{25}$$

Thus, combining equation (24) and equation (25) gives:

$$\Pr[\hat{c}_{\mathcal{R}}(\hat{S}) \le \alpha (1 - \varepsilon_3) c(\hat{S}^*)] \le 2\delta/3 \tag{26}$$

Now, based on Stopping Rule Algorithm (Section 2.1 [19]), we have:

$$\Pr[\mathsf{c}^* \ge (1 + \varepsilon_2)\mathsf{c}(S)] \le \delta/3 \tag{27}$$

 $\Pr[\mathbf{c}^* \geq (1+\varepsilon_2)\mathbf{c}(S)] \leq \delta/3 \tag{27}$ Combine equation (27) and the condition $\hat{\mathbf{c}}_{\mathcal{R}}(S) \leq (1+\varepsilon_1)\mathbf{c}^*$, we have:

$$\Pr[\hat{c}_{\mathcal{R}}(S) \ge (1 + \varepsilon_1)(1 + \varepsilon_2)\mathbf{c}(S)] \le \delta/3$$
 (28)
Finally, combining equation (26) and equation (28) gives

Finally, combining equation (26) and equation (28) gives
$$\Pr[\mathsf{c}(S) \geq \alpha \frac{1 - \varepsilon_3}{(1 + \varepsilon_1)(1 + \varepsilon_2)} \mathsf{c}(S^*)] \geq 1 - \delta$$
 Since $\epsilon \geq \varepsilon_1 + \varepsilon_2 + \varepsilon_3 + \varepsilon_1 \varepsilon_2$ (line 3 alg 5), we have:
$$\Pr[\mathsf{c}(S) \geq \alpha (1 - \epsilon) \mathsf{c}(S^*)] \geq 1 - \delta$$

which completes the proof.

How to integrate the MAXR algorithms? MAF or BT(MB) algorithms could be easily integrated into Alg. 5. However, UBG algorithm has a problem that: UBG's approximation guarantee depends on R, which leads to inconsistency on approximate ratio of MAXR on each stop stage in Alg 5. To deal with this problem, we only consider an upper-bound function $\nu(\cdot)$. Running SSA with objective function $\nu(\cdot)$ would produce a solution S_{ν} that guarantees $1 - 1/e - \epsilon$ approximation ratio. By returning $S = S_{\nu}$, we have:

$$\begin{split} \mathsf{c}(S) &= \frac{\mathsf{c}(S_\nu)}{\nu(S_\nu)} \nu(S_\nu) \geq \frac{\mathsf{c}(S_\nu)}{\nu(S_\nu)} (1 - 1/e - \epsilon) \nu(S_\nu^*) \\ &\geq \frac{\mathsf{c}(S_\nu)}{\nu(S_\nu)} (1 - 1/e - \epsilon) \nu(S^*) \geq \frac{\mathsf{c}(S_\nu)}{\nu(S_\nu)} (1 - 1/e - \epsilon) \mathsf{c}(S^*) \\ \text{where } S_\nu^*, S^* \text{ are optimal solutions to the objective functions } \\ \nu(\cdot), \mathsf{c}(\cdot) \text{ respectively.} \end{split}$$

VI. EXPERIMENTAL EVALUATION

In this section, we compare the performance of our proposed algorithms with several intuitive heuristic solutions.

A. Experimental Settings

The experiments were conducted on a Linux machine with 2.3GHz Xeon 18 cores processor and 256GB of RAM. We

TABLE I: Statistics of datasets

Data	Туре	Nodes	Edges
Facebook	Undirected Directed Directed Undirected Directed	747	60.05 K
Wiki-vote		7.1 K	103.6 K
Espinions		76 K	508.8 K
DBLP		317 K	1.05 M
Pokec		1.6 M	30.6 M

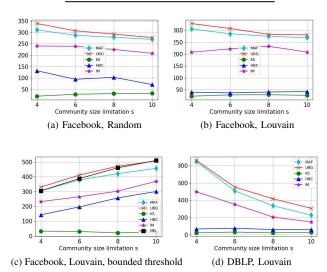


Fig. 4: Quality of solution with different community structures. Y-axis is the expected benefit of influenced communities. used the datasets from Stanford Network Analysis Project [20],

used the datasets from Stanford Network Analysis Project [20], summarized in Table I.

Settings. To partition the social networks into disjoint

Settings. To partition the social networks into disjoint communities, we utilized the well-known Louvain algorithm [21], [22], which extracts communities to optimize the network modularity. In order to see the impact of communities selection on IMC, we also use Random algorithm as the baseline. In the Random algorithm, we fix the number of communities and randomly put nodes into communities. To prevent cases in which some communities are significantly larger than the others, we limited the community size by a certain value s. If a community C was larger than s, we split it into $\lceil |C|/s \rceil$ communities. We set s=8 unless otherwise stated.

The benefit of a community was equal to its population $(b_i=|\mathtt{C}_i|)$. For each edge e=(u,v), let e's weight equal $\frac{1}{d(v)}$ where d(v) is the in-coming degree of node v (in case of undirected graph, the edge e=(u,v) will be considered as two directed edges (u,v) and (v,u)). The activation threshold of each community was set to be 2 in the experiments including MB algorithm. Otherwise, the activation threshold of each community was set to be 50% of its population $(\mathtt{h}_i=0.5|\mathtt{C}_i|)$. In all the experiments, we kept $\epsilon=\delta=0.2,\ \epsilon_1=\epsilon_2=\epsilon/2$ and $\epsilon_1=\epsilon_2=\epsilon_3=\epsilon/4$.

Baselines: To our extent, no existing algorithm was applicable to the IMC problem. Thus, we compared our proposed algorithms with the following heuristic methods:

• HBC (High Beneficial Connection). HBC selects k nodes that have high beneficial connection to the communities set. A beneficial connection of node u was defined as $B(u) = \sum_{v \in N^-(u)} w(u,v) \cdot \frac{b_{c(v)}}{\mathbf{h}_{c(v)}}$ where $N^-(u)$ is u's out-coming neighbors; $\mathbf{C}(v)$ is a community in which v

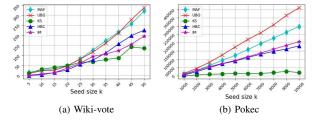


Fig. 5: Benefit of different algorithms on regular case. Y-axis is the expected benefit of influenced communities.

is a single user, and $h_{C(v)}$ is its activation threshold.

- KS (Knapsack-like Algorithm). KS considers the activation threshold of a community as a cost to influence it. Then KS finds communities set with cost constraint k to maximize the overall benefit. This is a Knapsack problem and it is possible to obtain an optimal solution in polynomial runtime. Define R as an optimal solution to the Knapsack problem: for each community $C \in R$, we selected k nodes in C and put them into the seed set k, where k is the activation threshold of C.
- IM (Influence Maximization). IM selects k nodes that maximize the influence spread. Then we estimate their expected benefit on influenced communities.

With the baseline algorithms, to evaluate the benefit of influenced communities, we used Dagum [19] estimation method with the same ϵ , θ we used for our proposed algorithm. Also, for simplicity, in this section UBG/MAF/MB means running the IMCAF framework with the corresponding MAXR algorithms. For each set of experiments, we ran ten times and report the average results. Due to space limits, we only report the important observations.

B. Performance comparison

Quality of solution. First, we compare the algorithms' performance with different community formations. Also, with each community formation, we changed value of s and recorded our results in Figure 4. In this experiment, we fix k=10. Figure 4a, 4b and 4d show the results of the experiments with Facebook and the DBLP dataset using the Louvain and Random community formation. Figure 4c shows the result in the case of bounded activation threshold. Generally, our algorithms always returned the best solutions regardless community formation methods. Also, the quality of all algorithms tended to decrease when s increases, which contradicts the experiment on bounded activation threshold. These experiments suggest that: there should be a relationship between s, the overall benefit, and the average activation threshold of influenced communities.

Next, we compare the algorithms' performance in regard to the change of k. Fig. 5 compares the quality of different algorithms without bounded activation threshold. Both of our proposed algorithms returned better solutions than any baseline algorithm. All algorithms performed similarly when k is small but the gap became noticeable when k increases. We could say UBG always returned the best solutions while KS was the worst one, which is explainable because KS does not consider the network topology and diffusion model of users when finding solution. Furthermore, the gap of quality between IM and our two solutions grew when k increased.

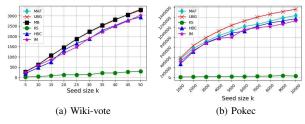


Fig. 6: Benefit of different algorithms on the bounded activation threshold case. Y-axis is the expected benefit of influenced communities.

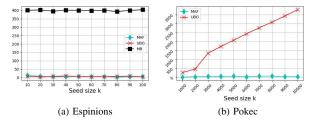


Fig. 7: Runtime of our algorithms. Y-axis is the average runtime in seconds

This is because when k is large, the number of activated users increases but scatters in the network. When we estimated the benefit of IM solution, we observed there was a huge amount of RIC samples that needed only one more activated node to be influenced, which occasionally happened in UBG and MAF. This further demonstrates the difference in objective between ICM and IM. The similar patterns were met in the experiments with bounded activation threshold (Fig. 6). Note that we discarded the MB's results in Fig. 6b because they exceeded the runtime limit.

Running time. We now compare the CPU runtime between our three solutions (Fig. 7) to evaluate the tradeoffs. Since the algorithms ran very fast (several seconds) on the small networks, we only show their performance on the large networks. As can be seen in Fig. 7b, MAF ran much faster than UBG. Also, the change in k did not affect much on MAF's performance, which totally contradicted to UBG. This could be explained by: MAF takes one-pass on all communities/nodes to calculate their impacts, another pass to sort them in order of this value and finally identify the k most influencing nodes. Meanwhile, the performance of the greedy algorithm in UBG was highly affected when k increases.

In case of bounded activation threshold, our three algorithms showed similar running time on the small-scale dataset. However, in the large networks shown in Fig. 7a, MAF and UBG outperformed MB by a huge margin. Furthermore, MB could not finish within the runtime limit in the experiments with the Pokec datasets. The poor performance of MB was because MB splits the IMC instance into O(|V|) subproblems and solves them all to get the best solution (Alg. 4). The runtime for each subproblem is equivalent to the runtime of UBG.

Evaluation of UBG. Even UBG always returned the best results among all algorithms, UBG's approximate guarantee depends on the ratio $\frac{c(S_{\nu})}{\nu(S_{\nu})}$ where S_{ν} is $1-1/e-\epsilon$ approximate solution to the submodular function $\nu(\cdot)$. Within our configura-

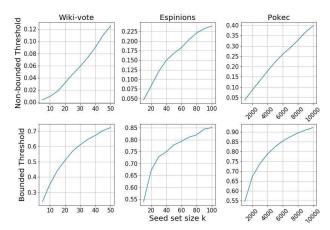


Fig. 8: Relation of UBG ratio with seed size k. Y-axis is the average value of $\frac{c(S_{\nu})}{\nu(S_{\nu})}$.

tion, we observed that: the ratio has a strong relationship with k and the average activation threshold of the communities. These relations are shown in Fig. 8. We calculated the ratio by obtaining S_{ν} and using Monte Carlo method to estimate $c(S_{\nu})$ and $\nu(S_{\nu})$. We observed that: the ratio increased and was closer to 1 when k grew. Also, the ratio increased when the average activation threshold of communities decreased. This could be observed by comparing the experiments in the same network with and without bounded activation threshold. For example in Pokec network, we obtained the ratio approximate to 1 with k = 10000 in the bounded threshold case but only get 0.4 ratio in the regular case. This is because when the activation threshold decreases, $c(\cdot)$ exhibits more "submodular" behaviors and becomes closer to $\nu(\cdot)$. $c(\cdot)$ would be totally submodular if the activation threshold is bounded by 1, which we already proved in subsection IV-C.

VII. CONCLUSION

In this paper, we introduced a novel IMC problem, on which we consider the collaborative impact of the influence spread in social networks. Solving this problem was shown to be very challenging. We proposed a new sampling method, which could be used as a baseline for any algorithm for IMC. We next devised several approaches in a trade off between effectiveness and running time. Experimental results demonstrated the effectiveness of our proposed algorithms and confirmed with the theoretical results.

VIII. ACKNOWLEDGEMENTS

This work was supported in part by NSF EFRI-1441231, NSF CNS-1814614, NSF CNS-1443905, and DTRA HDTRA1-14-1-0055.

REFERENCES

- D. Kempe, J. Kleinberg, and É. Tardos, "Maximizing the spread of influence through a social network," in *Proceedings of the ninth ACM SIGKDD international conference on Knowledge discovery and data mining*. ACM, 2003, pp. 137–146.
 H. T. Nguyen, M. T. Thai, and T. N. Dinh, "Stop-and-stare: Optimal
- [2] H. T. Nguyen, M. T. Thai, and T. N. Dinh, "Stop-and-stare: Optimal sampling algorithms for viral marketing in billion-scale networks," in Proceedings of the 2016 International Conference on Management of Data. ACM.
- [3] 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 SIGKDD international conference on Knowledge discovery and data mining. ACM, 2010, pp. 1029–1038.

- [4] Y. Tang, X. Xiao, and Y. Shi, "Influence maximization: Near-optimal time complexity meets practical efficiency," in Proceedings of the 2014
- ACM SIGMOD international conference on Management of data.

 [5] C. Borgs, M. Brautbar, J. Chayes, and B. Lucier, "Maximizing social influence in nearly optimal time," in *Proceedings of the twenty-fifth annual ACM-SIAM symposium on Discrete algorithms*. SIAM, 2014.
- B. Wellman, J. Salaff, D. Dimitrova, L. Garton, M. Gulia, and C. Haythornthwaite, "Computer networks as social networks: Collaborative work, telework, and virtual community," Annual review of sociology,
- vol. 22, no. 1, pp. 213–238, 1996. T. Pan, S. Mishra, L. N. Nguyen, G. Lee, J. Kang, J. Seo, and M. T. Thai, "Threat from being social: Vulnerability analysis of social network coupled smart grid," *IEEE Access*, vol. 5, pp. 16774–16783, 2017. H. T. Nguyen, T. N. Dinh, and M. T. Thai, "Cost-aware targeted
- viral marketing in billion-scale networks," in INFOCOM 2016-The 35th Annual IEEE International Conference on Computer Communications.
- R. Impagliazzo and R. Paturi, "On the complexity of k-sat," Journal of Computer and System Sciences, vol. 62, no. 2, pp. 367-375, 2001.
- [10] Y. Lin, W. Chen, and J. C. Liu, "Boosting information spread: An algorithmic approach," in *Data Engineering (ICDE)*, 2017 IEEE 33rd International Conference on. IEEE, 2017, pp. 883–894.
 [11] W. Lu, W. Chen, and L. V. Lakshmanan, "From competition to
- complementarity: comparative influence diffusion and maximization,' Proceedings of the VLDB Endowment, vol. 9, no. 2, pp. 60-71, 2015.
- [12] L. Ma, G. Cao, and L. Kaplan, "Graphical approach for influence maximization in social networks under generic threshold-based nonsubmodular model," in Big Data (Big Data), 2017 IEEE International Conference on. IEEE, 2017, pp. 970-975.
- [13] X. Li, J. D. Smith, T. N. Dinh, and M. T. Thai, "Why approximate when you can get the exact? optimal targeted viral marketing at scale," in IEEE INFOCOM 2017-IEEE Conference on Computer Communications.
- IEEE, 2017, pp. 1–9.
 [14] P. Manurangsi, "Almost-polynomial ratio eth-hardness of approximating densest k-subgraph," in Proceedings of the 49th Annual ACM SIGACT Symposium on Theory of Computing. ACM, 2017, pp. 954–961. "Markov's inequality," https://en.wikipedia.org/wiki/Markov%
- https://en.wikipedia.org/wiki/Markov%27s_
- inequality.

 [16] Y. Tang, Y. Shi, and X. Xiao, "Influence maximization in near-linear time: A martingale approach," in *Proceedings of the 2015 ACM SIG-MOD International Conference on Management of Data.* ACM, 2015.

 [17] K. Huang, S. Wang, G. Bevilacqua, X. Xiao, and L. V. Lakshmanan,
- "Revisiting the stop-and-stare algorithms for influence maximization," *Proceedings of the VLDB Endowment*, 2017.
- [18] H. T. Nguyen, T. N. Dinh, and M. T. Thai, "Revisiting of revisiting the stop-and-stare algorithms for influence maximization," in *Proceedings of* the 2018 International Conference on Computational Data and Social Networks, 2018.
- [19] P. Dagum, R. Karp, M. Luby, and S. Ross, "An optimal algorithm for monte carlo estimation," SIAM Journal on computing, 2000.
- "Stanford network analysis Leskovec et al., http://snap.stanford.edu, 2010.
- V. D. Blondel, J.-L. Guillaume, R. Lambiotte, and E. Lefebvre, "Fast unfolding of communities in large networks," Journal of statistical mechanics: theory and experiment, vol. 2008, no. 10, p. P10008, 2008.
- [22] N. Dugué and A. Perez, "Directed louvain: maximizing modularity in directed networks," Ph.D. dissertation, Université d'Orléans, 2015.