

# Revenue Maximization in Incentivized Social Advertising

Cigdem Aslay  
ISI Foundation  
Turin, Italy  
cigdem.aslay@isi.it

Francesco Bonchi  
ISI Foundation  
Turin, Italy  
francesco.bonchi@isi.it

Laks V.S. Lakshmanan  
Univ. of British Columbia  
Vancouver, Canada  
laks@cs.ubc.ca

Wei Lu  
LinkedIn Corp.  
Sunnyvale, CA, USA  
wlu@linkedin.com

## ABSTRACT

Incentivized social advertising, an emerging marketing model, provides monetization opportunities not only to the owners of the social networking platforms but also to their influential users by offering a “cut” on the advertising revenue. We consider a social network (the host) that sells ad-engagements to advertisers by inserting their ads, in the form of promoted posts, into the feeds of carefully selected “initial endorsers” or seed users: these users receive monetary incentives in exchange for their endorsements. The endorsements help propagate the ads to the feeds of their followers. Whenever any user of the platform engages with an ad, the host is paid some fixed amount by the advertiser, and the ad further propagates to the feed of her followers, potentially recursively. In this context, the problem for the host is to allocate ads to influential users, taking into account the propensity of ads for viral propagation, and carefully apportioning the monetary budget of each of the advertisers between incentives to influential users and ad-engagement costs, with the rational goal of maximizing its own revenue. In particular, we consider a monetary incentive for the influential users, which is proportional to their influence potential.

We show that, taking all important factors into account, the problem of revenue maximization in incentivized social advertising corresponds to the problem of monotone submodular function maximization, subject to a partition matroid constraint on the ads-to-seeds allocation, and submodular knapsack constraints on the advertisers’ budgets. We show that this problem is NP-hard and devise two greedy algorithms with provable approximation guarantees, which differ in their sensitivity to seed user incentive costs.

Our approximation algorithms require repeatedly estimating the expected marginal gain in revenue as well as in advertiser payment. By exploiting a connection to the recent advances made in scalable estimation of expected influence spread, we devise efficient and scalable versions of our two greedy algorithms. An extensive experimental assessment confirms the high quality of our proposal.

## 1. INTRODUCTION

The rise of online advertising platforms has generated new opportunities for advertisers in terms of personalizing and targeting their marketing messages. When users access a platform, they leave a trail of information that can be correlated with their consumption tastes, enabling better targeting options for advertisers. Social networking platforms particularly can gather large amounts of users’ shared posts that stretches beyond general demographic and geographic data. This offers more advanced interest, behavioral, and

connection-based targeting options, enabling a level of personalization that is not achievable by other online advertising channels. Hence, advertising on social networking platforms has been one of the fastest growing sectors in the online advertising landscape: a market that did not exist until Facebook launched its first advertising service in May 2005, is projected to generate \$11 billion revenue by 2017, almost doubling the 2013 revenue<sup>1</sup>.

**Social advertising.** Social advertising models are typically employed by platforms such as Twitter, Tumblr, and Facebook through the implementation of *promoted posts* that are shown in the “news feed” of their users.<sup>2</sup> A promoted post can be a video, an image, or simply a textual post containing an advertising message. Social advertising models of this type are usually associated with a *cost per engagement* (CPE) pricing scheme: the advertiser does not pay for the ad impressions, but pays the platform owner (hereafter referred to as the *host*) only when a user actively engages with the ad. The *engagement* can be in the form of a social action such as “like”, “share”, or “comment”: in this paper we blur the distinction between these different types of actions, and generically refer to them all as *engagements* or *clicks* interchangeably.

Similar to organic (i.e., non-promoted) posts, promoted posts can propagate from user to user in the network<sup>3</sup>, potentially triggering a viral contagion: whenever a user  $u$  engages with an ad  $i$ , the host is paid some fixed amount by the advertiser (the CPE). Furthermore,  $u$ ’s engagement with  $i$  appears in the feed of  $u$ ’s followers, who are then exposed to ad  $i$  and could in turn be influenced to engage with  $i$ , producing further revenue for the host [5, 35].

**Incentivized social advertising.** In this paper, we study the novel model of *incentivized social advertising*. Under this model, users selected by the host as *seeds* for the campaign on a specific ad  $i$ , can take a “cut” on the social advertising revenue. These users are typically selected because they are influential or authoritative on the specific topic, brand, or market of  $i$ .

A recent report<sup>4</sup> indicates that Facebook is experimenting with the idea of incentivizing users. YouTube launched a revenue-sharing program for prominent users in 2007. Twitch, the streaming platform of choice for gamers, lets partners make money through revenue sharing, subscriptions, and merchandise sales. YouNow, a streaming platform popular among younger users, earns money by taking a cut of the tips and digital gifts that fans give its stars. On platforms without partner deals, including Twitter and

<sup>1</sup><http://www.unified.com/historyofsocialadvertising/>

<sup>2</sup>According to a recent report, Facebook’s news feed ads have 21 times higher click-through rate than standard web retargeting ads and 49 times the click-through rate of Facebook’s right-hand side display ads: see <https://blog.adroll.com/trends/facebook-exchange-news-feed-numbers>.

<sup>3</sup>Tumblr’s CEO D. Karp reported (CES 2014) that a normal post is reposted on average 14 times, while promoted posts are on average reposted more than 10000 times: <http://yhoo.it/1vFfIAC>.

<sup>4</sup><http://www.theverge.com/2016/4/19/11455840/facebook-tip-jar-partner-program-monetization>

Snapchat, celebrity users often strike sponsored deals to include brands in their posts, which suggests potential monetization opportunities for Twitter and Snapchat<sup>5</sup>.

In this work, we consider incentives that are determined by the topical influence of the seed users for the specific ad. More concretely, given an ad  $i$ , the financial incentive that a seed user  $u$  would get for engaging with  $i$  is a function of the social influence that  $u$  has exhibited in the past in the topic of  $i$ . For instance, a user who often produces relevant content about long-distance running, capturing the attention of a relatively large audience, might be a good seed for endorsing a new model of running shoes. In this case, her past demonstrated influence on this very topic would be taken into consideration when defining the lumpsum amount for her engagement with the new model of running shoes. The same user could be considered as a seed for a new model of tennis shoes, but in that case the incentive might be lower, due to her lower past influence demonstrated. To summarize, incentives are paid by the host to users selected as seeds. These incentives count as seeding costs and depend on the topic of the ad and the user’s past demonstrated influence in the topic.

The incentive model above has several advantages. First, it captures in a uniform framework both the “celebrity-influencer”, whose incentives are naturally very high (like her social influence), and who are typically preferred by more traditional types of advertising such as TV ads; as well as the “ordinary-influencer” [6], a non-celebrity individual who is an expert in some specific topic, and thus has a relatively restricted audience, or tribe, that trust her. Second, incentives not only play their main role, i.e., encourage the seed users to endorse an advert campaign, but also, as a by-product, they incentivize users of the social media platform to become influential in some topics by actively producing *good-quality content*. This has an obvious direct benefit for the social media platform.

**Revenue maximization.** In the context of incentivized social advertising, we study the fundamental problem of revenue maximization from the host perspective: an advertiser enters into an agreement with the host to pay, following the CPE model, a fixed price  $cpe(i)$  for each engagement with ad  $i$ . The agreement also specifies the finite budget  $B_i$  of the advertiser for the campaign for ad  $i$ . The host has to carefully select the seed users for the campaign: given the maximum amount  $B_i$  that it can receive from the advertiser, the host must try to achieve as many engagements on the ad  $i$  as possible, while spending as little as possible on the incentives for “seed” users. The host’s task gets even more challenging by having to simultaneously accommodate multiple campaigns by different advertisers. Moreover, for a fixed time window (e.g., 1 day, or 1 week), the host can select each user as the seed endorser for at most one ad: this constraint maintains higher credibility for the endorsements and avoids the undesirable situation where, e.g., the same sport celebrity endorses Nike and Adidas in the same time window. Therefore two ads  $i$  and  $j$ , which are in the same topical area, naturally compete for the influential users in that area.

We show that, taking all important factors (such as topical relevance of ads, their propensity for social propagation, the topical influence of users, seed incentives and advertiser budgets) into account, the problem of revenue maximization in incentivized social advertising corresponds to the problem of *monotone submodular function maximization subject to a partition matroid constraint on the ads-to-seeds allocation, and submodular knapsack constraints on the advertisers’ budgets*. This problem is NP-hard and furthermore is far more challenging than the classical influence maximization problem (IM) [24] and its variants. For this problem, we de-

velop two natural greedy algorithms, for which we provide formal approximation guarantees. The two algorithms differ in their sensitivity to cost-effectiveness in the seed user selection:

- *Cost-Agnostic Greedy Algorithm* (CA-GREEDY), which greedily chooses the seed users based on the marginal gain in the revenue, without using any information about the users’ incentive costs;
- *Cost-Sensitive Greedy Algorithm* (CS-GREEDY), which greedily chooses the seed users based on the *rate* of marginal gain in revenue per marginal gain in the advertiser’s payment for each advertiser.

Our results generalize the results of Iyer *et al.* [22, 23] on submodular function maximization by (i) generalizing from a single submodular knapsack constraint to multiple submodular knapsack constraints, and (ii) by handling an additional partition matroid constraint. Our theoretical analysis leverages the notion of curvature of submodular functions.

Our approximation algorithms require repeatedly estimating the expected marginal gain in revenue as well in advertiser payment. We leverage recent advances in scalable estimation of expected influence spread and devise scalable algorithms for revenue maximization in our model.

### Contributions and roadmap.

- We propose *incentivized* social advertising, and formulate a fundamental problem of revenue maximization from the host perspective, when the incentives paid to the seed users are determined by their demonstrated past influence in the topic of the specific ad (Section 2).
- We prove the hardness of our problem and we devise two greedy algorithms with approximation guarantees. The first (CA-GREEDY) is agnostic to users’ incentives during the seed selection while the other (CS-GREEDY) is not (Section 3).
- We devise scalable versions of our approximation algorithms (Section 4). Our comprehensive experimentation on real-world datasets (Section 5) confirms the scalability of our methods and shows that the scalable version of CS-GREEDY consistently outperforms that of CA-GREEDY, and is far superior to natural baselines, thanks to a mindful allocation of budget on incentives.

Related work is discussed in Section 6 while Section 7 concludes the paper discussing future work.

## 2. PROBLEM STATEMENT

**Business model: the advertiser.** An advertiser<sup>6</sup>  $i$  enters into an agreement with the *host*, the owner of the social networking platform, for an incentivized social advertising campaign on its ad. The advertiser agrees to pay the host:

1. an incentive  $c_i(u)$  for each seed user  $u$  chosen to endorse ad  $i$ ; we let  $S_i$  denote the set of users selected to endorse ad  $i$ ;
2. a cost-per-engagement amount  $cpe(i)$  for each user that engages with (e.g., clicks) its ad  $i$ .

An advertiser  $i$  has a finite budget  $B_i$  that limits the amount it can spend on the campaign for its ad.

**Business model: the host.** The host receives from advertiser  $i$ :

1. a description of the ad  $i$  (e.g., a set of keywords) which allows the host to map the ad to a distribution  $\gamma_i$  over a latent topic space (described in more detail later);

<sup>5</sup><http://www.wsj.com/articles/more-marketers-offer-incentives-for-watching-ads-1451991600>

<sup>6</sup>We assume each advertiser has one ad to promote per time window, and use  $i$  to refer to the  $i$ -th advertiser and its ad interchangeably.

2. a commercial agreement that specifies the cost-per-engagement amount  $cpe(i)$  and the campaign budget  $B_i$ .

The host is in charge of running the campaign, by selecting which users and how many to allocate as a seed set  $S_i$  for each ad  $i$ , and by determining their incentives. Given that these decisions must be taken *before* the campaign is started, the host has to reason in terms of *expectations* based on past performance. Let  $\sigma_i(S_i)$  denote the *expected number of clicks* ad  $i$  receives when using  $S_i$  as the seed set of incentivized users. The host models the total payment that advertiser  $i$  needs to make for its campaign, denoted  $\rho_i(S_i)$ , as the sum of its total costs for the *expected* ad-engagements (e.g., clicks), and for incentivizing its seed users: i.e.,  $\rho_i(S_i) = \pi_i(S_i) + c_i(S_i)$  where  $\pi_i(S_i) = cpe(i) \cdot \sigma_i(S_i)$  and  $c_i(S_i) := \sum_{u \in S_i} c_i(u)$ , where  $c_i(u)$  denotes the incentive paid to a candidate seed user  $u$  for ad  $i$ . We assume  $c_i(u)$  is a monotone function  $f$  of the influence potential of  $u$ , capturing the intuition that seeds with higher expected spread cost more: i.e.,  $c_i(u) := f(\sigma_i(\{u\}))$ .

Notice that the expected revenue of the host from the engagements to ad  $i$  is just  $\pi_i(S_i)$ , as the cost  $c_i(S_i)$  paid by the advertiser to the host for the incentivizing influential users, is in turn paid by the host to the seeds. In this setting, the host faces the following trade-off in trying to maximize its revenue. Intuitively, targeting influential seeds would increase the expected number of clicks, which in turn could yield a higher revenue. However, influential seeds cost more to incentivize. Since the advertiser has a fixed overall budget for its campaign, the higher seeding cost may come at the expense of reduced revenue for the host. Finally, an added challenge is that the host has to serve many advertisers at the same time, with potentially competitive ads, i.e., ads which are very close in the topic space.

**Data model, topic model, and propagation model.** The host owns: a *directed graph*  $G = (V, E)$  representing the social network, where an arc  $(u, v)$  means that user  $v$  follows user  $u$ , and thus  $v$  can see  $u$ 's posts and may be influenced by  $u$ . The host also owns a *topic model* for ads and users' interests, defined by a hidden variable  $Z$  that can range over  $L$  latent topics. A topic distribution thus abstracts the interest pattern of a user and the relevance of an ad to those interests. More precisely, the topic model maps each ad  $i$  to a distribution  $\gamma_i^z$  over the latent topic space:

$$\gamma_i^z = \Pr(Z = z|i), \text{ with } \sum_{z=1}^L \gamma_i^z = 1.$$

Finally, the host uses a topic-aware influence propagation model defined on the social graph  $G$  and the topic model. The propagation model governs the way in which ad impressions propagate in the social network, driven by topic-specific influence. In this work, we adopt the *Topic-aware Independent Cascade* model<sup>7</sup> (TIC) proposed by Barbieri et al. [8] which extends the standard *Independent Cascade* (IC) model [24]: In TIC, an ad is represented by a topic distribution, and the influence strength from user  $u$  to  $v$  is also topic-dependent, i.e., there is a probability  $p_{u,v}^z$  for each topic  $z$ . In this model, when a node  $u$  clicks an ad  $i$ , it gets one chance of influencing each of its out-neighbors  $v$  that has not clicked  $i$ . This event succeeds with a probability equal to the weighted average of the arc probabilities w.r.t. the topic distribution of ad  $i$ :

$$p_{u,v}^i = \sum_{z=1}^L \gamma_i^z \cdot p_{u,v}^z. \quad (1)$$

<sup>7</sup>Note that the use of the topic-based model is orthogonal to the technical development and contributions of our work. Specifically, if we assume that the topic distributions of all ads and users are identical, the TIC model reduces to the standard IC model. The techniques and results in the paper remain intact.

Using this stochastic propagation model the host can determine the *expected spread*  $\sigma_i(S_i)$  of a given campaign for ad  $i$  when using  $S_i$  as seed set. For instance, the influence value of a user  $u$  for ad  $i$  is defined as the expected spread of the singleton seed  $\{u\}$  for the given the description for ad  $i$ , under the TIC model, i.e.,  $\sigma_i(\{u\})$ : this is the quantity that is used to determine the incentive for a candidate seed user  $u$  to endorse the ad  $i$ .

**The revenue maximization problem.** Hereafter we assume a fixed time window (say a 24-hour period) in which the revenue maximization problem is defined. Within this time window we have  $h$  advertisers with ad description  $\gamma_i$ , cost-per-engagement  $cpe(i)$ , and budget  $B_i$ ,  $i \in [h]$ . We define an *allocation*  $\vec{S}$  as a vector of  $h$  *pair-wise disjoint* sets  $(S_1, \dots, S_h) \in 2^V \times \dots \times 2^V$ , where  $S_i$  is the seed set assigned to advertiser  $i$  to start the ad-engagement propagation process. Within the time window, each user in the platform can be selected to be seed for at most one ad, that is,  $S_i \cap S_j = \emptyset$ ,  $i, j \in [h]$ . We denote the total revenue of the host from advertisers as the sum of the ad-specific revenues:

$$\pi(\vec{S}) = \sum_{i \in [h]} \pi_i(S_i).$$

Next, we formally define the revenue maximization problem for incentivized social advertising from the host perspective. Note that given an instance of the TIC model on a social graph  $G$ , for each ad  $i$ , the ad-specific influence probabilities are determined by Eq. (1).

**Problem 1 (REVENUE-MAXIMIZATION (RM)).** *Given a social graph  $G = (V, E)$ ,  $h$  advertisers, cost-per-engagement  $cpe(i)$  and budget  $B_i$ ,  $i \in [h]$ , ad-specific influence probabilities  $p_{u,v}^i$  and seed user incentive costs  $c_i(u)$ ,  $u, v \in V$ ,  $i \in [h]$ , find a feasible allocation  $\vec{S}$  that maximizes the host's revenue:*

$$\begin{aligned} & \underset{\vec{S}}{\text{maximize}} && \pi(\vec{S}) \\ & \text{subject to} && \rho_i(S_i) \leq B_i, \forall i \in [h], \\ & && S_i \cap S_j = \emptyset, i \neq j, \forall i, j \in [h]. \end{aligned}$$

In order to avoid degenerate problem instances, we assume that no single user incentive exceeds any advertiser's budget. This ensures that every advertiser can afford at least one seed node.

### 3. HARDNESS AND APPROXIMATION

**Hardness.** We first show that Problem 1 (RM) is NP-hard. We recall that a set function  $f : 2^U \rightarrow \mathbb{R}_{\geq 0}$  is monotone if for  $S \subset T \subseteq U$ ,  $f(S) \leq f(T)$ . We define the marginal gain of an element  $x$  w.r.t.  $S \subset U$  as  $f(x|S) := f(S \cup \{x\}) - f(S)$ . A set function  $f$  is submodular if for  $S \subset T \subset U$  and  $x \in U \setminus T$ ,  $f(x|T) \leq f(x|S)$ , i.e., the marginal gains diminish with larger sets.

It is well known that the influence spread function  $\sigma_i(\cdot)$  is monotone and submodular [24], from which it follows that the ad-specific revenue function  $\pi_i(\cdot)$  is monotone and submodular. Finally, since the total revenue function,  $\pi(\vec{S}) = \sum_{i \in [h]} \pi_i(S_i)$ , is a non-negative linear combination of monotone and submodular functions, these properties carry over to  $\pi(\vec{S})$ . Likewise, for each ad  $i$ , the payment function  $\rho_i(\cdot)$  is a non-negative linear combination of two monotone and submodular functions,  $\pi_i(\cdot)$  and  $c_i(\cdot)$ , and so is also monotone and submodular. Thus, the constraints  $\rho_i(S_i) \leq B_i$ ,  $i \in [h]$ , in Problem 1 are submodular knapsack constraints. We start with our hardness result.

**Theorem 1.** *Problem 1 (RM) is NP-hard.*

*Proof.* Consider the special case with one advertiser, i.e.,  $h = 1$ . Then we have one submodular knapsack constraint and no partition

matroid constraint. This corresponds to maximizing a submodular function subject to a submodular knapsack constraint, the so-called Submodular Cost Submodular Knapsack (SCSK) problem, which is known to be NP-hard [23]. Since this is a special case of Problem 1, the claim follows.  $\square$

Next, we characterize the constraint that the allocation  $\vec{S} = (S_1, \dots, S_h)$  should be composed of pairwise disjoint sets, i.e.,  $S_i \cap S_j = \emptyset, i \neq j, \forall i, j \in [h]$ . We will make use of the following notions on matroids.

**Definition 1 (Independence System).** A set system  $(\mathcal{E}, \mathcal{I})$  defined with a finite ground set  $\mathcal{E}$  of elements, and a family  $\mathcal{I}$  of subsets of  $\mathcal{E}$  is an independence system if  $\mathcal{I}$  is non-empty and if it satisfies downward closure axiom, i.e.,  $X \in \mathcal{I} \wedge Y \subseteq X \rightarrow Y \in \mathcal{I}$ .

**Definition 2 (Matroid).** An independence system  $(\mathcal{E}, \mathcal{I})$  is a matroid  $\mathfrak{M} = (\mathcal{E}, \mathcal{I})$  if it also satisfies the augmentation axiom: i.e.,  $X \in \mathcal{I} \wedge Y \in \mathcal{I} \wedge |Y| > |X| \rightarrow \exists e \in Y \setminus X : X \cup \{e\} \in \mathcal{I}$ .

**Definition 3 (Partition Matroid).** Let  $\mathcal{E}_1, \dots, \mathcal{E}_l$  be a partition of the ground set  $\mathcal{E}$  into  $l$  non-empty disjoint subsets. Let  $d_i$  be an integer,  $0 \leq d_i \leq |\mathcal{E}_i|$ . In a partition matroid  $\mathfrak{M} = (\mathcal{E}, \mathcal{I})$ , a set  $X$  is defined to be independent iff, for every  $i, 1 \leq i \leq l, |X \cap \mathcal{E}_i| \leq d_i$ . That is,  $\mathcal{I} = \{X \subseteq \mathcal{E} : |X \cap \mathcal{E}_i| \leq d_i, \forall i = 1, \dots, l\}$ .

**Lemma 1.** The constraint that in an allocation  $\vec{S} = (S_1, \dots, S_h)$ , the seed sets  $S_i$  are pairwise disjoint is a partition matroid constraint over the ground set  $\mathcal{E}$  of all (node, advertiser) pairs.

*Proof.* Given  $G = (V, E)$ ,  $|V| = n$ , and a set  $A = \{i : i \in [h]\}$  of advertisers, let  $\mathcal{E} = V \times A$  denote the ground set of all (node, advertiser) pairs. Define  $\mathcal{E}_u = \{(u, i) : i \in A\}, u \in V$ . Then the set  $\{\mathcal{E}_u : \forall u \in V\}$  forms a partition of  $\mathcal{E}$  into  $n$  disjoint sets, i.e.,  $\mathcal{E}_u \cap \mathcal{E}_v = \emptyset, u \neq v$ , and  $\bigcup_{u \in V} \mathcal{E}_u = \mathcal{E}$ . Given a subset  $\mathcal{X} \subseteq \mathcal{E}$ , define

$$S_i = \{u : (u, i) \in \mathcal{X}\}.$$

Then it is easy to see that the sets  $S_i, i \in [h]$  are pairwise disjoint iff the set  $\mathcal{X}$  satisfies the constraint

$$|\mathcal{X} \cap \mathcal{E}_u| \leq 1, \forall u \in V.$$

The lemma follows on noting that the set system  $\mathfrak{M} = (\mathcal{E}, \mathcal{I})$ , where  $\mathcal{I} = \{\mathcal{X} \subseteq \mathcal{E} : |\mathcal{X} \cap \mathcal{E}_u| \leq 1, \forall u \in V\}$  is actually a partition matroid.  $\square$

Therefore, the RM problem corresponds to the problem of submodular function maximization subject to a partition matroid constraint  $\mathfrak{M} = (\mathcal{E}, \mathcal{I})$ , and  $h$  submodular knapsack constraints.

**Approximation analysis.** Next lemma states that the constraints of the RM problem together form an independence system defined on the ground set  $\mathcal{E}$ . This property will be leveraged later in developing approximation algorithms. Given the partition matroid constraint  $\mathfrak{M} = (\mathcal{E}, \mathcal{I})$ , and  $h$  submodular knapsack constraints, let  $\mathcal{C}$  denote the family of subsets, defined on  $\mathcal{E}$ , that are feasible solutions to the RM problem.

**Lemma 2.** The system  $(\mathcal{E}, \mathcal{C})$  is an independence system.

*Proof.* For each knapsack constraint  $\rho_i(\cdot) \leq B_i$ , let  $\mathcal{F}_i \subseteq 2^V$  denote the collection of feasible subsets of  $V$ , i.e.,

$$\mathcal{F}_i = \{S_i \subseteq V : \rho_i(S_i) \leq B_i\}.$$

The set system  $(V, \mathcal{F}_i)$  defined by the set of feasible solutions to any knapsack constraint is downward-closed, hence is an independence system. Given  $\mathcal{F}_i, \forall i \in [h]$  and the partition matroid constraint  $\mathfrak{M} = (\mathcal{E}, \mathcal{I})$ , we can define the family of subsets of  $\mathcal{E}$  that are feasible solutions to the RM problem as follows:

$$\mathcal{C} = \{\mathcal{X} : \mathcal{X} \in \mathcal{I} \text{ and } S_i \in \mathcal{F}_i, \forall i \in [h]\}$$

where  $S_i = \{u : (u, i) \in \mathcal{X}\}$ . Let  $\mathcal{X} \in \mathcal{C}$  and  $\mathcal{X}' \subseteq \mathcal{X}$ . In order to show that  $\mathcal{C}$  is an independence system, it suffices to show that  $\mathcal{X}' \in \mathcal{C}$ .

Let  $S'_i = \{u : (u, i) \in \mathcal{X}'\}, i \in [h]$ . Clearly,  $S'_i \subseteq S_i$ . As each single knapsack constraint  $\rho_i(\cdot) \leq B_i$  is associated with the independence system  $(V, \mathcal{F}_i)$ , we have  $S'_i \in \mathcal{F}_i$  for any  $S'_i \subseteq S_i, i \in [h]$ .

Next, as  $\mathcal{X} \in \mathcal{I}$ , we have  $S_i \cap S_j = \emptyset$ . Since  $\mathfrak{M} = (\mathcal{E}, \mathcal{I})$  is a partition matroid, by downward closure,  $\mathcal{X}' \in \mathcal{I}$ , and hence  $S'_i \cap S'_j = \emptyset, i \neq j$ . We just proved  $\mathcal{X}' \in \mathcal{C}$ , verifying that  $\mathcal{C}$  is an independence system.  $\square$

Our theoretical guarantees for our approximation algorithms to the RM problem depend on the notion of *curvature* of submodular functions. Recall that  $f(j|S), j \notin S$ , denotes the marginal gain  $f(S \cup \{j\}) - f(S)$ .

**Definition 4 (Curvature).** [15] Given a submodular function  $f$ , the total curvature  $\kappa_f$  of  $f$  is defined as

$$\kappa_f = 1 - \min_{j \in V} \frac{f(j|V \setminus \{j\})}{f(\{j\})},$$

and the curvature  $\kappa_f(S)$  of  $f$  wrt a set  $S$  is defined as

$$\kappa_f(S) = 1 - \min_{j \in S} \frac{f(j|S \setminus \{j\})}{f(\{j\})}.$$

It is easy to see that  $0 \leq \kappa_f = \kappa_f(V) \leq 1$ . Intuitively, the curvature of a function measures the deviation of  $f$  from *modularity*: modular functions have a curvature of 0, and the further away  $f$  is from modularity, the larger  $\kappa_f$  is. Similarly, the curvature  $\kappa_f(S)$  of  $f$  wrt a set  $S$  reflects how much the marginal gains  $f(j|S)$  can decrease as a function of  $S$ , measuring the deviation from modularity, given the context  $S$ . Iyer et al. [22] introduced the notion of *average curvature*  $\hat{\kappa}_f(S)$  of  $f$  wrt a set  $S$  as

$$\hat{\kappa}_f(S) = 1 - \frac{\sum_{j \in S} f(j|S \setminus \{j\})}{\sum_{j \in S} f(\{j\})},$$

and showed the following relation between these several forms of curvature:

$$0 \leq \hat{\kappa}_f(S) \leq \kappa_f(S) \leq \kappa_f(V) = \kappa_f \leq 1.$$

In the next subsections, we propose two greedy approximation algorithms for the RM problem. The first of these, Cost-Agnostic Greedy Algorithm (CA-GREEDY), greedily chooses the seed users solely based on the marginal gain in the revenue, without considering seed user incentive costs. The second, Cost-Sensitive Greedy Algorithm (CS-GREEDY), greedily chooses the seed users based on the *rate* of marginal gain in revenue per marginal gain in the advertiser's payment for each advertiser.

We note that Iyer et al. [22, 23] study a restricted special case of the RM problem, referred as Submodular-Cost Submodular-Knapsack (SCSK), and propose similar cost-agnostic and cost-sensitive algorithms. Our results extend theirs in two major ways. First, we extend from a single advertiser to multiple advertisers (i.e., from a single submodular knapsack constraint to multiple submodular knapsack constraints). Second, unlike SCSK, our RM problem is subject to an additional partition matroid constraint on the ads-to-seeds allocation, which naturally arises when multiple advertisers are present.

### 3.1 Cost-Agnostic Greedy Algorithm

The Cost-Agnostic Greedy Algorithm (CA-GREEDY) for the RM problem, whose pseudocode is provided in Algorithm 1, chooses at each iteration a (node, advertiser) pair that provides the maximum increase in the revenue of the host. Let  $\mathcal{X}_g \subseteq \mathcal{E}$  denote the greedy solution set of (node, advertiser) pairs, returned by CA-GREEDY, having one-to-one correspondence with the greedy allocation  $\vec{S}_g$ , i.e.,  $S_i = \{u : (u, i) \in \mathcal{X}_g\}$ ,  $\forall S_i \in \vec{S}_g$ . Let  $\mathcal{X}_g^t$  denote the greedy solution after  $t$  iterations of CA-GREEDY. At each iteration  $t$ , CA-GREEDY first finds the (node, advertiser) pair  $(u^*, i^*)$  that maximizes  $\pi_i(u \mid S_i^{t-1})$ , and tests whether adding this pair to the current greedy solution  $\mathcal{X}_g^{t-1}$  would violate any constraint: if  $\mathcal{X}_g^{t-1} \cup \{(u^*, i^*)\}$  is feasible, the pair  $(u^*, i^*)$  is added to the greedy solution as the  $t$ -th (node, advertiser) pair. Otherwise,  $(u^*, i^*)$  is removed from the current ground set of (node, advertiser) pairs  $\mathcal{E}^{t-1}$ . CA-GREEDY terminates when there is no feasible (node, advertiser) pair left in the current ground set  $\mathcal{E}^{t-1}$ .

**Observation 1.** *Being monotone and submodular, the total revenue function  $\pi(\vec{S}_g)$  has a total curvature  $\kappa_\pi$ , given by:*

$$\kappa_\pi = 1 - \min_{(u, i) \in \mathcal{E}} \frac{\pi_i(u \mid V \setminus \{u\})}{\pi_i(\{u\})}.$$

*Proof.* Let  $g : 2^\mathcal{E} \mapsto \mathbb{R}_{\geq 0}$  be monotone and submodular. Then, the total curvature  $\kappa_g$  of  $g$  is defined as follows:

$$\kappa_g = 1 - \min_{x \in \mathcal{E}} \frac{g(x \mid \mathcal{E} \setminus \{x\})}{g(\{x\})},$$

where  $x = (u, i) \in \mathcal{E}$ . Using the one-to-one correspondence between  $\mathcal{X}_g$  and  $\vec{S}_g$ , we can alternatively formulate the RM problem as follows:

$$\begin{aligned} & \underset{\mathcal{X} \subseteq \mathcal{E}}{\text{maximize}} && g(\mathcal{X}) \\ & \text{subject to} && \mathcal{X} \in \mathcal{C}. \end{aligned}$$

where  $g(\mathcal{X}) = \sum_{i \in [h]} \pi_i(S_i)$  with  $S_i = \{u : (u, i) \in \mathcal{X}\}$ .

Using this correspondence, we can rewrite  $\kappa_g$  as  $\kappa_\pi$  as follows:

$$\kappa_g = \kappa_\pi = 1 - \min_{(u, i) \in \mathcal{E}} \frac{\pi_i(\{u\} \mid V \setminus \{u\})}{\pi_i(\{u\})}.$$

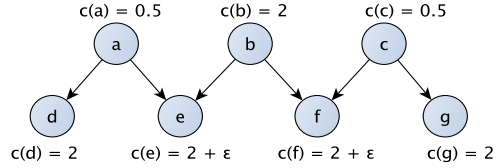
□

We will make use of the following notions in our results on approximation guarantees.

**Definition 5** (Upper and lower rank). *Let  $(\mathcal{E}, \mathcal{C})$  be an independence system. Its upper rank  $R$  and lower rank  $r$  are defined as the cardinalities of the smallest and largest maximal independent sets:*

$$r = \min\{|X| : X \in \mathcal{C} \text{ and } X \cup \{(u, i)\} \notin \mathcal{C}, \forall (u, i) \notin X\},$$

$$R = \max\{|X| : X \in \mathcal{C} \text{ and } X \cup \{(u, i)\} \notin \mathcal{C}, \forall (u, i) \notin X\}.$$



**Figure 1:** Instance illustrating tightness of bound in Theorem 2.

When the independence system is a matroid,  $r = R$ , as all maximal independent sets have the same cardinality.

**Theorem 2.** *CA-GREEDY achieves an approximation guarantee of  $\frac{1}{\kappa_\pi} \left[ 1 - \left( \frac{R - \kappa_\pi}{R} \right)^r \right]$  to the optimum, where  $\kappa_\pi$  is the total curvature of the total revenue function  $\pi(\cdot)$ ,  $r$  and  $R$  are respectively the lower and upper rank of  $(\mathcal{E}, \mathcal{C})$ . This bound is tight.*

*Proof.* We note that the family  $\mathcal{C}$  of subsets that constitute feasible solutions to the RM problem form an independence system defined on  $\mathcal{E}$  (Lemma 2). Given this, the approximation guarantee of CA-GREEDY directly follows from the result of Conforti et al. [15, Theorem 5.4] for submodular function maximization subject to an independence system constraint. However, the tightness does not directly follow from the tightness result in [15], which we address next.

We now exhibit an instance to show that the bound is tight. Consider one advertiser, i.e.,  $h = 1$ . The network is shown in Figure 1, where all influence probabilities are 1. The incentive costs for nodes are as shown in the figure, while  $cpe(\cdot) = 1$ . The budget is  $B = 7$ . It is easy to see that the lower rank is  $r = 1$ , corresponding to the maximal feasible seed set  $S = \{b\}$ , while the upper rank is  $R = 2$ , e.g., corresponding to maximal feasible seed sets such as  $T = \{a, c\}$ . Furthermore, the total curvature is  $\kappa_\pi = 1$ . On this instance, the optimal solution is  $T$  which achieves a revenue of 6. In its first iteration, CA-GREEDY could choose  $b$  as a seed. Once it does, it is forced to the solution  $S = \{b\}$  as no more seeds can be added to  $S$ . The revenue of CA-GREEDY is  $3 = \frac{1}{\kappa_\pi} \left[ 1 - \left( \frac{R - \kappa_\pi}{R} \right)^r \right] OPT = \frac{1}{2} \cdot 6$ . □

**Discussion.** We next discuss the significance and the meaning of the bound in Theorem 2. Notice that when there is just one advertiser, TIC reduces to IC. Even for this simple setting, the bound on CA-GREEDY is tight. By a simple rearrangement of the terms, we have:

$$\frac{1}{\kappa_\pi} \left[ 1 - \left( \frac{R - \kappa_\pi}{R} \right)^r \right] \geq \frac{1}{\kappa_\pi} \left( 1 - e^{-\kappa_\pi \frac{r}{R}} \right).$$

Clearly, the cost-agnostic approximation bound improves as  $\frac{r}{R}$  approaches 1, achieving the best possible value when  $r = R$ . As a special case, the cost-agnostic approximation further improves when the independence system  $(\mathcal{E}, \mathcal{C})$  is a matroid since for a matroid  $r = R$  always holds: e.g., consider the standard IM problem [24] which corresponds to submodular function maximization subject to a uniform matroid. Here,  $\pi(\cdot) = \sigma(\cdot)$ . Then the approximation guarantee becomes  $\frac{1}{\kappa_\pi} (1 - e^{-\kappa_\pi})$ , providing a slight improvement over the usual  $(1 - 1/e)$ -approximation, thanks to the curvature term  $\kappa_\pi$ .<sup>8</sup> This remark is also valid for budgeted influence maximization [26] with uniform seed costs. For more general

<sup>8</sup>Note that  $\kappa_\pi \leq 1$  always. Hence, the extent of improvement increases as the total curvature  $\kappa_\pi$  decreases.

instances of the problem, the guarantee depends on the characteristics of the instance, specifically, the lower and upper ranks and the curvature. This kind of instance dependent bound is characteristic of submodular function maximization over an independence system [15, 25]. Specifically for the RM problem, given its constraints, the values of  $r$  and  $R$  are dictated by the values of  $h$  payment functions over all feasible allocations. For instance, given our assumption that every advertiser can afford at least one seed, we always have  $r \geq h$ . The worst-case value  $r = h$  corresponds to the case in which each advertiser  $i$  is allocated a single seed node  $u_i$  whose payment  $\rho_i(u_i)$  exhausts its budget  $B_i$ . Similarly for  $R$ , without using any particular assumption on  $B_i, \forall i \in [h]$ , we always have  $R \leq \min(n, \sum_{i \in [h]} \lfloor B_i / cpe(i) \rfloor)$ . Notice also that:

$$\frac{1}{\kappa_\pi} \left[ 1 - \left( \frac{R - \kappa_\pi}{R} \right)^r \right] = \frac{1}{\kappa_\pi} \left[ 1 - \left( 1 - \frac{\kappa_\pi}{R} \right)^r \right] \quad (2)$$

$$\geq \frac{1}{\kappa_\pi} \left[ 1 - \left( 1 - \frac{\kappa_\pi}{R} \right) \right] = \frac{1}{\kappa_\pi} \frac{\kappa_\pi}{R} = \frac{1}{R} \quad (3)$$

Hence, the worst-case approximation is always bounded by  $1/R$ .

---

#### Algorithm 1: CA-GREEDY

---

**Input :**  $G = (V, E), B_i, cpe(i), \tilde{\gamma}_i, \forall i \in [h],$   
 $c_i(u), \forall i \in [h], \forall u \in V$   
**Output:**  $\vec{S}_g = (S_1, \dots, S_h)$   
1  $t \leftarrow 1, \mathcal{E}^0 \leftarrow \mathcal{E}, \mathcal{X}_g^0 \leftarrow \emptyset$   
2  $S_i^0 \leftarrow \emptyset, \forall i \in [h]$   
3 **while**  $\mathcal{E}^{t-1} \neq \emptyset$  **do**  
4    $(u^*, i^*) \leftarrow \operatorname{argmax}_{(u,i) \in \mathcal{E}^{t-1}} \pi_i(u \mid S_i^{t-1})$   
5   **if**  $(\mathcal{X}_g^{t-1} \cup \{(u^*, i^*)\}) \in \mathcal{C}$  **then**  
6      $S_{i^*}^t \leftarrow S_{i^*}^{t-1} \cup \{u^*\}$   
7      $S_j^t \leftarrow S_j^{t-1}, \forall j \neq i^*$   
8      $\mathcal{X}_g^t \leftarrow \mathcal{X}_g^{t-1} \cup \{(u^*, i^*)\}$   
9      $\mathcal{E}^t \leftarrow \mathcal{E}^{t-1} \setminus \{(u^*, i^*)\}$   
10     $t \leftarrow t + 1$   
11 **else**  
12     $\mathcal{E}^{t-1} \leftarrow \mathcal{E}^{t-1} \setminus \{(u^*, i^*)\}$   
13  $S_i \leftarrow S_i^{t-1}, \forall i \in [h]$   
14 **return**  $\vec{S}_g = (S_1, \dots, S_h)$

---

### 3.2 Cost-Sensitive Greedy Algorithm

The Cost-sensitive greedy algorithm (CS-GREEDY) for the RM problem is similar to CA-GREEDY. The main difference is that at each iteration  $t$ , CS-GREEDY first finds the (node, advertiser) pair  $(u^*, i^*)$  that maximizes  $\frac{\pi_i(u \mid S_i^{t-1})}{\rho_i(u \mid S_i^{t-1})}$ , and tests whether the addition of this pair to the current greedy solution set  $\mathcal{X}_g^{t-1}$  would violate any matroid or knapsack independence constraint: if the addition is feasible, the pair  $(u^*, i^*)$  is added to the greedy solution as the  $t$ -th (node, advertiser) pair. Otherwise,  $(u^*, i^*)$  is removed from the current ground set  $\mathcal{E}^{t-1}$ . CS-GREEDY terminates when there is no (node, advertiser) pair left in the current ground set  $\mathcal{E}^{t-1}$ . CS-GREEDY can be obtained by simply replacing Line 4 of Algorithm 1 with

$$(u^*, i^*) \leftarrow \operatorname{argmax}_{(u,i) \in \mathcal{E}^{t-1}} \frac{\pi_i(u \mid S_i^{t-1})}{\rho_i(u \mid S_i^{t-1})}.$$

**Theorem 3.** CS-GREEDY achieves an approximation guarantee of

$$1 - \frac{R \cdot \rho_{\max}}{R \cdot \rho_{\max} + (1 - \max_{i \in [h]} \kappa_{\rho_i}) \cdot \rho_{\min}}$$

to the optimum where  $R$  is the upper rank of  $(\mathcal{E}, \mathcal{C})$ ,  $\kappa_{\rho_i}$  is the total curvature of  $\rho_i(\cdot)$ ,  $\forall i \in [h]$ ,  $\rho_{\max} := \max_{(u,i) \in \mathcal{E}} \rho_i(u)$  and  $\rho_{\min} :=$

$\min_{(u,i) \in \mathcal{E}} \rho_i(u)$  are respectively the maximum and minimum singleton payments over all (node, advertiser) pairs.

*Proof.* We use  $\vec{S}^* = (S_1^*, \dots, S_h^*)$  and  $\vec{S}_g = (S_1, \dots, S_h)$  to denote the optimal and greedy allocations respectively, and  $\mathcal{X}^*$  and  $\mathcal{X}_g$  to denote the corresponding solution sets. Specifically,  $S_i^* = \{u : (u, i) \in \mathcal{X}^*\}$ , and  $S_i = \{u : (u, i) \in \mathcal{X}_g\}$ . We denote by  $\mathcal{X}_g^t$  the result of the greedy solution after  $t$  iterations. Let  $K = |\mathcal{X}_g|$  denote the size of the greedy solution. Thus,  $\mathcal{X}_g = \mathcal{X}_g^K$ . By submodularity and monotonicity:

$$\pi(\vec{S}^*) \leq \pi(\vec{S}_g) + \sum_{(u,i) \in \mathcal{X}^* \setminus \mathcal{X}_g} \pi_i(u \mid S_i) \leq \pi(\vec{S}_g) + \sum_{(u,i) \in \mathcal{X}^*} \pi_i(u \mid S_i).$$

At each iteration  $t$ , the greedy algorithm first finds the (node, advertiser) pair  $(u^*, i^*) \leftarrow \operatorname{argmax}_{(u,i) \in \mathcal{E}^{t-1}} \frac{\pi_i(u \mid S_i^{t-1})}{\rho_i(u \mid S_i^{t-1})}$ , and tests whether the addition of this pair to the current greedy solution set  $\mathcal{X}_g^{t-1}$  would violate any independence constraint. If  $(u^*, i^*)$  is feasible, i.e., if  $\mathcal{X}_g^{t-1} \cup \{(u^*, i^*)\} \in \mathcal{C}$ , then the pair  $(u^*, i^*)$  is added to the greedy solution as the  $t$ -th (node, advertiser) pair; otherwise,  $(u^*, i^*)$  is removed from the current ground set  $\mathcal{E}^{t-1}$ . In what follows, for clarity, we use the notation  $(u_t, i_t)$  to denote the (node, advertiser) pair that is *successfully* added by the greedy algorithm to  $\mathcal{X}_g^{t-1}$  in iteration  $t$ .

Let  $U^t$  denote the set of (node, advertiser) pairs that the greedy algorithm *tested* for possible addition to the greedy solution in the first  $(t+1)$  iterations *before* the addition of the  $(t+1)$ -st pair  $(u_{t+1}, i_{t+1})$  into  $\mathcal{X}_g^t$ . Thus,  $U^t \setminus U^{t-1}$  includes the  $t$ -th pair  $(u_t, i_t)$  that was successfully added to  $\mathcal{X}_g^{t-1}$ , as well as all the pairs that were tested for addition into  $\mathcal{X}_g^t$  but failed the independence test. Thus,  $\forall (u, i) \in U^t \setminus U^{t-1}$ , we have  $\frac{\pi_i(u \mid S_i^t)}{\rho_i(u \mid S_i^t)} \geq$

$$\frac{\pi_{i_{t+1}}(u_{t+1} \mid S_{i_{t+1}}^t)}{\rho_{i_{t+1}}(u_{t+1} \mid S_{i_{t+1}}^t)}, \text{ since they were tested for addition to } \mathcal{X}_g^t \text{ before } (u_{t+1}, i_{t+1}), \text{ but failed the independence test. For all } (u, i) \in U^t \setminus U^{t-1}, \text{ we have } \frac{\pi_i(u \mid S_i^{t-1})}{\rho_i(u \mid S_i^{t-1})} \leq \frac{\pi_{i_t}(u_t \mid S_{i_t}^{t-1})}{\rho_{i_t}(u_t \mid S_{i_t}^{t-1})}.$$

since they were not good enough to be added to  $\mathcal{X}_g^{t-1}$  as the  $t$ -th pair. Note that, the greedy algorithm terminates when there is no feasible pair left in the ground set. Hence after  $K$  iterations,  $\mathcal{E}^K$  contains only the *infeasible* pairs that violate some matroid or knapsack constraint. Thus, we have  $\mathcal{X}^* = \bigcup_{t=1}^K [\mathcal{X}^* \cap (U^t \setminus U^{t-1})]$ . Let  $U_t^* := \mathcal{X}^* \cap (U^t \setminus U^{t-1})$ . Notice that  $\mathcal{X}^* = \bigcup_{t=1}^K U_t^*$ . Then, we have:

$$\begin{aligned} \pi(\vec{S}^*) &\leq \pi(\vec{S}_g) + \sum_{(u,i) \in \mathcal{X}^*} \pi_i(u \mid S_i) \\ &= \pi(\vec{S}_g) + \sum_{t=1}^K \sum_{(u,i) \in U_t^*} \pi_i(u \mid S_i) \\ &\leq \pi(\vec{S}_g) + \sum_{t=1}^K \sum_{(u,i) \in U_t^*} \frac{\pi_{i_t}(u_t \mid S_{i_t}^{t-1})}{\rho_{i_t}(u_t \mid S_{i_t}^{t-1})} \cdot \rho_i(u \mid S_i^{t-1}). \end{aligned}$$

The last inequality is due to the fact that  $\forall(u, i) \in \mathcal{U}_t^*$ :

$$\pi_i(u | S_i) \leq \pi_i(u | S_i^{t-1}) \leq \frac{\pi_{i_t}(u_t | S_{i_t}^{t-1})}{\rho_{i_t}(u_t | S_{i_t}^{t-1})} \cdot \rho_i(u | S_i^{t-1}),$$

where the first inequality follows from submodularity and the second follows from the greedy choice of (node, advertiser) pairs. Continuing, we have:

$$\begin{aligned} \pi(\vec{S}^*) &\leq \pi(\vec{S}_g) + \sum_{t=1}^K \sum_{(u,i) \in \mathcal{U}_t^*} \frac{\pi_{i_t}(u_t | S_{i_t}^{t-1})}{\rho_{i_t}(u_t | S_{i_t}^{t-1})} \cdot \rho_i(u | S_i^{t-1}) \\ &= \pi(\vec{S}_g) + \sum_{t=1}^K \frac{\pi_{i_t}(u_t | S_{i_t}^{t-1})}{\rho_{i_t}(u_t | S_{i_t}^{t-1})} \sum_{(u,i) \in \mathcal{U}_t^*} \rho_i(u | S_i^{t-1}) \\ &\leq \pi(\vec{S}_g) + \sum_{t=1}^K \frac{\pi_{i_t}(u_t | S_{i_t}^{t-1})}{\rho_{i_t}(u_t | S_{i_t}^{t-1})} \cdot \sum_{t=1}^K \sum_{(u,i) \in \mathcal{U}_t^*} \rho_i(u) \\ &= \pi(\vec{S}_g) + \sum_{t=1}^K \frac{\pi_{i_t}(u_t | S_{i_t}^{t-1})}{\rho_{i_t}(u_t | S_{i_t}^{t-1})} \cdot \sum_{(u,i) \in \mathcal{X}^*} \rho_i(u) \\ &\leq \pi(\vec{S}_g) + \pi(\vec{S}_g) \cdot \frac{R \cdot \max_{(u,i) \in \mathcal{X}^*} \rho_i(u)}{\min_{t \in [1, K]} \rho_{i_t}(u_t | S_{i_t}^{t-1})} \end{aligned} \quad (4)$$

where the last inequality follows from the fact that  $\pi(\vec{S}_g) = \sum_{t=1}^K \pi_{i_t}(u_t | S_{i_t}^{t-1})$  and  $|\mathcal{X}^*| \leq R$  since  $\mathcal{X}^* \in \mathcal{C}$ . Let  $(u_{t_m}, i_{t_m}) := \argmin_{t \in [1, K]} \rho_{i_t}(u_t | S_{i_t}^{t-1})$  and let  $(u_{min}, i_{min}) := \argmin_{(u,i) \in \mathcal{E}} \rho_i(u | V \setminus \{u\})$ . Being monotone and submodular, each

$\rho_i(\cdot)$  has the total curvature  $\kappa_{\rho_i} = 1 - \min_{u \in V} \frac{\rho_i(u | V \setminus \{u\})}{\rho_i(u)}$ .

Hence, for  $\rho_{i_{min}}(\cdot)$ , we have:

$$1 - \kappa_{\rho_{i_{min}}} = \min_{u \in V} \frac{\rho_{i_{min}}(u | V \setminus \{u\})}{\rho_{i_{min}}(u)} \leq \frac{\rho_{i_{min}}(u_{min} | V \setminus \{u_{min}\})}{\rho_{i_{min}}(u_{min})} \quad (5)$$

where the inequality above follows from the definition of total curvature. Then, using submodularity and Eq.5, we obtain:

$$\begin{aligned} \min_{t \in [1, K]} \rho_{i_t}(u_t | S_{i_t}^{t-1}) &= \rho_{i_{t_m}}(u_{t_m} | S_{i_{t_m}}^{t_m-1}) \\ &\geq \rho_{i_{t_m}}(u_{t_m} | V \setminus \{u_{t_m}\}) \\ &\geq \min_{(u,i) \in \mathcal{E}} \rho_i(u | V \setminus \{u\}) \\ &= \rho_{i_{min}}(u_{min} | V \setminus \{u_{min}\}) \\ &\geq (1 - \kappa_{\rho_{i_{min}}}) \cdot \rho_{i_{min}}(u_{min}) \\ &\geq (1 - \max_{i \in [h]} \kappa_{\rho_i}) \cdot \min_{(u,i) \in \mathcal{E}} \rho_i(u). \end{aligned} \quad (6)$$

Continuing from where we left in Eq.4 and using Eq.6, we have:

$$\begin{aligned} \pi(\vec{S}^*) &\leq \pi(\vec{S}_g) + \pi(\vec{S}_g) \cdot \frac{R \cdot \max_{(u,i) \in \mathcal{X}^*} \rho_i(u)}{\min_{t \in [1, K]} \rho_{i_t}(u_t | S_{i_t}^{t-1})} \\ &\leq \pi(\vec{S}_g) \cdot \left( 1 + \frac{R \cdot \max_{(u,i) \in \mathcal{E}} \rho_i(u)}{(1 - \max_{i \in [h]} \kappa_{\rho_i}) \cdot \min_{(u,i) \in \mathcal{E}} \rho_i(u)} \right) \\ &= \pi(\vec{S}_g) \cdot \left( 1 + \frac{R \cdot \rho_{max}}{(1 - \max_{i \in [h]} \kappa_{\rho_i}) \cdot \rho_{min}} \right) \end{aligned} \quad (7)$$

Rearranging the terms we obtain:

$$\begin{aligned} \pi(\vec{S}_g) &\geq \pi(\vec{S}^*) \cdot \frac{(1 - \max_{i \in [h]} \kappa_{\rho_i}) \cdot \rho_{min}}{(1 - \max_{i \in [h]} \kappa_{\rho_i}) \cdot \rho_{min} + R \cdot \rho_{max}} \\ &= \pi(\vec{S}^*) \cdot \left( 1 - \frac{R \cdot \rho_{max}}{R \cdot \rho_{max} + (1 - \max_{i \in [h]} \kappa_{\rho_i}) \cdot \rho_{min}} \right). \end{aligned}$$

□

**Discussion.** We next discuss the significance and the meaning of the bounds. Notice that the value of the cost-sensitive approximation bound improves as the ratio  $\frac{\rho_{max}}{\rho_{min}}$  decreases, as Eq. 7 shows. Since  $\rho_{max} \leq \min_{i \in [h]} B_i$ , we can see that as the value of  $\rho_{max}$  decreases, intuitively  $r$  would increase, for the corresponding maximal independent set of minimum size could pack more seeds under the knapsack constraints. Similarly, if the value of  $\rho_{min}$  increases,  $R$  would decrease since the corresponding maximal independent set of maximum size could pack fewer seeds under the knapsack constraints. Thus, intuitively as  $\frac{\rho_{max}}{\rho_{min}}$  decreases,  $\frac{r}{R}$  would increase. When this happens, both cost-agnostic and cost-sensitive approximations improve.

At one extreme, when  $\kappa_{\rho_i} = 0, \forall i \in [h]$ , i.e., when  $\rho_i(\cdot)$  is modular  $\forall i \in [h]$ , we have linear knapsack constraints. Thus, Theorem 2 and Theorem 3 respectively provide cost-agnostic and cost-sensitive approximation guarantees for the *Budgeted Influence Maximization* problem [26, 31] for the case of multiple advertisers, with an additional matroid constraint. At the other extreme, when  $\max_{i \in [h]} \kappa_{\rho_i} = 1$ , which is the case for totally normalized and saturated functions (e.g., matroid rank functions), the approximation guarantee of CS-GREEDY is unbounded, i.e., it becomes degenerate. This is similar to the result of [22] for the SCSK problem whose cost-sensitive approximation guarantee becomes unbounded. Nevertheless, combining the results of the cost-agnostic and cost-sensitive cases, we can obtain a bounded approximation.

On the other hand, while CA-GREEDY always has a bounded worst-case guarantee, our experiments show that CS-GREEDY empirically obtains higher revenue<sup>9</sup>.

## 4. SCALABLE ALGORITHMS

While Algorithms CA-GREEDY and CS-GREEDY provide approximation guarantees, their efficient implementation is a challenge, as both of them require a large number of influence spread computations: in each iteration  $t$ , for each advertiser  $i$  and each node  $u \in V \setminus S_i^{t-1}$ , the algorithms need to compute  $\pi_i(u | S_i^{t-1})$  and  $\pi_i(u | S_i^{t-1})/\rho_i(u | S_i^{t-1})$ , respectively.

Computing the exact influence spread  $\sigma(S)$  of a given seed set  $S$  under the IC model is #P-hard [13], and this hardness carries over to the TIC model. In recent years, significant advances have been made in efficiently estimating  $\sigma(S)$ . A natural question is whether they can be adapted to our setting, an issue we address next.

### 4.1 Scalable Influence Spread Estimation

Tang et al. [34] proposed a near-linear time randomized algorithm for influence maximization, called *Two-phase Influence Maximization (TIM)*, building on the notion of “reverse-reachable”

<sup>9</sup>It remains open whether the approximation bound for CS-GREEDY is tight. Interestingly, on the instance (Fig. 1) used in the proof of Theorem2, CS-GREEDY obtains the optimal solution  $T = \{a, c\}$ .

(RR) sets proposed by Borgs et al. [10]. Random RR-sets are critical in the efficient estimation of influence spread. Tang et al. [33] subsequently proposed an algorithm called IMM that improves upon TIM by tightening the lower bound on the number of random RR-sets required to estimate influence with high probability. The difference between TIM and IMM is that the lower bound used by TIM ensures that the number of random RR-sets it uses is sufficient to estimate the spread of *any* seed set of a given size  $s$ . By contrast, IMM uses a lower bound that is tailored for the seed that is greedily selected by the algorithm. Nguyen et al. [32], adapting ideas from TIM [34], and the sequential sampling design proposed by Dagum *et al.* [16], proposed an algorithm called SSA that provides significant run-time improvement over TIM and IMM.

These algorithms are designed for the basic influence maximization problem and hence require knowing the number of seeds as input. In our problem, the number of seeds is not fixed, but is dynamic and depends on the budget and partition matroid constraints. Thus a direct application of these algorithms is not possible.

Aslay et al. [4] recently proposed a technique for efficient seed selection for IM when the number of seeds required is not predetermined but can change dynamically. However, their technique cannot handle the presence of seed user incentives which, in our setting, directly affects the number of seeds required to solve the RM problem. In this section, we derive inspiration from their technique. First, though we note that for CA-GREEDY, in each iteration, for each advertiser, we need to find a feasible node that yields the maximum marginal gain in revenue, and hence the maximum marginal spread. By contrast, in CS-GREEDY, we need to find the node that yields the maximum *rate* of marginal revenue per marginal gain in payment, i.e.,  $\pi_i(u | S_i^{t-1})/\rho_i(u | S_i^{t-1})$ .

To find such node  $u_i^t$  we must compute  $\sigma_i(v | S_i^{t-1}), \forall v : (v, i) \in \mathcal{E}^{t-1}$ : notice that node  $u_i^t$  might even correspond to the node that has the *minimum* marginal gain in influence spread for iteration  $t$ . Thus, any scalable realization of CS-GREEDY should be capable of working as an influence spread oracle that can efficiently compute  $\pi_i(u | S_i^{t-1})/\rho_i(u | S_i^{t-1})$  for all  $u \in \{v : (v, i) \in \mathcal{E}^{t-1}\}$ .

Among the state-of-the-art IM algorithms [32–34], only TIM [34] can be adapted to serve as an influence oracle. For a given set size  $s$ , the derivation of the number of random RR-sets that TIM uses is done such that the influence spread of *any set of at most  $s$  nodes can be accurately estimated*. On the other hand, even though IMM [33] and SSA [32] provide significant run-time improvements over TIM, they inherently cannot perform this estimation task accurately: the sizes of the random RR-sets sample that these algorithms use are tuned just for accurately estimating the influence spread of *only* the approximate greedy solutions; the sample sizes used are inadequate for estimating the spread of arbitrary seed sets of a given size. Thus, we choose to extend TIM to devise scalable realizations of CA-GREEDY and CS-GREEDY, namely, TI-CARM and TI-CSR. Next, we describe how to extend the ideas of RR-sets sampling and TIM’s sample size determination technique to obtain scalable approximation algorithms for the RM problem: TI-CARM and TI-CSR.

## 4.2 Scalable Revenue Maximization

For the scalable estimation of influence spread, in this section we devise TI-CARM and TI-CSR, scalable realizations of CA-GREEDY and CS-GREEDY, based on the notion of Reverse-Reachable sets [10] and adapt the sample size determination procedure employed by TIM [34] to achieve a certain estimation accuracy with high confidence.

**Reverse-Reachable (RR) sets [10].** Under the IC model, a random RR-set  $R$  from  $G$  is generated as follows. First, for every edge

$(u, v) \in E$ , remove it from  $G$  w.p.  $1 - p_{u,v}$ : this generates a possible world (deterministic graph)  $X$ . Second, pick a *target* node  $w$  uniformly at random from  $V$ . Then,  $R$  consists of the nodes that can reach  $w$  in  $X$ . For a sufficient sample  $\mathbf{R}$  of random RR-sets, the fraction  $F_{\mathbf{R}}(S)$  of  $\mathbf{R}$  covered by  $S$  is an unbiased estimator of  $\sigma(S)$ , i.e.,  $\sigma(S) = \mathbb{E}[n \cdot F_{\mathbf{R}}(S)]$ .

**Sample Size Determination of TIM [34].** Let  $\mathbf{R}_i$  be a collection of  $\theta_i$  random RR-sets. Given any seed set size  $s_i$  and  $\varepsilon > 0$ , define  $L_i(s_i, \varepsilon)$  to be:

$$L_i(s_i, \varepsilon) = (8 + 2\varepsilon)n \cdot \frac{\ell \log n + \log \binom{n}{s_i} + \log 2}{OPT_{i,s_i} \cdot \varepsilon^2}, \quad (8)$$

where  $\ell > 0, \varepsilon > 0$  and  $OPT_{i,s_i} = \max_{S \subseteq V, |S| \leq s_i} \sigma_i(S)$ . Let  $\theta_i$  be a number no less than  $L_i(s_i, \varepsilon)$ . Then, for any seed set  $S$  with  $|S| \leq s_i$ , the following inequality holds w.p. at least  $1 - n^{-\ell} / \binom{n}{s_i}$ :

$$|n \cdot F_{\mathbf{R}_i}(S_i) - \sigma_i(S_i)| < \frac{\varepsilon}{2} \cdot OPT_{i,s_i}. \quad (9)$$

**Estimated Payments and Budget Feasibility.**<sup>10</sup> Let  $\tilde{S} = (\tilde{S}_1, \dots, \tilde{S}_h)$  denote the approximately greedy solution that TI-CARM (resp. TI-CSR) returns. Since the algorithm operates on the estimation of influence spread, the revenue and payment computed for each advertiser  $i$  will also be estimations of the actual revenue and payment for seed set  $\tilde{S}_i$ . Let  $\tilde{\pi}_i(\tilde{S}_i) = cpe(i) \cdot n \cdot F_{\mathbf{R}_i}(\tilde{S}_i)$  and  $\tilde{\rho}_i(\tilde{S}_i) = c_i(\tilde{S}_i) + \tilde{\pi}_i(\tilde{S}_i)$  denote the estimated revenue and estimated payment for advertiser  $i$ , respectively. As TI-CARM (resp. TI-CSR) performs budget feasibility check on the estimated payments, it is possible to encounter scenarios in which  $\tilde{\rho}_i(\tilde{S}_i) \leq B_i$  while  $\rho_i(\tilde{S}_i) > B_i$ . Thus, to ensure that the approximate greedy allocation results in actual payments that do not violate any budget constraints with high probability, one could consider to use a refined budget  $\tilde{B}_i < B_i$ , for each advertiser  $i$ , by taking into account the error introduced by spread estimation. Next, we provide details on how to set  $\tilde{B}_i$  so that  $\tilde{S}_i$  is budget feasible with high probability.

First, notice that, following Eq.9, we have  $\sigma_i(\tilde{S}_i) \leq n \cdot F_{\mathbf{R}_i}(\tilde{S}_i) + \frac{\varepsilon}{2} \cdot OPT_{i,s_i}$ . Thus, to ensure that  $c_i(\tilde{S}_i) + cpe(i) \cdot \sigma_i(\tilde{S}_i) \leq B_i$ , w.h.p., we need to have:

$$c_i(\tilde{S}_i) + cpe(i) \cdot \left( n \cdot F_{\mathbf{R}_i}(\tilde{S}_i) + \frac{\varepsilon}{2} \cdot OPT_{i,s_i} \right) \leq B_i$$

which implies that the budget constraint on the estimated payment  $\tilde{\rho}_i(\tilde{S}_i)$  should be refined as:

$$\tilde{\rho}_i(\tilde{S}_i) \leq B_i - cpe(i) \cdot \frac{\varepsilon}{2} \cdot OPT_{i,s_i}. \quad (10)$$

While using a refined budget of  $B_i - cpe(i) \cdot \frac{\varepsilon}{2} \cdot OPT_{i,s_i}$  would ensure w.h.p. that  $\rho_i(\tilde{S}_i) \leq B_i$ , such refinement requires to compute  $OPT_{i,s_i}$  which is unknown and NP-hard to compute. To circumvent this difficulty, one could consider an upper bound  $\eta_{i,s_i}$  on  $OPT_{i,s_i}$  so that

$$\begin{aligned} \tilde{B}_i &= B_i - cpe(i) \cdot \frac{\varepsilon}{2} \cdot \eta_{i,s_i} \\ &\leq B_i - cpe(i) \cdot \frac{\varepsilon}{2} \cdot OPT_{i,s_i}. \end{aligned}$$

Following [37], an upper bound  $\eta_{i,s_i}$  on  $OPT_{i,s_i}$  can be obtained as follows.

<sup>10</sup>We would like to thank to Kai Han and Jing Tang for bringing the budget feasibility issue into our attention, which we now address in this section.



**Lemma 3** (Restated from Lemma 4.3 [37]). Let  $\mathbf{R}_i$  be a sample of  $\theta_i$  RR-sets, such that,  $\theta_i \geq L_i(s_i, \varepsilon)$ , and let  $\tilde{A}_i \subseteq V$ ,  $|\tilde{A}_i| = s_i$  denote the greedy solution to maximum coverage problem on the sample  $\mathbf{R}_i$ . Define  $\eta_{i,s_i}$  to be:

$$\eta_{i,s_i} := \left( \sqrt{\frac{\theta_i \cdot F_{\mathbf{R}_i}(\tilde{A}_i)}{1 - 1/e}} + \frac{\ln n^\ell}{2} + \sqrt{\frac{\ln n^\ell}{2}} \right)^2 \cdot \frac{n}{\theta_i} \quad (11)$$

Then, we have:

$$\Pr[OPT_{i,s_i} \leq \eta_{i,s_i}] \geq 1 - n^{-\ell}.$$

Following Lemma 3, for a given seed set size  $s_i$ , we can define  $\tilde{B}_i$  for  $i$  as:

$$\tilde{B}_i = B_i - cpe(i) \cdot \frac{\varepsilon}{2} \cdot \eta_{i,s_i}. \quad (12)$$

**Latent Seed Set Size Estimation.** The derivation of the sufficient sample size, depicted in Eq. 8, requires the number of seeds as input for each  $i$ , which is not available for RM problem. Let  $s_i^* = |S_i^*|$  denote the true number of seeds that the optimal allocation would assign to  $i$ . From the advertisers' budgets, there is no obvious way to determine  $s_i^*$  for each  $i$ . This poses a challenge as the required number of RR-sets ( $\theta_i$ ) for advertiser  $i$  depends on  $s_i^*$ .

To circumvent this difficulty, one can use a safe upper bound  $\bar{s}_i = \lceil \frac{B_i}{\rho_{min}^i} \rceil$  on  $s_i^*$ , where  $\rho_{min}^i$  is the minimum singleton payment for  $i$  so that, by using a sample of at least  $L_i(\bar{s}_i, \varepsilon)$  RR-sets, we can quantify how the approximation guarantee of TI-CARM (resp, TI-CSRM) deteriorate from the guarantee of CA-GREEDY (resp., CS-GREEDY) as a function of the estimation accuracy that the sample size ensures for all seed sets of size at most  $\bar{s}_i$  (Eq.9). However, when  $\rho_{min}$  is very small w.r.t.  $B_i$ , a direct application of TIM's sample size derivation technique for  $\bar{s}_i$  seeds could result in a large estimation error  $\frac{\varepsilon}{2} \cdot OPT_{i,\bar{s}_i}$ , due to  $\bar{s}_i$  being a very loose upper bound on  $s_i^*$ . Such large estimation error could translate to working with a refined budget  $\tilde{B}_i$  that is very small w.r.t.  $B_i$ , resulting in greatly under-utilizing the budget for the sake of budget feasibility. Now, we explain how to derive a sample size that can estimate the spread of any seed set of size at most  $\bar{s}_i$  while using a more stringent estimation error  $\frac{\varepsilon}{2} \cdot OPT_{i,\tilde{s}_i}$  with  $\tilde{s}_i < \bar{s}_i$ , where  $\tilde{s}_i$  is the latent seed set size estimation obtained during the execution of TI-CARM (resp., TI-CSRM) as we will explain next.

**Lemma 4.** Let  $\mathbf{R}_i$  be a collection of  $\theta_i$  random RR-sets. Given  $\bar{s}_i$ ,  $\tilde{s}_i$ , and  $\varepsilon > 0$ , define  $L_i(\bar{s}_i, \tilde{s}_i, \varepsilon)$  to be:

$$L_i(\bar{s}_i, \tilde{s}_i, \varepsilon) = (8\lambda + 2\varepsilon)n \cdot \frac{\ell \log n + \log \binom{n}{\tilde{s}_i} + \log 2}{OPT_{i,\tilde{s}_i} \cdot \varepsilon^2}, \quad (13)$$

where  $\ell > 0, \varepsilon > 0$ ,  $OPT_{i,s} = \max_{S \subseteq V, |S| \leq s} \sigma_i(S)$ , for any integer  $s$ , and  $\lambda = \frac{OPT_{i,\bar{s}_i}}{OPT_{i,\tilde{s}_i}}$ . Let  $\theta_i$  be a number no less than  $L_i(\bar{s}_i, \tilde{s}_i, \varepsilon)$ . Then, for any seed set  $S$  with  $|S| \leq \bar{s}_i$ , the following inequality holds w.p. at least  $1 - n^{-\ell}/\binom{n}{\tilde{s}_i}$ :

$$|n \cdot F_{\mathbf{R}_i}(S) - \sigma_i(S)| < \frac{\varepsilon}{2} \cdot OPT_{i,\tilde{s}_i}. \quad (14)$$

*Proof.* Let  $S$  be any seed set of size at most  $\bar{s}_i$  and let  $\tau_i$  denote the probability that  $S$  overlaps with a random RR set, i.e.,

$$\tau_i = \mathbb{E}[F_{\mathbf{R}_i}(S)] = \frac{\sigma_i(S)}{n}.$$

Then, we have:

$$\begin{aligned} \Pr[|n \cdot F_{\mathbf{R}_i}(S) - \sigma_i(S)| < \frac{\varepsilon}{2} \cdot OPT_{i,\tilde{s}_i}] \\ &= \Pr\left[|\theta_i \cdot F_{\mathbf{R}_i}(S) - \tau_i \theta_i| < \frac{\varepsilon \theta_i}{2n} \cdot OPT_{i,\tilde{s}_i}\right] \\ &= \Pr\left[|\theta_i \cdot F_{\mathbf{R}_i}(S) - \tau_i \theta_i| < \frac{\varepsilon \cdot OPT_{i,\tilde{s}_i}}{2n\tau_i} \cdot \tau_i \theta_i\right]. \end{aligned} \quad (15)$$

Letting  $\delta = \frac{\varepsilon \cdot OPT_{i,\tilde{s}_i}}{2n\tau_i}$ , by Chernoff bounds, we have:

$$\begin{aligned} \text{r.h.s. of Eq.15} &< 2 \exp\left(-\frac{\delta^2}{2 + \delta} \cdot \tau_i \theta_i\right) \\ &= 2 \exp\left(-\frac{\varepsilon^2 \cdot OPT_{i,\tilde{s}_i}^2}{8n^2\tau_i + 2\varepsilon n \cdot OPT_{i,\tilde{s}_i}} \cdot \theta_i\right) \\ &< 2 \exp\left(-\frac{\varepsilon^2 \cdot OPT_{i,\tilde{s}_i}^2}{8n \cdot OPT_{i,\tilde{s}_i} + 2\varepsilon n \cdot OPT_{i,\tilde{s}_i}} \cdot \theta_i\right) \\ &= 2 \exp\left(-\frac{\varepsilon^2 \cdot OPT_{i,\tilde{s}_i}}{8n \cdot \frac{OPT_{i,\tilde{s}_i}}{OPT_{i,\tilde{s}_i}} + 2\varepsilon n} \cdot \theta_i\right) \end{aligned}$$

where the last inequality follows from the fact that  $\tau_i \leq OPT_{i,\tilde{s}_i}$ . Finally, we obtain the lower bound on  $\theta_i$  by solving

$$2 \exp\left(-\frac{\varepsilon^2 \cdot OPT_{i,\tilde{s}_i}}{8n \cdot \frac{OPT_{i,\tilde{s}_i}}{OPT_{i,\tilde{s}_i}} + 2\varepsilon n} \cdot \theta_i\right) \leq \frac{n^{-\ell}}{\binom{n}{\tilde{s}_i}}.$$

□

An upper bound on the  $\lambda$  term required for the sample size derivation in Eq. 13 can be obtained by using an upper bound on  $OPT_{i,\bar{s}_i}$ , as given by Lemma 3, and a lower bound on  $OPT_{i,\tilde{s}_i}$  by using the lower bounding technique provided in [34] for TIM's sample size derivation (Eq. 8).

We now explain the “latent seed set size estimation” procedure which first makes an initial guess at the true number of seeds required to maximize cost-agnostic (cost-sensitive) revenue and then iteratively revises the estimated value, until no more seeds are needed, while concurrently selecting seeds and allocating them to advertisers. For ease of exposition, let us first consider a single advertiser  $i$ . We start with an initial estimate, denoted by  $\tilde{s}_i^1$ , and use it to obtain a corresponding sample size  $\theta_i^1 = L_i(\bar{s}_i, \tilde{s}_i^1, \varepsilon)$  using Eq. 8, an upper bound  $\eta_{i,\tilde{s}_i^1}$  using Eq. 11, and a refined budget  $\tilde{B}_i^1$  using Eq. 12. As it is #P-hard to compute  $\rho_{min}^i$ , we also compute in this iteration a safe upper bound  $\bar{s}_i$  from

$$\bar{s}_i = \left\lceil \frac{B_i}{\tilde{\rho}_{min}^i + cpe(i) \cdot \frac{\varepsilon}{2} \cdot \eta_{i,\tilde{s}_i}} \right\rceil$$

where  $\tilde{\rho}_{min}^i = \min_{u \in V} c_i(u) + cpe(i) \cdot n \cdot F_{\mathbf{R}_i}(u)$ . At iteration  $t > 1$ , we compute the sample size from  $\theta_i^t = L_i(\bar{s}_i, \tilde{s}_i^t, \varepsilon)$ , and if  $\theta_i^t > \theta_i^{t-1}$ , we will need to sample additional  $(\theta_i^t - \theta_i^{t-1})$  RR-sets, and use all RR-sets sampled up to this iteration to select  $(\tilde{s}_i^t - \tilde{s}_i^{t-1})$  additional seeds into the seed set  $\tilde{S}_i$  of advertiser  $i$ , while revising the upper bound  $\eta_{i,\tilde{s}_i^t}$  and the corresponding refined budget  $\tilde{B}_i^t$ . After adding those seeds, if the current payment estimate  $\tilde{\rho}_i(\tilde{S}_i)$  is

**Algorithm 2: TI-CSRSM**

**Input :**  $G = (V, E), B_i, cpe(i), \tilde{\gamma}_i, \forall i \in [h], c_i(u), \forall i \in [h], \forall u \in V$

**Output:**  $\tilde{S} = (\tilde{S}_1, \dots, \tilde{S}_h)$

```

1 foreach  $j = 1, 2, \dots, h$  do
2    $\tilde{S}_j \leftarrow \emptyset; Q_j \leftarrow \emptyset$ ; // a priority queue
3    $\tilde{s}_j \leftarrow 1; \theta_j \leftarrow L_j(\tilde{s}_j, \varepsilon); \mathbf{R}_j \leftarrow \text{Sample}(G, \gamma_j, \theta_j)$ ;
4    $\bar{s}_j \leftarrow \left\lfloor \frac{B_j}{\tilde{\rho}_{min}^{max} + cpe(j) \cdot \frac{\varepsilon}{2} \cdot \eta_j, \tilde{s}_j} \right\rfloor$ ;
5    $\tilde{B}_j \leftarrow B_j - cpe(i) \cdot \frac{\varepsilon}{2} \cdot \eta_j, \tilde{s}_j$ ;
6    $\text{assigned}[u] \leftarrow \text{false}, \forall u \in V$ ;
7 while true do
8   foreach  $j = 1, 2, \dots, h$  do
9      $(v_j, cov_j(v_j)) \leftarrow \text{SelectBestCSNode}(\mathbf{R}_j)$  (Alg 5)
10     $F_{\mathbf{R}_j}(v_j) \leftarrow cov_j(v_j)/\theta_j$ ;
11     $\tilde{\pi}_j(\tilde{S}_j \cup \{v_j\}) \leftarrow \tilde{\pi}_j(\tilde{S}_j) + cpe(j) \cdot n \cdot F_{\mathbf{R}_j}(v_j)$ ;
12     $i \leftarrow \arg\max_{j=1}^h \frac{\tilde{\pi}_j(v_j|\tilde{S}_j)}{\tilde{\rho}_j(v_j|\tilde{S}_j)}$  subject to:
13     $\tilde{\rho}_j(\tilde{S}_j \cup \{v_j\}) \leq \tilde{B}_j \wedge \text{assigned}[v_j] = \text{false}$ ;
14    if  $i \neq \text{NULL}$  then
15       $\tilde{S}_i \leftarrow \tilde{S}_i \cup \{v_i\}$ ;
16       $\text{assigned}[v_i] = \text{true}$ ;
17       $Q_i.\text{insert}(v_i, cov_i(v_i))$ ;
18       $\mathbf{R}_i \leftarrow \mathbf{R}_i \setminus \{R \mid v_i \in R \wedge R \in \mathbf{R}_i\}$ ;
19      //remove RR-sets that are covered;
20      else return //all advertisers exhausted;;
21      if  $|\tilde{S}_i| = \tilde{s}_i$  then
22         $\tilde{s}_i \leftarrow \tilde{s}_i + \left\lfloor \frac{\tilde{B}_i - \tilde{\rho}_i(\tilde{S}_i)}{c_i^{max} + cpe(i) \cdot (n \cdot F_{\mathbf{R}_i}^{max} + \frac{\varepsilon}{2} \cdot \eta_i, \tilde{s}_i)} \right\rfloor$ ;
23         $\mathbf{R}_i \leftarrow \mathbf{R}_i \cup \text{Sample}(G, \gamma_i, \max\{0, L_i(\tilde{s}_i, \varepsilon) - \theta_i\})$ ;
24         $\theta_i \leftarrow \max\{L_i(\tilde{s}_i, \varepsilon), \theta_i\}$ ;
25         $\tilde{\pi}_i(\tilde{S}_i) \leftarrow \text{UpdateEstimates}(\mathbf{R}_i, \theta_i, \tilde{S}_i, Q_i)$ ;
26         $\tilde{B}_i \leftarrow B_i - cpe(i) \cdot \frac{\varepsilon}{2} \cdot \eta_i, \tilde{s}_i$ ;
27        //revise estimates to reflect newly added RR-sets;
28         $\tilde{\rho}_i(\tilde{S}_i) \leftarrow \tilde{\pi}_i(\tilde{S}_i) + c_i(\tilde{S}_i)$ ;

```

**Algorithm 3: UpdateEstimates( $\mathbf{R}_i, \theta_i, \tilde{S}_i, Q_i$ )**

**Output:**  $\tilde{\pi}_i(\tilde{S}_i)$

```

1  $\tilde{\pi}_i(\tilde{S}_i) \leftarrow 0$ ;
2 for  $j = 0, \dots, |\tilde{S}_i| - 1$  do
3    $(v, cov_i(v)) \leftarrow Q_i[j]$ ;
4    $cov'_i(v) \leftarrow |\{R \mid v \in R, R \in \mathbf{R}_i\}|$ ;
5    $Q_i.\text{insert}(v, cov_i(v) + cov'_i(v))$ ;
6    $\tilde{\pi}_i(\tilde{S}_i) \leftarrow cpe(i) \cdot n \cdot ((cov_i(v) + cov'_i(v))/\theta_i)$ ; //update coverage of existing seeds w.r.t. new RR-sets added to collection.

```

still less than  $\tilde{B}_i^t$ , more seeds can be assigned to advertiser  $i$ . Thus, we will need another iteration and we further revise our estimation of  $\tilde{s}_i^*$ . The new value,  $\tilde{s}_i^{t+1}$ , is obtained as follows:

$$\tilde{s}_i^{t+1} \leftarrow \tilde{s}_i^t + \left\lfloor \frac{\tilde{B}_i^t - \tilde{\rho}_i(\tilde{S}_i)}{c_i^{max} + cpe(i) \cdot (n \cdot F_{\mathbf{R}_i}^{max} + \frac{\varepsilon}{2} \cdot \eta_i, \tilde{s}_i^t)} \right\rfloor \quad (16)$$

where  $c_i^{max} := \max_{v \in V} c_i(v)$  is the maximum seed user incentive cost for advertiser  $i$ , and  $F_{\mathbf{R}_i}^{max} := \max_{u \in V \setminus \tilde{S}_i} F_{\mathbf{R}_i}(u)$ . This ensures we do not overestimate as future seeds have diminishing marginal gains, thanks to submodularity, and incentives bounded by  $c_i^{max}$ .

While the core logic of TI-CSRSM (resp. TI-CARM) is still based on the greedy seed selection outlined for CS-GREEDY (resp. CA-GREEDY), TI-CSRSM (resp. TI-CARM) uses random RR-sets samples for the scalable estimation of influence spread. Since

**Algorithm 4: SelectBestCANode( $\mathbf{R}_j$ )**

**Output:**  $(u, cov_j(u))$

```

1  $u \leftarrow \arg\max_{v \in V} |\{R \mid v \in R \wedge R \in \mathbf{R}_j\}|$ 
   subject to:  $\text{assigned}[v] = \text{false}$ ;
2  $cov_j(u) \leftarrow |\{R \mid u \in R \wedge R \in \mathbf{R}_j\}|$ ; //find best cost-agnostic seed for ad  $j$  as well as its coverage.

```

**Algorithm 5: SelectBestCSNode( $\mathbf{R}_j$ )**

**Output:**  $(u, cov_j(u))$

```

1  $u \leftarrow \arg\max_{v \in V} \frac{|\{R \mid v \in R \wedge R \in \mathbf{R}_j\}|}{c_j(v)}$ 
   subject to:  $\text{assigned}[v] = \text{false}$ ;
2  $cov_j(u) \leftarrow |\{R \mid u \in R \wedge R \in \mathbf{R}_j\}|$ ; //find best cost-sensitive seed for ad  $j$  as well as its coverage.

```

TI-CARM and TI-CSRSM are very similar, differing only in their greedy seed selection criteria, we only provide the pseudocode of TI-CSRSM (Algorithm 2). Algorithm TI-CSRSM works as follows. For every advertiser  $j$ , we initially set the latent seed set size  $\tilde{s}_j = 1$  (a conservative but safe estimate), create a sample  $\mathbf{R}_j$  of  $\theta_j = L_j(\tilde{s}_j, \varepsilon)$  RR-sets, compute the refined budget  $\tilde{B}_j$  for  $\tilde{s}_j$ , and the safe upper bound  $\bar{s}_j$  (lines 1 – 6). In the main loop, we follow the greedy selection logic of CS-GREEDY. That is, in each round, we first invoke Algorithm 5 to find an unassigned candidate node  $v_j$  that has the largest coverage-to-cost ratio<sup>11</sup> for each advertiser  $j$  whose budget is not yet exhausted. Then, we select, among these (node, advertiser) pairs, the feasible pair  $(v_i, i)$  that has the largest rate of marginal gain in revenue per marginal gain in payment and add it to the solution set, and remove from  $\mathbf{R}_i$  the RR-sets that are covered by node  $v_i$  (lines 10 – 15). While doing so, whenever  $|\tilde{S}_i| = \tilde{s}_i$ , we update the latent seed set size  $\tilde{s}_i$  using Eq. 16, hence  $\tilde{B}_i$ , and sample  $\max\{0, L_i(\tilde{s}_i, \varepsilon) - \theta_i\}$  additional RR-sets into  $\mathbf{R}_i$ . Note that, after adding additional RR-sets, we update the influence spread estimation of current  $\tilde{S}_i$  w.r.t. the updated sample  $\mathbf{R}_i$  by invoking Algorithm 3 to ensure that future marginal gain estimations are accurate (line 22). The main loop executes until the budget of each advertiser is exhausted or no more eligible seed can be found.

For TI-CARM, there are only two differences. First, line 9 of Algorithm 2 is replaced by

$$(v_j, cov_j(v_j)) \leftarrow \text{SelectBestCANode}(\mathbf{R}_j) \quad (\text{Algorithm 4}).$$

Second, line 11 of Algorithm 2 is replaced by

$$i \leftarrow \arg\max_{j=1}^h \pi_j(v_j|\tilde{S}_j) \quad \text{subject to: } \rho_j(\tilde{S}_j \cup \{v_j\}) \leq B_j \wedge \text{assigned}[v_j] = \text{false}.$$

**Deterioration of approximation guarantees.** Since TI-CARM and TI-CSRSM use random RR-sets for the accurate estimation of  $\sigma_i(\cdot), \forall i \in [h]$ , their approximation guarantees slightly deteriorate from the ones of CA-GREEDY and CS-GREEDY (see Theorems 2 and 3). Such deterioration is common to all the state-of-the-art IM algorithms [10, 32–34] that similarly use random RR-sets for influence spread estimation. Our next result provides the deteriorated approximation guarantees for TI-CARM and TI-CSRSM.

<sup>11</sup> Following the definition of  $\tilde{\rho}_j(\cdot)$  as a function of  $\tilde{\pi}_j(\cdot)$ , the node with the largest rate of marginal gain in revenue per marginal gain in payment for a given ad  $j$  corresponds to the node  $u$  with the largest coverage-to-cost ratio for ad  $j$ .

**Theorem 4.** *W.p. at least  $1 - n^{-\ell}$ , TI-CARM (resp. TI-CSRM) returns a solution  $\vec{S} = (\tilde{S}_1, \dots, \tilde{S}_h)$  that satisfies*

$$\pi(\vec{S}) \geq \pi(\vec{S}^*) \cdot \beta - \sum_{i \in [h]} cpe(i) \cdot \varepsilon \cdot OPT_{\tilde{s}_i}.$$

where  $\vec{S}^* = (S_1^*, \dots, S_h^*)$  is the optimal allocation,  $\tilde{s}_i$  is the final latent seed set size estimated for each  $i$  upon termination of TI-CARM (resp. TI-CSRM), and  $\beta$  is the approximation guarantee given in Theorem 2 (resp. Theorem 3).

*Proof.* Let  $\vec{S}^+ = (S_1^+, \dots, S_h^+)$  denote the optimal solution to RM problem on the sample with refined budget constraints, i.e., the feasible allocation that maximizes  $\sum_{i \in [h]} \tilde{\pi}_i(S_i)$  subject to  $\tilde{\rho}_i(S_i) \leq \tilde{B}_i, \forall i \in [h]$ . Since  $\vec{S}$  is the cost-agnostic (resp., cost-sensitive) greedy solution to RM on the sample, we have:

$$\sum_i cpe(i) \cdot n \cdot F_{\mathbf{R}_i}(\tilde{S}_i) \geq \beta \cdot \sum_i cpe(i) \cdot n \cdot F_{\mathbf{R}_i}(S_i^+). \quad (17)$$

Given that  $\vec{S}^+$  is the optimal solution to solving RM on the sample, we also have:

$$\sum_i cpe(i) \cdot n \cdot F_{\mathbf{R}_i}(S_i^+) \geq \sum_i cpe(i) \cdot n \cdot F_{\mathbf{R}_i}(S_i^*). \quad (18)$$

Furthermore, it follows from Lemma 4 that, for any set  $S$  of at most  $\tilde{s}_i$  seeds, we have  $|n \cdot F_{\mathbf{R}_i}(S) - \sigma_i(S)| \geq \frac{\varepsilon}{2} \cdot OPT_{i, \tilde{s}_i}$  w.p. at most  $\frac{n^{-\ell}}{\binom{n}{\tilde{s}_i}}$ . Notice that, we also have  $|S_i^*| \leq \tilde{s}_i$  by definition. Thus, by using Eqs.17 and 18 and a union bound over all  $\binom{n}{\tilde{s}_i}$  estimations, w.p. at least  $1 - n^{-\ell}$  we have:

$$\begin{aligned} & \sum_i cpe(i) \cdot \sigma_i(\tilde{S}_i) \\ & \geq \sum_i cpe(i) \cdot \left( n \cdot F_{\mathbf{R}_i}(\tilde{S}_i) - \frac{\varepsilon}{2} \cdot OPT_{i, \tilde{s}_i} \right) \\ & = \sum_i cpe(i) \cdot n \cdot F_{\mathbf{R}_i}(\tilde{S}_i) - \sum_i cpe(i) \cdot \frac{\varepsilon}{2} \cdot OPT_{i, \tilde{s}_i} \\ & \geq \beta \cdot \sum_i cpe(i) \cdot n \cdot F_{\mathbf{R}_i}(S_i^+) - \sum_i cpe(i) \cdot \frac{\varepsilon}{2} \cdot OPT_{i, \tilde{s}_i} \\ & \geq \beta \cdot \sum_i cpe(i) \cdot n \cdot F_{\mathbf{R}_i}(S_i^*) - \sum_i cpe(i) \cdot \frac{\varepsilon}{2} \cdot OPT_{i, \tilde{s}_i} \\ & \geq \beta \cdot \sum_i cpe(i) \cdot \left( \sigma_i(S_i^*) - \frac{\varepsilon}{2} \cdot OPT_{i, \tilde{s}_i} \right) \\ & \quad - \sum_i cpe(i) \cdot \frac{\varepsilon}{2} \cdot OPT_{i, \tilde{s}_i} \\ & \geq \beta \cdot \pi(\vec{S}^*) - \sum_{i \in [h]} cpe(i) \cdot \varepsilon \cdot OPT_{i, \tilde{s}_i}, \end{aligned}$$

where the last inequality follows upon noting that  $\beta < 1$ .  $\square$

As a corollary to Theorem 4, Lemma 3 and Lemma 4, the following result is immediate.

**Theorem 5.** *W.p. at least  $1 - n^{-\ell}$ , TI-CARM (resp. TI-CSRM) returns an approximate greedy solution  $\vec{S} = (\tilde{S}_1, \dots, \tilde{S}_h)$  that*

**Table 1: Statistics of network datasets.**

	FLIXSTER	EPINIONS	DBLP	LIVEJOURNAL
#nodes	30K	76K	317K	4.8M
#edges	425K	509K	1.05M	69M
type	directed	directed	undirected	directed

**Table 2: Advertiser budgets and cost-per-engagement values.**

Dataset	Budgets			CPEs		
	mean	max	min	mean	max	min
FLIXSTER	10.1K	20K	6K	1.5	2	1
EPINIONS	8.5K	12K	6K	1.5	2	1

is budget feasible, i.e.,  $\rho_i(\tilde{S}_i) \leq B_i$ , for all  $i$ , and achieves an approximation that satisfies

$$\pi(\vec{S}) \geq \pi(\vec{S}^*) \cdot \beta - \sum_{i \in [h]} cpe(i) \cdot \varepsilon \cdot OPT_{\tilde{s}_i}.$$

$\beta$  is the approximation guarantee given in Theorem 2 (resp. Theorem 3).

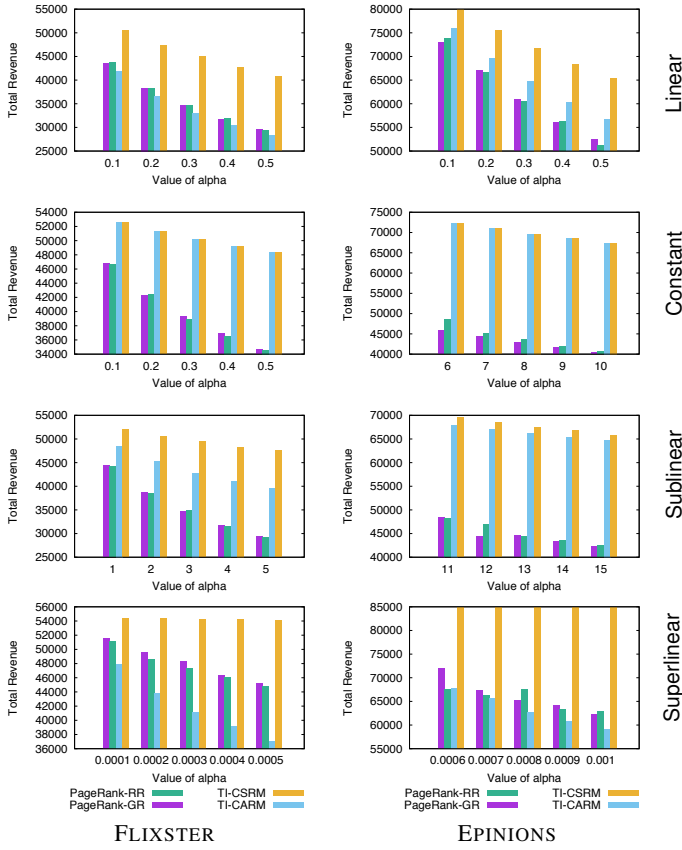
## 5. EXPERIMENTS

We conducted extensive experiments to evaluate (i) the quality of our proposed algorithms, measured by the revenue achieved vis à vis the incentives paid to seed users, and (ii) the efficiency and scalability of the algorithms w.r.t. advertiser budgets, which indirectly control the number of seeds required, and w.r.t. the number of advertisers, which effectively controls the size of the graph. All experiments were run on a 64-bit OpenSUSE Linux server with Intel Xeon 2.90GHz CPU and 264GB memory. As a preview, our largest configuration is LIVEJOURNAL with 20 ads, which effectively yields a graph with  $69M \times 20 \approx 1.4B$  edges; this is comparable with [34], whose largest dataset has 1.5B edges.

**Data.** Our experiments were conducted on four real-world social networks, whose basic statistics are summarized in Table 1. We used FLIXSTER and EPINIONS for quality experiments and DBLP and LIVEJOURNAL for scalability experiments. FLIXSTER is from a social movie-rating website (<http://www.flixster.com/>), which contains movie ratings by users along with timestamps. We use the topic-aware influence probabilities and the item-specific topic distributions provided by Barbieri et al. [8], who learned the probabilities using MLE for the TIC model, with  $L = 10$  latent topics. We set the default number of advertisers  $h = 10$  and used five of the learned topic distributions from the provided FLIXSTER dataset, in such a way that every two ads are in pure competition, i.e., have the same topic distribution, with probability 0.91 in one randomly selected latent topic, and 0.01 in all others. This way, among  $h = 10$  ads, every two ads are in pure competition with each other while having a completely different topic distribution than the rest, representing a diverse marketplace of ads. EPINIONS is a who-trusts-whom network taken from a consumer review website (<http://www.epinions.com/>). Likewise, we set  $h = 10$  and use the Weighted-Cascade model [24], where  $p_{u,v}^i = 1/|N^{in}(v)|$  for all ads  $i$ . Notice that this corresponds to  $L = 1$  topic for EPINIONS dataset, hence, all the ads are in pure competition.

For scalability experiments, we used two large networks<sup>12</sup> DBLP and LIVEJOURNAL. DBLP is a co-authorship graph (undirected) where nodes represent authors and there is an edge between two nodes if they have co-authored a paper indexed by DBLP. We direct all edges in both directions. LIVEJOURNAL is an online blogging site where users can declare which other users are their friends. In all datasets, advertiser budgets and CPEs were chosen

<sup>12</sup> Available at <http://snap.stanford.edu/>.



**Figure 2: Total revenue as a function of  $\alpha$ , on FLIXSTER (left) and EPINIONS (right), for linear, constant, sublinear, and superlinear incentive models.**

in such a way that the total number of seeds required for all ads to meet their budgets is less than  $n$ . This ensures that no ad is assigned an empty seed set. For lack of space, instead of enumerating all CPEs and budgets, we give a statistical summary in Table 2. The same information for DBLP and LIVEJOURNAL is provided later.

**Seed incentive models.** In order to understand how the algorithms perform w.r.t. different seed user incentive assignments, we used four different methods that directly control the range between the minimum and maximum singleton payments:

- **Linear incentives:** proportional to the ad-specific singleton influence spread of the nodes, i.e.,  $c_i(u) = \alpha \cdot \sigma_i(\{u\})$ ,  $\forall u \in V, i \in [h]$ ,
- **Constant incentives:** the average of the ad-specific total linear seed user incentives, i.e.,  $c_i(u) = \alpha \cdot \frac{\sum_{v \in V} \sigma_i(\{v\})}{n}$ ,  $\forall u \in V, i \in [h]$ ,
- **Sublinear incentives:** obtained by taking the logarithm of the ad-specific singleton influence spread of the nodes, i.e.,  $c_i(u) = \alpha \cdot \log(\sigma_i(\{u\}))$ ,  $\forall u \in V, i \in [h]$ ,
- **Superlinear incentives:** obtained by using the squared ad-specific singleton influence spread of the nodes, i.e.,  $c_i(u) = \alpha \cdot (\sigma_i(\{u\}))^2$ ,  $\forall u \in V, i \in [h]$ ,

where  $\alpha > 0$  denotes a fixed amount in dollar cents set by the host, which controls how expensive the seed user incentives are.

On FLIXSTER and EPINIONS we used Monte Carlo simulations (5K runs<sup>13</sup>) to compute  $\sigma_i(\{u\})$ . On DBLP and LIVEJOURNAL,

we use the out-degree of the nodes as a proxy to  $\sigma_i(\{u\})$  due to the prohibitive computational cost of Monte Carlo simulations.

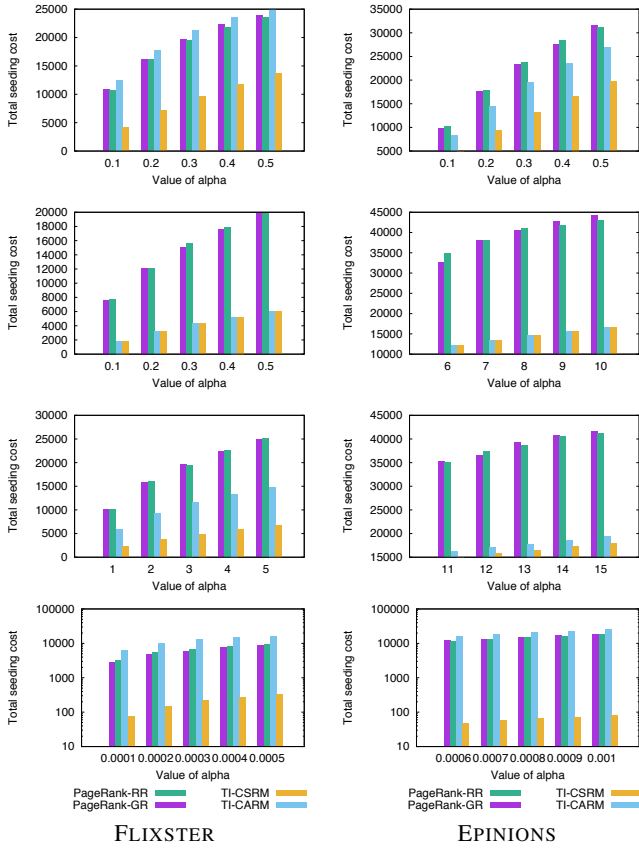
**Algorithms.** We compared four algorithms in total. Wherever applicable, we set the parameter  $\varepsilon$  to be 0.1 for quality experiments on FLIXSTER and EPINIONS, and 0.3 for scalability experiments on DBLP and LIVEJOURNAL, following the settings used in [34].

- **TI-CSRM** (Algorithm 2) that uses Algorithm 5 to find the best (cost-sensitive) candidate node for each advertiser (line 9), and selects among those the (node, advertiser) pair that provides the maximum rate of marginal gain in revenue per marginal gain in advertiser’s payment (line 11).
- **TI-CARM:** Cost-agnostic version of Algorithm 2 that uses Algorithm 4 to find the best (cost-agnostic) candidate node for each advertiser (replacing line 9), and selects among those the (node, advertiser) pair with the maximum increase in the revenue of the host (replacing line 11).
- **PageRank-GR:** A baseline that selects a candidate node for each advertiser based on the ad-specific PageRank ordering of the nodes (replacing line 9), and selects among those the (node, advertiser) pair that provides the maximum increase in the revenue of the host (replacing line 11). Since the selection is made greedily, we refer to this algorithm as PageRank-GR.
- **PageRank-RR:** Another PageRank-based baseline that selects a candidate node for each advertiser based on the ad-specific PageRank ordering of the nodes (replacing line 9), and uses a Round-Robin (RR in short) ordering of the advertisers for the assignment of their candidates into their seed sets.

**Revenue vs.  $\alpha$ .** We first compare the total revenue achieved by the four algorithms for four different seed incentive models and with varying levels of  $\alpha$  (Figure 2). Recall that by definition, a smaller  $\alpha$  value indicates lower seed costs for all users. Across all different values of  $\alpha$  and all seed incentive models, it can be seen that TI-CSRM consistently achieves the highest revenue, often by a large margin, which increases as  $\alpha$  grows. For instance, on EPINIONS, when  $\alpha = 0.5$ , TI-CSRM achieved 15.3%, 24.3%, 27.6% more revenue than TI-CARM, PageRank-RR, and PageRank-GR respectively on the linear incentive model, while these values for superlinear incentive model respectively are 25.2%, 25.8%, 18.1%. Notice that for the constant incentive model, the advantage of being cost-sensitive is nullified, hence TI-CARM and TI-CSRM end up performing identically as expected. Figure 3 reports the cost-effectiveness of the algorithms. Across all different values of  $\alpha$  and all incentive models, it can be seen that TI-CSRM consistently achieves the lowest total seed costs. This is as expected, since its seed allocation strategy takes into account revenue obtained per seed user cost.

Notice that in three of the test cases, i.e., linear seed incentives on FLIXSTER and superlinear seed incentives on both datasets, TI-CARM has slightly worse performance than the two PageRank-based heuristics (e.g., about 4–7% drop in revenue). This can be explained by the fact that, while TI-CARM picks seeds of high spreading potential (i.e., highest marginal revenue) without considering costs, the two PageRank-based heuristics may instead select seeds of low quality (i.e., low marginal revenue), but also of very low cost. This might create a situation in which the PageRank-based heuristics may select many more seeds, but with a smaller total seed cost than TI-CARM, hence, allowing the budget to be spent more on engagements that translate to higher revenue, mimicking the cost-sensitive behavior. On the other hand TI-CSRM always spends the given budget judiciously by selecting seeds with the best rate of marginal revenue per cost. Thus, it is able to use the budget more intelligently, which explains its superiority in all test cases. This hypothesis is confirmed by our experiments. E.g.,

<sup>13</sup>We didn’t observe any significant change in the influence spread estimation beyond 5K runs for both datasets.



**Figure 3: Total seeding cost as a function of  $\alpha$ , on FLIXSTER (left) and EPINIONS (right), for linear, constant, sublinear, and superlinear cost models.**

on FLIXSTER with linear seed incentives, we observed that the average values of marginal gain in revenue, seed user cost, and rate of marginal gain per cost obtained by PageRank-GR were respectively 2.67, 0.44, and 7.48, while the corresponding numbers for TI-CARM were 13.47, 2.7, and 4.89, and those for TI-CSR were 1.28, 0.12, and 9.95 respectively. While the two PageRank-based heuristics could obtain higher revenue than TI-CARM on FLIXSTER with linear and superlinear incentives, and on EPINIONS with superlinear incentives, they were greatly outperformed by TI-CARM, hence TI-CSR, in the other incentive models, showing that such heuristics are not robust to different seed incentive models, and can only get “lucky” to the extent they can mimic the cost-sensitive behavior.

Finally, as shown in Figure 2, the extent to which TI-CSR outperforms TI-CARM on both datasets is higher with linear incentives than with sublinear incentives. For instance, on FLIXSTER, TI-CSR achieved 45% more revenue than TI-CARM in the linear model, while this improvement drops to 20% in the sublinear model. To understand how the seeds’ expensiveness levels affect this improvement, we checked the values of singleton payments and found that the maximum singleton payment ( $\rho_{max}$ ) is 1347 times more expensive than the minimum singleton payment ( $\rho_{min}$ ) in the linear model, while it is 725 times more expensive in the sublinear model that has lower improvement rate. This relation is expected as higher variety in the expensiveness levels of the seeds require to use the budget more cleverly, hence, with more cost-effective strategies. Notice that this finding is also in line with our discussion following the proof of Theorem 3.

It is also worth noting that, from Figure 3, TI-CSR is two to

three orders of magnitude more cost-efficient than the rest in the superlinear model, and this gap is larger than that attained in linear, constant, and sublinear scenarios.

**Revenue & running time vs. window size.** Hereafter all presented results will be w.r.t. linear seed incentives, unless otherwise noted. As stated before in Section 4, TI-CSR needs to compute  $\sigma_i(v|S_i^{t-1})$ ,  $\forall v : (v, i) \in \mathcal{E}^{t-1}$  while  $u_i^t$  might even correspond to the node that has the *minimum* marginal gain in influence spread for iteration  $t$ . To have a closer look at how the revenue evolves when the seed selection criterion changes from cost-agnostic to cost-sensitive, we restrict TI-CSR to find the best cost-sensitive candidate nodes for each advertiser (line 9) among only the  $w$  nodes that have the highest marginal gain in revenue at each iteration. We refer to  $w$  as the “window size”. Notice that TI-CARM corresponds to the case when  $w = 1$ , i.e., in this case, TI-CSR inspects only the node with the maximum marginal gain in revenue.

We report the results of TI-CSR with various window sizes in Fig. 4, which depicts the revenue vs. running time tradeoff. Each figure corresponds to one dataset and one particular  $\alpha$  value. The X-axis is in log-scale. As expected, the maximum revenue is achieved when TI-CSR implements the full window  $w = n$ , i.e., when all the (feasible) nodes are inspected at each iteration for each advertiser. The running time can go up quickly as the window size increases to  $n$ . This is expected as the seed nodes selected do not necessarily provide high marginal gain in revenue, thus, TI-CSR needs to use higher number of seed nodes, hence, much more RR-sets to achieve accuracy, compared to TI-CARM.

**Scalability.** We tested the scalability of TI-CARM and TI-CSR on two larger graphs, DBLP and LIVEJOURNAL. In all scalability experiments, we use a window size of  $w = 5000$  nodes for TI-CSR due to its good revenue vs running time trade-off. For simplicity, all CPEs were set to 1. The influence probability on each edge  $(u, v) \in E$  was computed using the Weighted-Cascade model [24], where  $p_{u,v}^i = 1/|N^{in}(v)|$  for all ads  $i$ . We set  $\alpha = 0.2$  and  $\varepsilon = 0.3$ . This setting is well-suited for testing scalability as it simulates a fully competitive case: all advertisers compete for the same set of influential users (due to all ads having the same distribution over the topics), and hence it will “stress-test” the algorithms by prolonging the seed selection process.

Figure 5(a) and 5(b) depict the running time of TI-CARM and TI-CSR as the number of advertisers goes up from 1 to 20, while the budget is fixed (10K for DBLP and 100K for LIVEJOURNAL). As can be seen, the running time increases mostly in a linear manner, and TI-CSR is only slightly slower than TI-CARM. Figure 5(c) and 5(d) depict the running time of TI-CARM and TI-CSR as the budget increases, while the number of advertisers is fixed at  $h = 5$ . We can also see that the increasing trend is mostly linear for TI-CSR, while TI-CARM’s time goes in a flatter fashion. All in all, both algorithms exhibit decent scalability.

Table 3 shows the memory usage of TI-CARM and TI-CSR when  $h$  increases. TI-CSR in general needs to use higher memory than TI-CARM due to its requirement to generate more RR sets that ensures accuracy for using higher seed set size than TI-CARM. On DBLP, TI-CARM and TI-CSR respectively uses a total of 4676 and 7276 seed nodes for  $h = 20$ . On LIVEJOURNAL TI-CSR used typically between 20% to 40% more memory than TI-CARM: TI-CARM and TI-CSR respectively uses a total of 4327 and 6123 seed nodes for  $h = 20$ .

## 6. RELATED WORK

**Computational advertising.** Considerable work has been done in sponsored search and display ads [18–21, 28, 30]. In sponsored

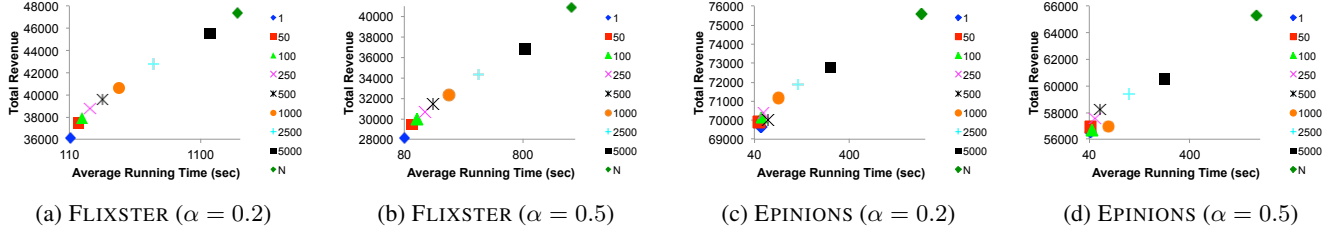


Figure 4: Revenue vs running time tradeoff on FLIXSTER and EPINIONS for two different value of  $\alpha$ .

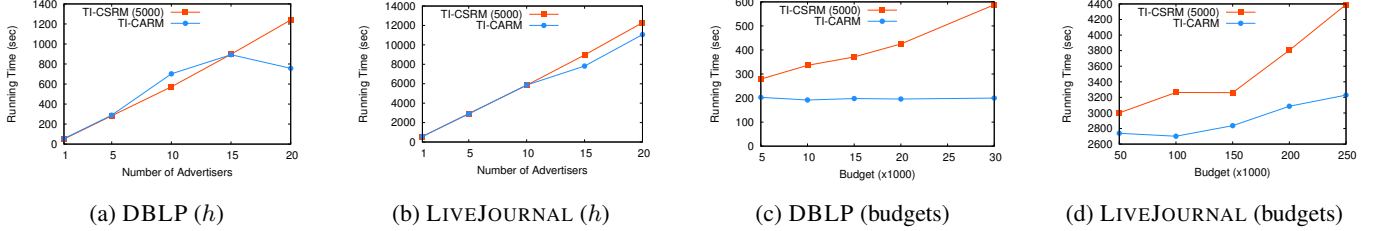


Figure 5: Running time of TI-CARM and TI-CSRM on DBLP and LIVEJOURNAL

Table 3: Memory usage (GB).

DBLP	$h = 1$	5	10	15	20
TI-CARM	1.6	7.5	14.9	22.4	29.8
TI-CSRM (5000)	1.6	7.6	15.1	22.7	30.2
LIVEJOURNAL	$h = 1$	5	10	15	20
TI-CARM	2.5	12.1	25.3	39.4	54.4
TI-CSRM (5000)	3.4	15.9	31.2	49.1	67.5

search, revenue maximization is formalized as the well-known *Ad-words* problem [29]. Given a set of keywords and bidders with their daily budgets and bids for each keyword, words need to be assigned to bidders upon arrival, to maximize the revenue for the day, while respecting bidder budgets. This can be solved with a competitive ratio of  $(1 - 1/e)$  [29].

**Social advertising.** In comparison with computational advertising, social advertising is in its infancy. Recent efforts, including Tucker [35] and Bakshy et al. [5], have shown, by means of field studies on sponsored posts in Facebook’s News Feed, the importance of taking social influence into account when developing social advertising strategies. However, literature on exploiting social influence for social advertising is rather limited. Bao and Chang have proposed *AdHeat* [7], a social ad model considering social influence in addition to relevance for matching ads to users. Their experiments show that AdHeat significantly outperforms the relevance model on click-through-rate (CTR). Wang et al. [36] propose a new model for learning relevance and apply it for selecting relevant ads for Facebook users. Neither of these works studies viral ad propagation or revenue maximization.

Chalermsook et al. [12] study revenue maximization for the host, when dealing with multiple advertisers. In their setting, each advertiser pays the host an amount for each product adoption, up to a budget. In addition, each advertiser also specifies the maximum size of its seed set. This additional constraint considerably simplifies the problem compared to our setting, where the absence of a prespecified seed set size is a *significant challenge*.

Aslay et al. [4] study regret minimization for a host supporting campaigns from multiple advertisers. Here, regret is the difference between the monetary budget of an advertiser and the value of expected number of engagements achieved by the campaign, based on the CPE pricing model. They share with us the pricing model and advertiser budget. However, they do not consider seed user costs.

Besides they attack a very different optimization problem and their algorithms and results do not carry over to our setting.

Abbassi et al. [2] study a cost-per-mille (CPM) model in display advertising. The host enters into a contract with each advertiser to show their ad to a fixed number of users, for an agreed upon CPM amount per thousand impressions. The problem is that of selecting the sequence of users to show the ads to, in order to maximize the expected number of clicks. This is a substantially different problem which they show is APX-hard and propose heuristic solutions.

Alon et al. [3] study budget allocation among channels and influential customers, with the intuition that a channel assigned a higher budget will make more attempts at influencing customers. They do not take into account viral propagation. Their main result is that for some influence models the budget allocation problem can be approximated, while for others it is inapproximable. Notably, none of these previous works studies *incentivized social advertising* where the seed users are paid monetary incentives.

**Viral marketing.** Kempe et al. [24] formalize the influence maximization problem which requires to select  $k$  seed nodes, where  $k$  is a cardinality budget, such that the expected spread of influence from the selected seeds is maximized. Of particular note are the recent advances (already reviewed in Section 4) that have been made in designing scalable approximation algorithms [10, 14, 32–34] for this hard problem. Numerous variants of the influence maximization problem have been studied over the years, including competition [9, 11], host perspective [4, 27], non-uniform cost model for seed users [26, 31], and fractional seed selection [17]. However, to our knowledge, there has been no previous work that addresses incentivized social advertising, while leveraging viral propagation of social ads and handling advertiser budgets.

## 7. CONCLUSIONS

In this paper, we initiate the investigation of incentivized social advertising, by formalizing the fundamental problem of revenue maximization from the host perspective. In our formulation, incentives paid to the seed users are determined by their demonstrated past influence in the topic of the specific ad. We show that, keeping all important factors – topical relevance of ads, their propensity for social propagation, the topical influence of users, seed users’ incentives, and advertiser budgets – in consideration, the problem



of revenue maximization in incentivized social advertising is NP-hard and it corresponds to the problem of monotone submodular function maximization subject to a partition matroid constraint on the ads-to-seeds allocation and multiple submodular knapsack constraints on the advertiser budgets. For this problem, we devise two natural greedy algorithms that differ in their sensitivity to seed user incentive costs, provide formal approximation guarantees, and achieve scalability by adapting to our context recent advances made in scalable estimation of expected influence spread.

Our work takes an important first step toward enriching the framework of incentivized social advertising with powerful ideas from viral marketing, while making the latter more applicable to real-world online marketing. It opens up several interesting avenues for further research: (i) it remains open whether our winning algorithm TI-CSR can be made more memory efficient hence more scalable; (ii) it remains open whether the approximation bound for CS-GREEDY provided in Theorem 3 is tight; (iii) it is interesting to integrate hard competition constraints into the influence propagation process; (iv) it is worth studying our problem in an online adaptive setting where the partial results of the campaign can be taken into account while deciding the next moves. All these directions offer a wealth of possibilities for future work.

## 8. REFERENCES

- [1] <https://arxiv.org/abs/1612.00531>.
- [2] Z. Abbassi, A. Bhaskara, and V. Misra. Optimizing display advertising in online social networks. In *WWW 2015*.
- [3] N. Alon, I. Gamzu, and M. Tennenholtz. Optimizing budget allocation among channels and influencers. In *WWW 2012*.
- [4] Ç. Aslay, W. Lu, F. Bonchi, A. Goyal, and L. V. S. Lakshmanan. Viral marketing meets social advertising: Ad allocation with minimum regret. *PVLDB*, 8(7):822–833, 2015.
- [5] E. Bakshy, D. Eckles, R. Yan, and I. Rosenn. Social influence in social advertising: evidence from field experiments. In *EC 2012*.
- [6] E. Bakshy, J. M. Hofman, W. A. Mason, and D. J. Watts. Everyone’s an influencer: quantifying influence on twitter. In *WSDM 2011*.
- [7] H. Bao and E. Y. Chang. Adheat: An influence-based diffusion model for propagating hints to match ads. In *WWW 2010*.
- [8] N. Barbieri, F. Bonchi, and G. Manco. Topic-aware social influence propagation models. In *ICDM 2012*.
- [9] S. Bharathi, D. Kempe, and M. Salek. Competitive influence maximization in social networks. In *WINE 2007*.
- [10] C. Borgs, M. Brautbar, J. T. Chayes, and B. Lucier. Maximizing social influence in nearly optimal time. In *SODA 2014*.
- [11] T. Carnes, C. Nagarajan, S. M. Wild, and A. van Zuylen. Maximizing influence in a competitive social network: a follower’s perspective. In *ICEC 2007*.
- [12] P. Chalermsook, A. D. Sarma, A. Lall, and D. Nanongkai. Social network monetization via sponsored viral marketing. In *SIGMETRICS 2015*.
- [13] W. Chen, C. Wang, and Y. Wang. Scalable influence maximization for prevalent viral marketing in large-scale social networks. In *KDD 2010*.
- [14] E. Cohen, D. Delling, T. Pajor, , and R. F. Werneck. Sketch-based influence maximization and computation: Scaling up with guarantees. In *CIKM 2014*.
- [15] M. Conforti and G. Cornuéjols. Submodular set functions, matroids and the greedy algorithm: tight worst-case bounds and some generalizations of the rado-edmonds theorem. *Discrete applied mathematics*, 7(3):251–274, 1984.
- [16] P. Dagum, R. Karp, M. Luby, and S. Ross. An optimal algorithm for monte carlo estimation. *SIAM Journal on computing*, 29(5):1484–1496, 2000.
- [17] E. D. Demaine, M. Hajiaghayi, H. Mahini, D. L. Malec, S. Raghavan, A. Sawant, and M. Zadimoghaddam. How to influence people with partial incentives. In *WWW 2014*.
- [18] N. R. Devanur, B. Sivan, and Y. Azar. Asymptotically optimal algorithm for stochastic adwords. In *EC 2012*.
- [19] J. Feldman, M. Henzinger, N. Korula, V. S. Mirrokni, and C. Stein. Online stochastic packing applied to display ad allocation. In *ESA 2010*.
- [20] J. Feldman, N. Korula, V. S. Mirrokni, S. Muthukrishnan, and M. Pál. Online ad assignment with free disposal. In *WINE 2009*.
- [21] G. Goel and A. Mehta. Online budgeted matching in random input models with applications to adwords. In *SODA 2008*.
- [22] R. Iyer. *Submodular Optimization and Machine Learning: Theoretical Results, Unifying and Scalable Algorithms, and Applications*. PhD thesis, Univ. of Washington, 2015.
- [23] R. K. Iyer and J. A. Bilmes. Submodular optimization with submodular cover and submodular knapsack constraints. In *NIPS 2013*.
- [24] D. Kempe, J. M. Kleinberg, and É. Tardos. Maximizing the spread of influence through a social network. In *KDD 2003*.
- [25] B. Korte and D. Hausmann. An analysis of the greedy heuristic for independence systems. In *Annals of Discrete Mathematics 1978*.
- [26] J. Leskovec, A. Krause, C. Guestrin, C. Faloutsos, J. M. VanBriesen, and N. S. Glance. Cost-effective outbreak detection in networks. In *KDD 2007*.
- [27] W. Lu, F. Bonchi, A. Goyal, and L. V. Lakshmanan. The bang for the buck: fair competitive viral marketing from the host perspective. In *KDD 2013*.
- [28] A. Mehta. Online matching and ad allocation. *Foundations and Trends in Theoretical Computer Science*, 8(4):265–368, 2013.
- [29] A. Mehta, A. Saberi, U. V. Vazirani, and V. V. Vazirani. Adwords and generalized online matching. *J. ACM*, 54(5), 2007.
- [30] V. S. Mirrokni, S. O. Gharan, and M. Zadimoghaddam. Simultaneous approximations for adversarial and stochastic online budgeted allocation. In *SODA 2012*.
- [31] H. Nguyen and R. Zheng. On budgeted influence maximization in social networks. *IEEE Journal on Selected Areas in Communications*, 31(6):1084–1094, 2013.
- [32] H. T. Nguyen, M. T. Thai, and T. N. Dinh. Stop-and-stare: Optimal sampling algorithms for viral marketing in billion-scale networks. In *SIGMOD 2016*.
- [33] Y. Tang, Y. Shi, and X. Xiao. Influence maximization in near-linear time: A martingale approach. In *SIGMOD 2015*.
- [34] Y. Tang, X. Xiao, and Y. Shi. Influence maximization: Near-optimal time complexity meets practical efficiency. *SIGMOD 2014*.
- [35] C. Tucker. Social advertising. Available at SSRN 1975897, 2012.
- [36] C. Wang, R. Raina, D. Fong, D. Zhou, J. Han, and G. Badros. Learning relevance from heterogeneous social network and its application in online targeting. In *SIGIR 2011*.
- [37] Y. Tang, X. Tang, X. Xiao, and Y. Junsong. Online processing algorithms for influence maximization. *SIGMOD 2018*.