

Optimal Price Profile for Influential Nodes in Online Social Networks

Yuqing Zhu¹ · Jing Tang² · Xueyan Tang¹

Received: date / Accepted: date

Abstract Influential nodes with rich connections in online social networks (OSNs) are of great values to initiate marketing campaigns. However, the potential influence spread that can be generated by these influential nodes is hidden behind the structures of OSNs, which are often held by OSN providers and unavailable to advertisers for privacy concerns. A social advertising model known as influencer marketing is to have OSN providers offer and price candidate nodes for advertisers to purchase for seeding marketing campaigns. In this setting, a reasonable price profile for the candidate nodes should effectively reflect the expected influence gain they can bring in a marketing campaign.

In this paper, we study the problem of pricing the influential nodes based on their expected influence spread to help advertisers select the initiators of marketing campaigns without the knowledge of OSN structures. We design a function characterizing the divergence between the price and the expected influence of the initiator sets. We formulate the problem to minimize the divergence and derive an optimal price profile. An advanced algorithm is developed to estimate the price profile with accuracy guarantees. Experiments with real

OSN datasets show that our pricing algorithm can significantly outperform other baselines.

Keywords Pricing · Optimization · Influence Spread · Online Social Network

1 Introduction

Online Social Networks (OSNs) attract billions of users to share information and bring new approaches to promote product sales or activity engagement. Real-world examples of web-based social networks include Facebook, Twitter, Orkut, etc. According to Facebook's official statistics¹, it has 2.13 billion monthly active users as of December 31, 2017. Given the tremendous number of active users, information can be propagated widely and rapidly through OSNs with the *word of mouth* effects. The interpersonal connections between individuals can strongly impact their decisions and behaviors. Applications like social advertising naturally emerge to make use of OSNs for information diffusion [12]. Nowadays, the advertisement market in OSNs is growing at an amazing speed. For example, eMarketer [13] estimates that advertisers are expected to spend \$35.98 billion on social media to promote their products. Fortune [28] claims that the expenditure of advertisement on social media will exceed traditional newspapers by 2020, which will be over \$50 billion.

In online social advertising, some influencers accept free products as rewards for running a marketing campaign. The well-known influence maximization problem emerges from giving a limited number of free samples to a subset of individuals to trigger a cascade of influence [21]. Meanwhile, some influencers charge a certain amount of money from advertisers. According to

Yuqing Zhu
yuqing002@e.ntu.edu.sg

✉ Jing Tang
jingtang@ust.hk

Xueyan Tang
asxytang@ntu.edu.sg

¹ School of Computer Science and Engineering, Nanyang Technological University, Singapore

² Data Science and Analytics Thrust, The Hong Kong University of Science and Technology, China

¹ <https://newsroom.fb.com/company-info/>

a survey conducted by an influencer platform named Klear, brands make an average payment of \$114 per video post on Instagram to nano-influencers who have between 500 and 5,000 followers, and \$775 to power users with followings between 30,000 and 50,000 for an Instagram video [15]. The rising cost for doing online advertising attracts investment in influencers (called influencer marketing), which often involves buying a list of influencer contracts and paying them to promote a product [14]. In 2016, an online celebrity named Papi Jiang with 10 million fans on Weibo, a Twitter-like micro blogging site, was valued at around 42 million dollars (300 million RMB) and received 2 million dollars investment for her potential market value without selling anything yet. According to a survey conducted by Influencer Marketing Hub, the majority of advertisers see influencer marketing as a direction and plan to increase their influencer marketing budget [20].

In practice, the social graph is normally possessed by OSN providers and kept secret for privacy reasons. Hence, it is difficult for advertisers to infer the values of influencers—only 39% of US marketers feel confident in identifying right influencers according to a Cision and PRWeek survey [16]. It remains an open question how to set reasonable prices on the influencers for their market values. Different from the advertisers, the OSN providers hold the structure of the networks and can identify reliable influencers and leverage data to set the marketing price properly. In other words, the OSN providers can offer the prices of the influencers to the advertisers. In fact, the OSN providers have started setting up platforms to facilitate the influencer search and selection process, as well as making the system more transparent and easier for both advertisers and influencers [20]. Recently, YouTube offers access to the set of influencers on the platform FameBit². On average, hiring an influencer costs \$20 a video per 1,000 subscribers.

Intuitively, the price of seeding any set of users should effectively reflect the expected influence spread that these users can generate in the campaigns. In this way, the advertisers can hold a clear view over the influence potential of the seeds selected and make more sensible business decisions. There are several intuitive ways to price users or nodes in OSNs. A simple strategy is to set the price of each node based on its degree in the OSN. This strategy is rather primitive since the degree of a node is not necessarily proportional to its actual influence spread. Another intuitive strategy is to set the price of each node according to its expected influence spread when selected as the only seed. However, when multiple nodes are selected to seed a campaign,

the influence spreads generated by different nodes may overlap substantially. Thus, the influence spreads of singleton seeds may not effectively reflect their influence contributions when they jointly initiate a campaign. A straightforward solution to precisely describe the influence contributions of the nodes in various seed sets is to derive a separate price of each node for including in each possible seed set. This method is unfortunately computationally expensive to implement due to the huge numbers of nodes and possible seed sets in real social networks.

In this paper, we propose a new pricing strategy that can effectively reflect the value of the nodes in any seed set. We define a function to measure the difference of the price of a randomly chosen seed set from its expected marketing value and formulate an optimization problem of pricing the nodes to minimize the difference. The optimization problem is challenging to solve in several aspects. First, in order to narrow down the divergence between the price and the expected influence, we need to calculate the expected influence spread of a seed set, which is #P-hard even for simple diffusion models [6, 7]. Second, as the number of possible seed sets grows exponentially with the number of candidate nodes offered by the OSN provider, it is computationally intractable to compute the expected influence spread for all possible seed sets. Furthermore, the nodes may have different contributions to the influence spreads of different seed sets, which makes it even more difficult to set a reasonable price for each node.

To tackle the pricing problem, we make the following major contributions in this paper:

- We propose a novel problem domain of pricing the nodes based on their expected influence spread to help advertisers select the initiators of marketing campaigns.
- We design a function to characterize the divergence between the price and the expected influence of seed sets and formulate and solve an optimization problem to minimize the divergence.
- We devise an efficient algorithm based on random reverse reachable sets [3] to compute the prices for the nodes. An advanced estimation algorithm is also developed to ensure that the estimated prices have accuracy guarantees.
- Extensive experiments based on real OSN datasets confirm that our pricing algorithm can yield high quality solutions and significantly outperform other baselines.

In our preliminary work [39], we studied the *budgeted* pricing problem where the total price of all can-

² <https://famebit.com/>

didate nodes equals a given budget, whereas the current paper focuses on the pricing problem without involving any budget constraint. Our techniques and analysis are tailored to the pricing problem. For completeness, the results for the budgeted pricing problem are also included in this paper, e.g., Sections 4.3 and 5.4.

The rest of this paper is organized as follows. Section 2 reviews the related work. Section 3 introduces the preliminaries. Section 4 presents the pricing problem and the solution. Section 5 elaborates the algorithm to estimate the node prices. Section 6 describes the experimental evaluation. Finally, Section 7 concludes the paper.

2 Related Work

Domingos and Richardson [29] were the first to study viral marketing as an algorithmic problem. They proposed approximation algorithms to determine the influential users and demonstrated that different sets of seed users in a marketing campaign can produce substantially different influence spreads. Kempe et al. [21] showed that the optimization problem of selecting the most influential seed set of a given size is NP-hard. They showed that the influence function is a submodular function under the Independent Cascade and Linear Threshold diffusion models [21]. They proposed a $(1 - 1/e)$ -approximation greedy algorithm utilizing Monte-Carlo simulations. Follow-up work has mostly focused on improving the efficiency of the algorithm implementation for large-scale OSNs based on the submodularity property or heuristics [3, 5–7, 10, 24, 26, 31–33, 37, 38]. Furthermore, some recent work utilized adaptive algorithms to improve the performance for various influence based optimization problems, including adaptive influence maximization [17, 18] and adaptive seed minimization [36], etc. There was also work studying profit maximization in OSNs to optimize the profit return of viral marketing. Lu et al. [25] extended the classical Linear Threshold model to incorporate product prices and user valuations, and factor them into the user’s decision process of adopting a product. They greedily chose the seeds with the greatest profit potential. Tang et al. [30, 34] defined the profit metric as the benefit of influence spread less the cost for seed selection and proposed a general problem of profit maximization for viral marketing. Huang et al. [19] studied adaptive profit maximization. All the above work focused on designing the seed selection algorithms for advertisers based on the premise that complete social network structures are available to advertisers. However, such information is normally kept secret by OSN providers for business and privacy reasons [4, 22]. Our

work in this paper aims to offer a seed pricing solution for seed selection without releasing the social network structures to the advertisers. Recently, Tang et al. [35] studied the profit maximization problem from the OSN provider’s perspective by taking the cost of information diffusion over the social network into account. In addition, Aslay et al. proposed and tackled other practical problems of regret minimization [1] and revenue maximization [2] in online social advertising. Nevertheless, they did not provide any value-based seed pricing solution for advertisers. Different from all the above studies, in this paper, we aim to tighten the relationship between the price setting and the seed’s influence spread.

3 Preliminaries

3.1 Influence Spread

An OSN can be modeled as a directed graph $G = (V, E)$ with a set V of nodes and a set E of edges. Users are represented by the nodes and connections between users are represented by edges. For each edge $(u, v) \in E$, we say that v is an *out-neighbor* of u and u is an *in-neighbor* of v .

Many models have been proposed to capture the diffusion process in the OSN. Our problem definition and solution are general and can be used for various diffusion models. For the purpose of illustrating basic concepts, we briefly introduce a widely used diffusion model known as *Independent Cascade* [21]. In this model, each edge (u, v) is associated with a propagation probability $p_{u,v}$ denoting the probability that v will be influenced by u . Initially, a set of seed nodes S are activated while all the other nodes are inactive. When a node u first becomes active, it is given a single chance to activate each inactive out-neighbor v with a probability $p_{u,v}$. The diffusion process stops when no more activation can be made.

Let $\sigma(S)$ denote the expected number of nodes activated by the diffusion process starting with a set of seed nodes S . $\sigma(S)$ is known as the *influence spread* of the seed set S .

3.2 Influence Spread Estimation

The computational complexity of the exact influence spread $\sigma(S)$ for a seed set S is proved to be #P-hard for several diffusion models including *Independent Cascade* [6, 7]. As a result, various sampling methods have been proposed for unbiased estimation of the influence spread. The RIS method proposed by Borgs et al. [3]

substantially improved the efficiency to estimate the influence spread compared to the naive Monte-Carlo simulation method. Thus, we adopt RIS for influence spread estimation.

Definition 1 ([3]) A random reverse reachable (RR) set R for a graph G is generated by the following steps:

1. Select a random node $v \in V$.
2. Sample a random graph g from G according to the diffusion model.
3. Take the set of nodes in g that can reach v as R .

For example, under the Independent Cascade model, a random RR set R on G can be constructed as follows:

1. Select a node $v \in V$ uniformly at random.
2. Starting from v , perform a stochastic breadth first search (BFS) following the incoming edges of each node. Specifically, for each node u encountered in the BFS, we examine the in-neighbors of u . For each in-neighbor w , we allow the BFS to traverse to w from u with probability $p_{w,u}$ (if w has not been traversed before).
3. Insert all the nodes traversed during the stochastic BFS into the RR set R .

Random RR sets have the following property [3].

Lemma 1 ([3]) *Given a seed set $S \subseteq V$, for a random RR set R , we have*

$$\sigma(S) = n \cdot \Pr[S \cap R \neq \emptyset], \quad (1)$$

where $n = |V|$ is the total number of nodes in the graph G .

According to Lemma 1, the influence spread of a seed set S is proportional to the probability that S intersects with a random RR set R . Thus, to estimate influence spread, we can generate a large number of RR sets \mathcal{R} . Given any seed set S , we can compute the number of RR sets in \mathcal{R} that intersect with S (denoted by $A(\mathcal{R}, S)$) and estimate the influence spread of S by $\frac{n}{|\mathcal{R}|} \cdot A(\mathcal{R}, S)$.

4 Divergence Function & Optimal Price Profile

In this section, we define a divergence function to measure the effectiveness of a price profile and derive an optimal price profile in terms of the function value. For ease of reference, we list the key notations used in this paper in Table 1.

Table 1 Frequently used notations.

Notations	Description
$G = (V, E)$	A graph G with a set V of nodes and a set E of edges
$\sigma(S)$	Influence spread of seed set S
R	A random RR set
\mathcal{R}	A set of RR sets where $ \mathcal{R} = \theta$
$A(\mathcal{R}, S)$	Number of RR sets in \mathcal{R} intersecting S
C	Candidate node set consisting of n_c nodes
n	Total number of nodes in V , i.e., $n = V $
n_c	Number of nodes in C , i.e., $n_c = C $
n_r	Number of nodes in $C \cap R$
p_i	Price of node $s_i \in C$
p_i°	Optimal price of node s_i to the pricing problem
p_i^*	Optimal price of node s_i to the budgeted pricing problem
b	Total price of the nodes in C

4.1 Divergence Function

The seed users generate revenue from the seed purchase of the advertiser for initiating the campaigns. The influence spread, on the other hand, is the reward gained by the advertiser in the campaigns. Thus, it is important to make sure that the influence spread is worth the cost of seed purchase. In this way, the prices set for seed purchase can not only minimize the regret in deriving the revenue for seeds but also give the advertiser a more predictable return for its purchase. Therefore, our objective of pricing is to match the price of any seed set with the expected marketing value of the seed set as closely as possible.

Consider a candidate node set C consisting of n_c nodes $\{s_1, s_2, \dots, s_{n_c}\}$ offered by the OSN provider for the advertisers to choose seeds. Let p_i be the price of the node s_i in C . We refer to $\langle p_1, p_2, \dots, p_{n_c} \rangle$ as the *price profile*. For any seed set $S \subseteq C$, the total price of the nodes in S is $\sum_{s_i \in S} p_i$, and the influence spread of S is $\sigma(S)$. Let c represent the revenue for influencing one user (e.g., purchasing a promoted product). Then, $c \cdot \sigma(S)$ is the expected market value for the seed set S . Thus, the divergence between the price and the influence spread can be characterized by $(c \cdot \sigma(S) - \sum_{s_i \in S} p_i)^2$. Since the advertisers can choose any subset of the nodes in C to initiate campaigns based on their preferences, we assume that all the subsets of C are equally likely to be chosen as the seed set. Therefore, the expected divergence between the price and the market value of a randomly chosen seed set is given by

$$\frac{1}{2^{n_c}} \sum_{S \subseteq C} \left(c \cdot \sigma(S) - \sum_{s_i \in S} p_i \right)^2 = \frac{c^2}{2^{n_c}} \sum_{S \subseteq C} \left(\sigma(S) - \sum_{s_i \in S} \frac{p_i}{c} \right)^2.$$

We aim to find a price profile to minimize the divergence function. Here, $\frac{p_i}{c}$ can be understood as the normalized individual node price. Without loss of generality, we can simply assume that $c = 1$. In the rest of this paper, we focus on the following divergence function

$$f(p_1, p_2, \dots, p_{n_c}) := \frac{1}{2^{n_c}} \sum_{S \subseteq C} \left(\sigma(S) - \sum_{s_i \in S} p_i \right)^2. \quad (2)$$

We then formally define the pricing problem that aims to minimize the divergence function defined in (2) as follows.

Definition 2 The pricing problem is to find a (non-negative) price profile for a set C of candidate nodes so that the divergence function is minimized, i.e.,

$$\arg \min_{\forall i, p_i \geq 0} f(p_1, p_2, \dots, p_{n_c}).$$

It is easy to verify that the optimal price profile for a general c value can be obtained by simply scaling the solution to the above problem by a multiplicative factor of c .

4.2 Optimal Price Profile

To solve the pricing problem, we first reformulate the divergence function.

Lemma 2 *Let*

$$g(p_i) := \frac{p_i^2}{4} - \frac{p_i}{2^{n_c-1}} \sum_{S \subseteq C \setminus \{s_i\}} \sigma(S \cup \{s_i\}). \quad (3)$$

Then, we have

$$f(p_1, p_2, \dots, p_{n_c}) = \frac{1}{2^{n_c}} \sum_{S \subseteq C} \sigma(S)^2 + \frac{b^2}{4} + \sum_{i=1}^{n_c} g(p_i), \quad (4)$$

where $b := \sum_{i=1}^{n_c} p_i$.

Proof We prove it by induction. When $n_c = 1$, we have

$$g(p_1) = \frac{p_1^2}{4} - p_1 \sigma(\{s_1\}).$$

Meanwhile, by definition, we have

$$f(p_1) = \frac{(\sigma(\{s_1\}) - p_1)^2}{2} = \frac{\sigma(\{s_1\})^2}{2} + \frac{p_1^2}{2} - p_1 \sigma(\{s_1\}).$$

Thus, $f(p_1) = \frac{\sigma(\{s_1\})^2}{2} + \frac{p_1^2}{4} + g(p_1)$, which indicates that (4) holds when $n_c = 1$.

Suppose that (4) holds when $n_c = N$ for an integer $N \geq 1$. In what follows, we will show that (4) holds when $n_c = N + 1$.

For any $i \in [2, n_c]$, let $\bar{\sigma}(S) := \sigma(S \cup \{s_1\}) - p_1$,

$$\bar{g}(p_i) := \frac{p_i^2}{4} - \frac{p_i}{2^{N-1}} \sum_{S \subseteq C \setminus \{s_1, s_i\}} \bar{\sigma}(S \cup \{s_i\}),$$

$$\text{and } \hat{g}(p_i) := \frac{p_i^2}{4} - \frac{p_i}{2^{N-1}} \sum_{S \subseteq C \setminus \{s_1, s_i\}} \sigma(S \cup \{s_i\}).$$

In addition, let

$$\bar{f} := \frac{1}{2^N} \sum_{S \subseteq C \setminus \{s_1\}} \left((\bar{\sigma}(S) - \sum_{s_i \in S} p_i)^2 + (\sigma(S) - \sum_{s_i \in S} p_i)^2 \right).$$

For any node set S , the node s_1 satisfies either $s_1 \in S$ or $s_1 \notin S$. Thus, we have $\bar{f} = 2f(p_1, p_2, \dots, p_{n_c})$.

Let $\bar{b} := \sum_{i=2}^{n_c} p_i = b - p_1$. Since $|C \setminus \{s_1\}| = n_c - 1 = N$, according to the hypothesis, we have

$$\bar{f} = \sum_{S \subseteq C \setminus \{s_1\}} \frac{\bar{\sigma}(S)^2 + \sigma(S)^2}{2^N} + \frac{\bar{b}^2}{2} + \sum_{i=2}^{n_c} (\bar{g}(p_i) + \hat{g}(p_i)).$$

For the first part, we have

$$\begin{aligned} & \frac{1}{2^N} \sum_{S \subseteq C \setminus \{s_1\}} (\bar{\sigma}(S)^2 + \sigma(S)^2) \\ &= \frac{1}{2^N} \sum_{S \subseteq C \setminus \{s_1\}} \left((\sigma(S \cup \{s_1\}) - p_1)^2 + \sigma(S)^2 \right) \\ &= \frac{1}{2^N} \sum_{S \subseteq C} \sigma(S)^2 + p_1^2 - \frac{p_1}{2^{N-1}} \sum_{S \subseteq C \setminus \{s_1\}} \sigma(S \cup \{s_1\}) \\ &= \frac{1}{2^N} \sum_{S \subseteq C} \sigma(S)^2 + 2g(p_1) + \frac{p_1^2}{2}. \end{aligned}$$

For the third part, we have

$$\begin{aligned} & \bar{g}(p_i) + \hat{g}(p_i) \\ &= \frac{p_i^2}{2} - \frac{p_i}{2^{N-1}} \sum_{S \subseteq C \setminus \{s_1, s_i\}} \left(\bar{\sigma}(S \cup \{s_i\}) + \sigma(S \cup \{s_i\}) \right) \\ &= \frac{p_i^2}{2} - \frac{p_i}{2^{N-1}} \sum_{S \subseteq C \setminus \{s_i\}} \sigma(S \cup \{s_i\}) + p_1 p_i \\ &= 2g(p_i) + p_1 p_i. \end{aligned}$$

Therefore, we have

$$\begin{aligned} \bar{f} &= \frac{1}{2^N} \sum_{S \subseteq C} \sigma(S)^2 + 2g(p_1) + \frac{p_1^2}{2} \\ &\quad + \frac{\bar{b}^2}{2} + \sum_{i=2}^{n_c} (2g(p_i) + p_1 p_i) \\ &= \frac{1}{2^N} \sum_{S \subseteq C} \sigma(S)^2 + \frac{p_1^2 + \bar{b}^2 + 2p_1 \bar{b}}{2} + 2 \sum_{i=1}^{n_c} g(p_i) \\ &= \frac{1}{2^N} \sum_{S \subseteq C} \sigma(S)^2 + \frac{b^2}{2} + 2 \sum_{i=1}^{n_c} g(p_i). \end{aligned}$$

This completes the proof. \square

Based on Lemma 2, we can derive the optimal price profile as follows.

Theorem 1 For each $i = 1, 2, \dots, n_c$, let

$$p_i^\circ := \frac{\sum_{S \subseteq C} \sigma(S) \left(\mathbf{1}_{\{s_i \in S\}} - \frac{|S|}{n_c+1} \right)}{2^{n_c-2}}, \quad (5)$$

where $\mathbf{1}_{\{s_i \in S\}}$ is an indicator function such that $\mathbf{1}_{\{s_i \in S\}} = 1$ if $s_i \in S$ and $\mathbf{1}_{\{s_i \in S\}} = 0$ otherwise. Then, $\langle p_1^\circ, p_2^\circ, \dots, p_{n_c}^\circ \rangle$ is an optimal price profile to the pricing problem.

Proof According to Lemma 2, for each i , taking the partial derivative of $f(p_1, p_2, \dots, p_{n_c})$ with respect to p_i gives

$$\frac{df(p_1, p_2, \dots, p_{n_c})}{dp_i} = \frac{b + p_i}{2} - c_i \triangleq 0, \quad (6)$$

where $b = \sum_{i=1}^{n_c} p_i$ and $c_i = \frac{\sum_{S \subseteq C \setminus \{s_i\}} \sigma(S \cup \{s_i\})}{2^{n_c-1}}$. It is trivial to see that the profile $\langle p_1, p_2, \dots, p_{n_c} \rangle$ satisfying (6) gives a lower bound $f(p_1, p_2, \dots, p_{n_c})$ to the pricing problem. Summing (6) from $i = 1$ to $i = n_c$ gives

$$\frac{(n_c + 1)b}{2} - \sum_{i=1}^{n_c} c_i = 0.$$

As a result, we have

$$\begin{aligned} p_i &= 2c_i - b = 2c_i - \frac{2}{n_c + 1} \cdot \sum_{i=1}^{n_c} c_i \\ &= \sum_{S \subseteq C \setminus \{s_i\}} \frac{\sigma(S \cup \{s_i\})}{2^{n_c-2}} - \frac{2}{n_c + 1} \cdot \sum_{S \subseteq C} \frac{\sigma(S) \cdot |S|}{2^{n_c-1}} \\ &= \sum_{S \subseteq C} \frac{\sigma(S) \left(\mathbf{1}_{\{s_i \in S\}} - \frac{|S|}{n_c+1} \right)}{2^{n_c-2}} \\ &= p_i^\circ. \end{aligned}$$

Then, it suffices to show that p_i° is non-negative.

Next, we utilize the RIS method [3] to show the non-negativity of p_i° . By Lemma 1, for a random RR set R , we have $\sigma(S) = n \cdot \Pr[S \cap R \neq \emptyset] = n \cdot \mathbb{E}[\mathbf{1}_{\{S \cap R \neq \emptyset\}}]$. Thus,

$$p_i^\circ = n \cdot \mathbb{E} \left[\frac{\sum_{S \subseteq C} \left(\mathbf{1}_{\{S \cap R \neq \emptyset\}} \cdot \left(\mathbf{1}_{\{s_i \in S\}} - \frac{|S|}{n_c+1} \right) \right)}{2^{n_c-2}} \right].$$

Denote by \tilde{A}_i° an estimate of p_i° using a random RR set R , i.e.,

$$\frac{\tilde{A}_i^\circ}{n} = \frac{\sum_{S \subseteq C} \left(\mathbf{1}_{\{S \cap R \neq \emptyset\}} \cdot \left(\mathbf{1}_{\{s_i \in S\}} - \frac{|S|}{n_c+1} \right) \right)}{2^{n_c-2}}.$$

Then p_i° can be written as $\mathbb{E}[\tilde{A}_i^\circ]$. Next, we show that for any RR set R , $\tilde{A}_i^\circ \geq 0$, which implies $p_i^\circ \geq 0$.

Given an RR set R , let $n_r := |C \cap R|$. Then, we have

$$\begin{aligned} \sum_{S \subseteq C} \left(\mathbf{1}_{\{S \cap R \neq \emptyset\}} \cdot |S| \right) &= \sum_{S \subseteq C} |S| - \sum_{S \subseteq (C \setminus R)} |S| \\ &= n_c \cdot 2^{n_c-1} - (n_c - n_r) \cdot 2^{n_c - n_r - 1}. \end{aligned}$$

In addition, if $s_i \in R$, we have

$$\sum_{S \subseteq C} \left(\mathbf{1}_{\{S \cap R \neq \emptyset\}} \cdot \mathbf{1}_{\{s_i \in S\}} \right) = 2^{n_c-1},$$

and thus

$$\begin{aligned} \frac{\tilde{A}_i^\circ}{n} &= \frac{2^{n_c-1} - \frac{1}{n_c+1} \cdot (n_c \cdot 2^{n_c-1} - (n_c - n_r) \cdot 2^{n_c - n_r - 1})}{2^{n_c-2}} \\ &= \frac{2}{n_c + 1} + \frac{n_c - n_r}{n_c + 1} \cdot 2^{1-n_r}. \end{aligned}$$

On the other hand, if $s_i \notin R$, we have

$$\begin{aligned} \sum_{S \subseteq C} \left(\mathbf{1}_{\{S \cap R \neq \emptyset\}} \cdot \mathbf{1}_{\{s_i \in S\}} \right) &= \sum_{S \subseteq C \setminus \{s_i\}} \mathbf{1}_{\{S \cap R \neq \emptyset\}} \\ &= \sum_{S \subseteq C \setminus \{s_i\}} 1 - \sum_{S \subseteq C \setminus (R \cup \{s_i\})} 1 \\ &= 2^{n_c-1} - 2^{n_c - n_r - 1}, \end{aligned}$$

and thus

$$\frac{\tilde{A}_i^\circ}{n} = \frac{2}{n_c + 1} - \frac{n_r + 1}{n_c + 1} \cdot 2^{1-n_r}.$$

To summarize,

$$\frac{\tilde{A}_i^\circ}{n} = \begin{cases} \frac{2}{n_c+1} + \frac{n_c - n_r}{n_c+1} \cdot 2^{1-n_r}, & \text{if } s_i \in R, \\ \frac{2}{n_c+1} - \frac{n_r+1}{n_c+1} \cdot 2^{1-n_r}, & \text{otherwise.} \end{cases} \quad (7)$$

It is easy to verify that

$$\begin{aligned} \frac{\tilde{A}_i^\circ}{n} &\geq \frac{2}{n_c + 1} - \frac{n_r + 1}{n_c + 1} \cdot 2^{1-n_r} \\ &= \frac{2 \cdot (1 - (n_r + 1) \cdot 2^{-n_r})}{n_c + 1} \geq 0, \end{aligned}$$

since $2^{n_r} - n_r - 1 \geq 0$. This completes the proof. \square

Hardness Analysis. We present a polynomial closed-form solution in Theorem 1 to the pricing problem defined in Definition 2. As can be seen from (5), for each $1 \leq i \leq n_c$, computing p_i° in our pricing problem can be reduced to calculating $\sigma(S)$ for any $S \subseteq C$, which is #P-hard for several diffusion models including the Independent Cascade and Linear Threshold models [6, 7]. Thus, our pricing problem is also #P-hard under these diffusion models.

4.3 Pricing Nodes with a Budget Constraint

Furthermore, we also study the budgeted pricing problem when the total price of all nodes in the candidate set C is equal to a given value b . The following theorem gives the optimal price profile to the budgeted pricing problem.

Theorem 2 *Given any price budget b , let*

$$c_i := \frac{\sum_{S \subseteq C \setminus \{s_i\}} \sigma(S \cup \{s_i\})}{2^{n_c-1}}, \quad \forall 1 \leq i \leq n_c, \quad (8)$$

$$\lambda \text{ be the root of } \sum_{i=1}^{n_c} \max\{0, 2(\lambda + c_i)\} = b, \quad (9)$$

$$\text{and } p_i^* := \max\{0, 2(\lambda + c_i)\}, \quad \forall 1 \leq i \leq n_c. \quad (10)$$

Then, $\langle p_1^, p_2^*, \dots, p_{n_c}^* \rangle$ is an optimal price profile to the budgeted pricing problem.*

Proof According to Karush-Kuhn-Tucker conditions, the optimal solution $\langle p_1^*, p_2^*, \dots, p_{n_c}^* \rangle$ satisfies that

$$\forall i, \frac{p_i^*}{2} - c_i - \lambda - \lambda_i = 0, \quad (11)$$

$$\forall i, \lambda_i p_i^* = 0, \quad (12)$$

$$\sum_{i=1}^{n_c} p_i^* = b, \quad (13)$$

$$\forall i, p_i^* \geq 0, \quad (14)$$

$$\forall i, \lambda_i \geq 0. \quad (15)$$

In the above, (11) represents stationarity, (12) shows complementary slackness, (13) and (14) ensure primal feasibility, and (15) ensures dual feasibility.

If $\lambda + c_i \leq 0$, by (11), we have $\frac{p_i^*}{2} \leq \lambda_i$. By (12), (14), and (15), we have $p_i^* = 0$. Similarly, if $\lambda + c_i > 0$, we have $\lambda_i = 0$, which indicates that $p_i^* = 2(\lambda + c_i)$. Therefore,

$$p_i^* = \max\{0, 2(\lambda + c_i)\}. \quad (16)$$

Then, by (13), λ is the solution for $\sum_{i=1}^{n_c} \max\{0, 2(\lambda + c_i)\} = b$, which completes the proof. \square

Theorem 2 states the optimal solution, where λ can be obtained via water-filling. Specifically, without loss of generality, we assume that $c_1 \geq c_2 \geq \dots \geq c_{n_c} > c_{n_c+1} = -\infty$. Then, we can find a unique $j \in [1, n_c]$ such that $\sum_{i=1}^{j-1} 2(c_i - c_j) < b$ and $\sum_{i=1}^j 2(c_i - c_{j+1}) \geq b$. This implies that $-c_j < \lambda \leq -c_{j+1}$. Hence, λ is the root of $\sum_{i=1}^j 2(\lambda + c_i) = b$.

Corollary 1 *Given any price budget b , for each $i = 1, 2, \dots, n_c$, let*

$$p_i^* := \frac{b}{n_c} + \frac{\sum_{S \subseteq C} \sigma(S) (\mathbf{1}_{\{s_i \in S\}} - \frac{|S|}{n_c})}{2^{n_c-2}}, \quad (17)$$

where $\mathbf{1}_{\{s_i \in S\}}$ is an indicator function such that $\mathbf{1}_{\{s_i \in S\}} = 1$ if $s_i \in S$ and $\mathbf{1}_{\{s_i \in S\}} = 0$ otherwise. If the budget b is no less than the influence spread $\sigma(C)$ of the candidate set C , i.e., $b \geq \sigma(C)$, $\langle p_1^, p_2^*, \dots, p_{n_c}^* \rangle$ is an optimal price profile to the budgeted pricing problem.*

Proof Let λ be the value such that

$$\sum_{i=1}^{n_c} 2(\lambda + c_i) = b.$$

Then,

$$\lambda = \frac{b}{2n_c} - \frac{1}{n_c} \sum_{i=1}^{n_c} c_i = \frac{b}{2n_c} - \frac{\sum_{S \subseteq C} (\sigma(S) \cdot |S|)}{n_c \cdot 2^{n_c-1}},$$

and

$$2(\lambda + c_i) = \frac{b}{n_c} + \frac{\sum_{S \subseteq C} \sigma(S) (\mathbf{1}_{\{s_i \in S\}} - \frac{|S|}{n_c})}{2^{n_c-2}}.$$

Again, we utilize the RIS method [3] to show that $2(\lambda + c_i) \geq 0$ to ensure that for each $1 \leq i \leq n_c$, p_i^* given in (17) is non-negative when $b \geq \sigma(C)$.

Specifically, define

$$A_i := 2(\lambda + c_i) - \frac{b}{n_c} = \frac{\sum_{S \subseteq C} \sigma(S) (\mathbf{1}_{\{s_i \in S\}} - \frac{|S|}{n_c})}{2^{n_c-2}}.$$

Let \tilde{A}_i be an estimate of A_i using a random RR set R , i.e.,

$$\frac{\tilde{A}_i}{n} = \frac{\sum_{S \subseteq C} (\mathbf{1}_{\{S \cap R \neq \emptyset\}} \cdot (\mathbf{1}_{\{s_i \in S\}} - \frac{|S|}{n_c}))}{2^{n_c-2}}.$$

Analogous to the proof of Theorem 1, it is easy to get that

$$\frac{\tilde{A}_i}{n} = \begin{cases} (1 - \frac{n_r}{n_c}) \cdot 2^{1-n_r}, & \text{if } s_i \in R, \\ -\frac{n_r}{n_c} \cdot 2^{1-n_r}, & \text{otherwise.} \end{cases} \quad (18)$$

It is trivial to see that

$$\frac{\tilde{A}_i}{n} \geq -\frac{n_r}{n_c} \cdot 2^{1-n_r} \geq -\frac{\mathbf{1}_{\{C \cap R \neq \emptyset\}}}{n_c}.$$

By Lemma 1, we have

$$\begin{aligned} 2(\lambda + c_i) &= \frac{b}{n_c} + A_i = \frac{b}{n_c} + \mathbb{E}[\tilde{A}_i] \\ &\geq \frac{b}{n_c} - \frac{n \cdot \Pr[C \cap R \neq \emptyset]}{n_c} = \frac{b - \sigma(C)}{n_c} \geq 0. \end{aligned}$$

This completes the proof. \square

Discussion. We also study the variants of the budgeted pricing problem when the total price of the candidate nodes is no larger than b or no less than b , which can be solved by combining the solution to the general pricing problem and the solution to the budgeted pricing problem. Specifically, under the assumption that the given budget b satisfies $b \geq \sigma(C)$, by Corollary 1, we have $\frac{dp_i^*}{db} = \frac{1}{n_c}$. Thus,

$$\begin{aligned} & \frac{df(p_1^*, p_2^*, \dots, p_{n_c}^*)}{db} \\ &= \frac{b}{2} + \sum_{i=1}^{n_c} \frac{dg(p_i^*)}{db} \\ &= \frac{b}{2} + \sum_{i=1}^{n_c} \frac{1}{4} \cdot \left(2p_i^* \cdot \frac{1}{n_c}\right) - \frac{1}{n_c} \frac{\sum_{S \subseteq C \setminus \{s_i\}} \sigma(S \cup \{s_i\})}{2^{n_c-1}} \\ &= \frac{b}{2} \left(\frac{1}{n_c} + 1\right) - \frac{\sum_{S \subseteq C} \sigma(S) \cdot |S|}{n_c \cdot 2^{n_c-1}}. \end{aligned}$$

Let $\frac{df(p_1^*, p_2^*, \dots, p_{n_c}^*)}{db} = 0$. Then, the divergence function achieves its minimal value at $b = \frac{1}{n_c+1} \cdot \frac{\sum_{S \subseteq C} \sigma(S) \cdot |S|}{2^{n_c-2}}$ while satisfying that p_i^* is non-negative for each $1 \leq i \leq n_c$. Let $b^\circ = \frac{1}{n_c+1} \cdot \frac{\sum_{S \subseteq C} \sigma(S) \cdot |S|}{2^{n_c-2}}$. Take the ‘‘less or equal to b ’’ pricing problem as an example. If the given budget b satisfies $b \leq b^\circ$, since the divergence function decreases with the growth of b , the minimal divergence value is achieved at value b and the optimal pricing solution is given by Theorem 2. Otherwise, when the given budget b satisfies $b \geq b^\circ$, the divergence function achieves its minimal value at value b° and the optimal pricing solution is given by Theorem 1. The solution to the pricing problem where the total price of the candidate nodes is larger or equal to b can be derived similarly.

4.4 Discussion on Privacy Issues

Privacy protection is critical for both OSN providers and influencers. On one hand, our pricing mechanism intrinsically protects the privacy of OSN providers since they do not need to unveil the network structures. On the other hand, as most information posts are publicly available and designed to attract followers on platforms such as Instagram and TikTok, many influencers are willing to monetize their public influence powers. To minimize the ethical issues, OSN providers can first confirm influencers’ willingness of engagement in marketing campaigns and then post their prices for marketing campaigns. Furthermore, to protect the privacy of the candidate seeds, the prices can be posted anonymously such that the personal information can be protected.

5 Estimation of Node Prices

In this section, we focus on the estimation of the node prices $p_1^\circ, p_2^\circ, \dots, p_{n_c}^\circ$ in the optimal price profile to the pricing problem defined in Definition 2, while the estimation of $p_1^*, p_2^*, \dots, p_{n_c}^*$ to the budgeted pricing problem is similar.

5.1 Unbiased Estimator via RR Sets

When θ RR sets are generated, according to (7), an RR set R contributes to the estimator \tilde{A}_i° by an additive factor of $\Delta(\tilde{A}_i^\circ, R)$, which is given as follows.

$$\Delta(\tilde{A}_i^\circ, R) = \begin{cases} \frac{n}{\theta} \cdot \left(\frac{2}{n_c+1} + \frac{n_c-n_r}{n_c+1} \cdot 2^{1-n_r}\right), & \text{if } s_i \in R, \\ \frac{n}{\theta} \cdot \left(\frac{2}{n_c+1} - \frac{n_r+1}{n_c+1} \cdot 2^{1-n_r}\right), & \text{otherwise.} \end{cases}$$

Recall that $\Delta(\tilde{A}_i^\circ, R)$ is ensured to be non-negative. Let $R_1, R_2, \dots, R_\theta$ be a sequence of random RR sets. Let $X_{i,j}$ be a random variable defined as

$$X_{i,j} = \begin{cases} 0 & \text{if } R_j \cap C = \emptyset, \\ \frac{2}{n_c+1} + \frac{n_c-n_r}{n_c+1} \cdot 2^{1-n_r} & \text{if } s_i \in R_j \cap C, \\ \frac{2}{n_c+1} - \frac{n_r+1}{n_c+1} \cdot 2^{1-n_r} & \text{otherwise,} \end{cases} \quad (19)$$

where $n_c = |C|$ and $n_r = |R_j \cap C|$. It is easy to verify that $0 \leq X_{i,j} \leq 1$. Then, by definition,

$$\Delta(\tilde{A}_i^\circ, R_j) = \frac{n}{\theta} \cdot X_{i,j}.$$

Thus, we have

$$p_i^\circ = \frac{n}{\theta} \cdot \mathbb{E} \left[\sum_{j=1}^{\theta} X_{i,j} \right].$$

To use $\frac{n}{\theta} \cdot \sum_{j=1}^{\theta} X_{i,j}$ as an estimator of p_i° , we need θ to be large enough in order to ensure that $\sum_{j=1}^{\theta} X_{i,j}$ does not deviate significantly from its expectation.

5.2 Stopping Rule Algorithm

We generalize the stopping rule algorithm [11] to get an (ε, δ) -approximation of p_i° . Similar to the work [27], we also use the martingale-based concentration bounds [9] to tighten the threshold setting in the stopping rule algorithm. The key differences of our algorithm from previous work [11, 27] are as follows:

- We invent a tighter threshold setting than previous work [11, 27] to improve the efficiency of the stopping rule algorithm.
- We construct an algorithm to estimate all the p_i° for every $i = 1, 2, \dots, n_c$ simultaneously.

Algorithm 1: Stopping Rule Algorithm

Input: number of nodes n , accuracy parameters $\varepsilon, \delta \in (0, 1)$;
Output: an (ε, δ) -approximation \tilde{A}_i° of p_i° for each $i \leq n_c$;

- 1 $\Upsilon \leftarrow (1 + \varepsilon)(1 + (2 + \frac{2}{3}\varepsilon) \ln(\frac{2}{\delta}) \frac{1}{\varepsilon^2})$ and $\theta \leftarrow 0$;
- 2 **foreach** node $s_i \in C$ **do**
- 3 initialize $S_i \leftarrow 0$;
- 4 **while** $\min S_i < \Upsilon$ **do**
- 5 $\theta \leftarrow \theta + 1$;
- 6 generate RR set R_θ ;
- 7 **foreach** node $s_i \in C$ **do**
- 8 **if** $S_i < \Upsilon$ **then**
- 9 $S_i \leftarrow S_i + X_{i,\theta}$, where $X_{i,\theta}$ is based on (19);
- 10 $\theta_i \leftarrow \theta$;
- 11 **return** $\{\tilde{A}_i^\circ = n \cdot \frac{\Upsilon}{\theta_i} : 1 \leq i \leq n_c\}$;

To obtain an (ε, δ) -approximation of the mean of a random variable, the stopping rule algorithm first computes a threshold Υ and then continuously generates samples according to the distribution until their sum exceeds Υ . Finally, the stopping rule algorithm returns the average of these samples as the estimate. The basic stopping rule algorithm can estimate the mean of only one random variable. In our pricing problem, we need to estimate all the values p_i° for every $i = 1, 2, \dots, n_c$ in order to derive the optimal price profile. Estimating each p_i° by a separate invocation of the stopping rule algorithm can result in generating an unnecessarily large number of samples (RR sets). In the following, we construct a stopping rule algorithm to estimate all the values p_i° for every $i = 1, 2, \dots, n_c$ simultaneously.

Algorithm 1 shows the details. The algorithm first calculates the threshold Υ based on the required approximation parameters ε and δ (line 1). After that, samples are generated and aggregated until the minimum sum among all the nodes s_i 's exceeds Υ (Lines 4–10). The number of samples θ_i is recorded for each node s_i until its sum S_i exceeds Υ (Line 10). Finally, the estimate \tilde{A}_i° of the price p_i° for each node s_i is returned (Line 11).

5.3 Theoretical Analysis

5.3.1 A Tighter Threshold Setting

The original stopping rule algorithm [11] sets the threshold Υ as

$$\begin{aligned} \Upsilon_D &= 1 + 4(1 + \varepsilon)(e - 2) \ln\left(\frac{2}{\delta}\right) \frac{1}{\varepsilon^2} \\ &> 1 + 2.87(1 + \varepsilon) \ln\left(\frac{2}{\delta}\right) \frac{1}{\varepsilon^2}. \end{aligned}$$

In Algorithm 1, we set the threshold Υ as

$$\begin{aligned} \Upsilon &= (1 + \varepsilon)(1 + (2 + \frac{2}{3}\varepsilon) \ln\left(\frac{2}{\delta}\right) \frac{1}{\varepsilon^2}) \\ &< (1 + \varepsilon)(1 + 2.67 \ln\left(\frac{2}{\delta}\right) \frac{1}{\varepsilon^2}), \end{aligned}$$

since $0 < \varepsilon \leq 1$. Hence, the Υ setting in our algorithm is tighter than that in [11] when $0.2(1 + \varepsilon) \ln(\frac{2}{\delta}) \frac{1}{\varepsilon^2} > \varepsilon$, which holds when $\varepsilon \leq 0.5 < \sqrt[3]{0.2 \cdot \ln 2}$.

Similar to our algorithm, the stopping rule algorithm proposed by Nguyen et al. [27] also uses the martingale-based concentration bounds [9] to set the threshold Υ . Since $0 \leq X_{i,j} \leq 1$ in our problem, applying the algorithm in [27], the threshold Υ is set as

$$\Upsilon_N = (1 + \varepsilon)(2 + \frac{2}{3}\varepsilon') \ln\left(\frac{2}{\delta}\right) \frac{1}{\varepsilon'^2},$$

where $\varepsilon' = \varepsilon(1 - \frac{\varepsilon}{(2 + \frac{2}{3}\varepsilon) \ln(\frac{2}{\delta})}) < \varepsilon$. In Algorithm 1, we replace ε' with ε in the setting of Υ and add an additive factor of $1 + \varepsilon$. Next, we show that the Υ value in our algorithm is smaller than that in Nguyen et al.'s work [27]. To prove

$$(1 + \varepsilon)(2 + \frac{2}{3}\varepsilon') \ln\left(\frac{2}{\delta}\right) \frac{1}{\varepsilon'^2} > (1 + \varepsilon)(1 + (2 + \frac{2}{3}\varepsilon) \ln\left(\frac{2}{\delta}\right) \frac{1}{\varepsilon^2}),$$

it is equivalent to show that

$$(1 + \frac{1}{3}\varepsilon') \frac{1}{\varepsilon'^2} - (1 + \frac{1}{3}\varepsilon) \frac{1}{\varepsilon^2} > \frac{1}{2 \ln(\frac{2}{\delta})}. \quad (20)$$

Let $\alpha := (2 + \frac{2}{3}\varepsilon) \ln(\frac{2}{\delta})$ and $\beta := \frac{\varepsilon'}{\varepsilon} = 1 - \frac{\varepsilon}{\alpha}$. We have

$$\begin{aligned} &(1 + \frac{1}{3}\varepsilon') \frac{1}{\varepsilon'^2} - (1 + \frac{1}{3}\varepsilon) \frac{1}{\varepsilon^2} \\ &= (1 + \frac{1}{3}\beta\varepsilon) \frac{1}{\beta^2\varepsilon^2} - (1 + \frac{1}{3}\varepsilon) \frac{1}{\varepsilon^2} \\ &= \left(\frac{1 + \beta}{\beta\varepsilon} + \frac{1}{3}\right) \cdot \frac{1 - \beta}{\beta\varepsilon} \\ &> \frac{2 - \frac{\varepsilon}{\alpha}}{\varepsilon(1 - \frac{\varepsilon}{\alpha})} \cdot \frac{\frac{\varepsilon}{\alpha}}{\varepsilon(1 - \frac{\varepsilon}{\alpha})} \\ &= \frac{2\alpha - \varepsilon}{\varepsilon(\alpha - \varepsilon)^2} \\ &> \frac{2}{\varepsilon(\alpha - \varepsilon)}. \end{aligned}$$

Since $4 \ln(\frac{2}{\delta}) > \frac{8}{3} \ln(\frac{2}{\delta}) \geq \alpha \geq \alpha - \varepsilon \geq \varepsilon(\alpha - \varepsilon)$, we have

$$\frac{2}{\varepsilon(\alpha - \varepsilon)} > \frac{1}{2 \ln(\frac{2}{\delta})}.$$

Thus, (20) holds and the Υ setting in our algorithm is tighter. So, our algorithm would generate less samples than that in Nguyen et al.'s work [27].

5.3.2 Theoretical Guarantee

In the following, we prove that our \mathcal{Y} setting can guarantee an (ε, δ) -approximation of estimation. The proof is similar in spirit to the original stopping rule algorithm [11], but we make use of the martingale-based concentration bounds [9] in the derivation.

Definition 3 ([9]) A sequence of random variables Y_1, Y_2, \dots is a martingale if and only if $\mathbb{E}[|Y_j|] < \infty$ and $\mathbb{E}[Y_j | Y_1, Y_2, \dots, Y_{j-1}] = Y_{j-1}$ for any j .

Since each RR set R_j is generated randomly and independently of all the prior RR sets, we have

$$\mathbb{E}[X_{i,j} | X_{i,1}, X_{i,2}, \dots, X_{i,j-1}] = \mathbb{E}[X_{i,j}] = \frac{p_i^\circ}{n}. \quad (21)$$

Let $\mu_i = \frac{p_i^\circ}{n}$ and $M_{i,j} = \sum_{k=1}^j (X_{i,k} - \mu_i)$. Then, we have $\mathbb{E}[|M_{i,j}|] \leq j < \infty$, and

$$\mathbb{E}[M_{i,j} | M_{i,1}, M_{i,2}, \dots, M_{i,j-1}] = M_{i,j-1}.$$

Therefore, $M_{i,1}, M_{i,2}, \dots, M_{i,\theta}$ is a martingale.

Lemma 3 ([9]) Let Y_1, Y_2, \dots, Y_j be a martingale, such that $Y_1 \leq a$, $Y_k - Y_{k-1} \leq a$ for any $2 \leq k \leq j$, and $\text{Var}[Y_1] + \sum_{k=2}^j \text{Var}[Y_k | Y_1, Y_2, \dots, Y_{k-1}] \leq b$. Then, for any $\eta > 0$,

$$\Pr[Y_j - \mathbb{E}[Y_j] \geq \eta] \leq \exp\left(-\frac{\eta^2}{\frac{2}{3}a\eta + 2b}\right). \quad (22)$$

Since $0 \leq X_{i,j} \leq 1$ for any $1 \leq j \leq \theta$, we have $M_{i,1} = X_{i,1} - \mu_i \leq 1$ and $M_{i,j} - M_{i,j-1} = X_{i,j} - \mu_i \leq 1$ for any $2 \leq j \leq \theta$. Let $\text{Var}[\cdot]$ denote the variance of a random variable. It follows that $\text{Var}[X_{i,j}] = \mathbb{E}[X_{i,j}^2] - \mathbb{E}[X_{i,j}]^2 = \mathbb{E}[X_{i,j}^2] - \mu_i^2 \leq \mathbb{E}[X_{i,j}] - \mu_i^2 \leq \mu_i(1 - \mu_i)$. Hence,

$$\begin{aligned} & \text{Var}[M_{i,1}] + \sum_{j=2}^{\theta} \text{Var}[M_{i,j} | M_{i,1}, M_{i,2}, \dots, M_{i,j-1}] \\ &= \sum_{j=1}^{\theta} \text{Var}[X_{i,j} - \mu_i] \leq \theta \mu_i \cdot (1 - \mu_i) \leq \theta \mu_i. \end{aligned}$$

By Lemma 3, we have $\Pr[M_{i,\theta} \geq \varepsilon \cdot \theta \mu_i] \leq \exp\left(-\frac{\varepsilon^2 \theta \mu_i}{2 + \frac{2}{3}\varepsilon}\right)$.

Similarly, $-M_{i,1}, -M_{i,2}, \dots, -M_{i,\theta}$ is a martingale such that

$$\begin{aligned} & -M_{i,1} \leq \mu_i, \\ & -M_{i,j} + M_{i,j-1} \leq \mu_i, \end{aligned}$$

$$\begin{aligned} & \sum_{j=2}^{\theta} \text{Var}[-M_{i,j} | -M_{i,1}, -M_{i,2}, \dots, -M_{i,j-1}] \\ &+ \text{Var}[-M_{i,1}] \leq \theta \mu_i \cdot (1 - \mu_i). \end{aligned}$$

Then, we have

$$\begin{aligned} \Pr[-M_{i,\theta} \geq \varepsilon \cdot \theta \mu_i] &= \Pr\left[\sum_{j=1}^{\theta} X_{i,j} - \theta \mu_i \leq -\varepsilon \cdot \theta \mu_i\right] \\ &\leq \exp\left(-\frac{\varepsilon^2 \theta^2 \mu_i^2}{\frac{2}{3}\varepsilon \theta \mu_i^2 + 2\theta \mu_i(1 - \mu_i)}\right) \\ &\leq \exp\left(-\frac{\varepsilon^2 \theta \mu_i}{2}\right). \end{aligned}$$

To summarize, we have the following corollary.³

Corollary 2 For any $\varepsilon > 0$,

$$\Pr\left[\sum_{j=1}^{\theta} X_{i,j} - \theta \mu_i \geq \varepsilon \cdot \theta \mu_i\right] \leq \exp\left(-\frac{\varepsilon^2 \theta \mu_i}{2 + \frac{2}{3}\varepsilon}\right), \quad (23)$$

$$\Pr\left[\sum_{j=1}^{\theta} X_{i,j} - \theta \mu_i \leq -\varepsilon \cdot \theta \mu_i\right] \leq \exp\left(-\frac{\varepsilon^2 \theta \mu_i}{2}\right). \quad (24)$$

According to Corollary 2, the following lemma describes the phenomenon that the failure probability is independent of the mean of the tested random variable in the multiplicative-additive error form of martingale-based concentration bounds.

Lemma 4 Let $Z_1 - \mathbb{E}[Z_1], \dots, Z_T - \mathbb{E}[Z_T]$ be a martingale difference sequence such that $Z_j \in [0, 1]$ for each j . Let $\bar{Z} = \frac{1}{T} \sum_{j=1}^T Z_j$. If $\mathbb{E}[Z_j]$ is identical for every j , i.e., $\mathbb{E}[Z_j] = \mathbb{E}[\bar{Z}]$, then,

$$\Pr[\bar{Z} \leq (1 - \varepsilon)\mathbb{E}[\bar{Z}] - \beta] \leq \exp(-2\varepsilon\beta T), \quad (25)$$

$$\Pr[\bar{Z} \geq (1 + \varepsilon)\mathbb{E}[\bar{Z}] + \beta] \leq \exp\left(-\frac{2\varepsilon\beta T}{(1 + \varepsilon/3)^2}\right). \quad (26)$$

Proof By (24) in Corollary 2, we have

$$\begin{aligned} \Pr[\bar{Z} \leq (1 - \varepsilon)\mathbb{E}[\bar{Z}] - \beta] &\leq \exp\left(-\frac{(\varepsilon\mathbb{E}[\bar{Z}] + \beta)^2 T}{2\mathbb{E}[\bar{Z}]}\right) \\ &\leq \exp\left(-\frac{(2\sqrt{\varepsilon\mathbb{E}[\bar{Z}]\beta})^2 T}{2\mathbb{E}[\bar{Z}]}\right) \\ &= \exp(-2\varepsilon\beta T). \end{aligned}$$

Similarly, by (23) in Corollary 2, we have

$$\Pr[\bar{Z} \geq (1 + \varepsilon)\mathbb{E}[\bar{Z}] + \beta] \leq \exp(-h(\lambda)),$$

where $h(\lambda) = \frac{(\lambda^2 T)}{2(\lambda - \beta)/\varepsilon + 2\lambda/3}$ and $\lambda = \varepsilon\mathbb{E}[\bar{Z}] + \beta$. Let

$$\frac{dh(\lambda)}{d\lambda} = \frac{(2\lambda((\lambda - \beta)/\varepsilon + \lambda/3) - (1/\varepsilon + 1/3)\lambda^2)T}{2((\lambda - \beta)/\varepsilon + \lambda/3)^2} \triangleq 0.$$

Thus, $h(\lambda)$ achieves its minimum at $\lambda = \frac{2\beta}{\varepsilon(1/\varepsilon + 1/3)}$ such that $h(\lambda) = \frac{2\varepsilon\beta T}{(1 + \varepsilon/3)^2}$. This completes the proof. \square

³ Tang et al. [38] directly gave the lower tail result without providing the detailed proof. Our analysis is based on Lemma 3 requiring $Y_1 \leq a$ and $Y_k - Y_{k-1} \leq a$, whereas Tang et al. [38] utilized a similar lemma requiring $|Y_1| \leq a$ and $|Y_k - Y_{k-1}| \leq a$, which might be insufficient for deriving the lower tail result.

Based on Corollary 2, we show that Algorithm 1 returns an (ε, δ) -approximate $\tilde{\mu}_i$ of each μ_i ($1 \leq i \leq n_c$).

Theorem 3 Algorithm 1 returns an (ε, δ) -approximate $\tilde{\mu}_i$ of μ_i for each $1 \leq i \leq n_c$, i.e.,

$$\Pr[(1 - \varepsilon)\mu_i \leq \tilde{\mu}_i \leq (1 + \varepsilon)\mu_i] \geq 1 - \delta. \quad (27)$$

Proof Given any i where $1 \leq i \leq n_c$, we will prove the following two probabilistic inequalities:

$$\Pr[\tilde{\mu}_i < (1 - \varepsilon)\mu_i] \leq \frac{\delta}{2}, \quad (28)$$

$$\Pr[\tilde{\mu}_i > (1 + \varepsilon)\mu_i] \leq \frac{\delta}{2}. \quad (29)$$

First, we prove (28). Since $0 \leq X_{i,j} \leq 1$, by the definition of Algorithm 1, when it terminates, we have

$$\Upsilon \leq S_i = \sum_{j=1}^{\theta_i} X_{i,j} \leq \Upsilon + 1.$$

Let $L_1 := \lceil \frac{\Upsilon}{(1-\varepsilon)\mu_i} \rceil$. Then,

$$L_1 \geq \frac{\Upsilon}{(1-\varepsilon)\mu_i},$$

and hence,

$$\frac{\Upsilon}{L_1} \leq (1 - \varepsilon)\mu_i.$$

Since θ_i is an integer, we have

$$\begin{aligned} & \Pr[\tilde{\mu}_i < (1 - \varepsilon)\mu_i] \\ &= \Pr[\Upsilon < (1 - \varepsilon)\mu_i \theta_i] \\ &= \Pr\left[\frac{\Upsilon}{(1 - \varepsilon) \cdot \mu_i} < \theta_i\right] \\ &= \Pr[L_1 \leq \theta_i] \\ &\leq \Pr\left[\sum_{j=1}^{L_1} X_{i,j} \leq \sum_{j=1}^{\theta_i} X_{i,j}\right] \\ &\leq \Pr\left[\sum_{j=1}^{L_1} X_{i,j} \leq \Upsilon + 1\right] \\ &= \Pr\left[\frac{\sum_{j=1}^{L_1} X_{i,j}}{L_1} \leq \frac{\Upsilon + 1}{L_1}\right] \\ &\leq \Pr\left[\frac{\sum_{j=1}^{L_1} X_{i,j}}{L_1} \leq (1 - \varepsilon)\mu_i + \frac{1}{L_1}\right]. \end{aligned}$$

Moreover, by the definition of L_1 , we have

$$\begin{aligned} \frac{1}{L_1} &\leq \frac{(1 - \varepsilon)\mu_i}{\Upsilon} \\ &= \frac{(1 - \varepsilon)\mu_i}{(1 + \varepsilon)(1 + (2 + \frac{2}{3}\varepsilon) \ln(\frac{2}{\delta}) \frac{1}{\varepsilon^2})} \\ &< \frac{\varepsilon^2 \mu_i}{1 + \varepsilon}. \end{aligned}$$

Therefore,

$$\begin{aligned} & \Pr[\tilde{\mu}_i < (1 - \varepsilon)\mu_i] \\ &\leq \Pr\left[\frac{\sum_{j=1}^{L_1} X_{i,j}}{L_1} < (1 - \varepsilon)\mu_i + \frac{\varepsilon^2 \cdot \mu_i}{1 + \varepsilon}\right] \\ &\leq \Pr\left[\frac{\sum_{j=1}^{L_1} X_{i,j}}{L_1} < (1 - \frac{\varepsilon}{1 + \varepsilon}) \cdot \mu_i\right] \\ &= \Pr\left[\sum_{j=1}^{L_1} X_{i,j} - L_1 \mu_i < -\frac{\varepsilon}{1 + \varepsilon} \cdot L_1 \mu_i\right]. \end{aligned}$$

Meanwhile,

$$L_1 \geq \frac{\Upsilon}{(1 - \varepsilon)\mu_i} > \frac{2(1 + \varepsilon) \ln(\frac{2}{\delta})}{(1 - \varepsilon)\varepsilon^2 \mu_i} > \frac{2(1 + \varepsilon)^2 \ln(\frac{2}{\delta})}{\varepsilon^2 \mu_i}.$$

Note that $\frac{\sum_{j=1}^{L_1} X_{i,j}}{L_1}$ is an estimate of μ_i using the first L_1 random samples. Applying (24), we obtain

$$\begin{aligned} \Pr[\tilde{\mu}_i \leq (1 - \varepsilon)\mu_i] &\leq \exp\left(-\frac{\varepsilon^2 L_1 \mu_i}{2(1 + \varepsilon)^2}\right) \\ &< \exp\left(-\frac{\varepsilon^2 2(1 + \varepsilon)^2 \ln(\frac{2}{\delta}) \mu_i}{\varepsilon^2 \mu_i 2(1 + \varepsilon)^2}\right) \\ &= \frac{\delta}{2}. \end{aligned}$$

This completes the proof of (28).

Next, we prove (29), which is similar. Let $L_2 := \lceil \frac{\Upsilon}{(1+\varepsilon)\mu_i} \rceil$. Then, we have

$$\begin{aligned} \Pr[\tilde{\mu}_i > (1 + \varepsilon)\mu_i] &= \Pr[\Upsilon > (1 + \varepsilon)\mu_i \theta_i] \\ &= \Pr\left[\frac{\Upsilon}{(1 + \varepsilon)\mu_i} > \theta_i\right] \\ &= \Pr[L_2 \geq \theta_i] \\ &\leq \Pr\left[\sum_{j=1}^{L_2} X_{i,j} \geq \sum_{j=1}^{\theta_i} X_{i,j}\right] \\ &\leq \Pr\left[\sum_{j=1}^{L_2} X_{i,j} \geq \Upsilon\right] \\ &= \Pr\left[\frac{\sum_{j=1}^{L_2} X_{i,j}}{L_2} \geq \frac{\Upsilon}{L_2}\right]. \end{aligned}$$

By the definition of L_2 , we have

$$L_2 \leq \frac{\Upsilon}{(1 + \varepsilon)\mu_i},$$

which indicates that

$$\frac{\Upsilon}{L_2} \geq (1 + \varepsilon)\mu_i.$$

In addition, $L_2 > \frac{\gamma}{(1+\varepsilon)\mu_i} - 1 = \frac{1}{\mu_i} + (2 + \frac{2}{3}\varepsilon) \ln(\frac{2}{\delta}) \frac{1}{\varepsilon^2} - 1 > (2 + \frac{2}{3}\varepsilon) \ln \frac{2}{\delta} \cdot \frac{1}{\varepsilon^2} \cdot \frac{1}{\mu_i}$. By (23), we obtain

$$\begin{aligned} \Pr[\tilde{\mu}_i > (1 + \varepsilon)\mu_i] &\leq \Pr\left[\frac{\sum_{j=1}^{L_2} X_{i,j}}{L_2} \geq (1 + \varepsilon)\mu_i\right] \\ &= \Pr\left[\sum_{j=1}^{L_2} X_{i,j} - L_2\mu_i \geq \varepsilon \cdot L_2\mu_i\right] \\ &\leq \exp\left(-\frac{\varepsilon^2 L_2 \mu_i}{2 + \frac{2}{3}\varepsilon}\right) \\ &< \exp\left(-\frac{\varepsilon^2 (2 + \frac{2}{3}\varepsilon) \ln \frac{2}{\delta} \frac{1}{\varepsilon^2} \frac{1}{\mu_i} \mu_i}{2 + \frac{2}{3}\varepsilon}\right) \\ &= \frac{\delta}{2}. \end{aligned}$$

Combining (28) and (29) gives rise to (27). \square

With the values \tilde{A}_i° returned by Algorithm 1, according to Theorem 3, each value p_i° for $i = 1, 2, \dots, n_c$ can be estimated accurately by \tilde{A}_i° with a high probability, i.e., $\Pr[(1 - \varepsilon)p_i^\circ \leq \tilde{A}_i^\circ \leq (1 + \varepsilon)p_i^\circ] \geq 1 - \delta$. To ensure all the values p_i° are estimated accurately, by a union bound, we have

$$\Pr\left[\bigwedge_{i=1}^{n_c} (1 - \varepsilon)p_i^\circ \leq \tilde{A}_i^\circ \leq (1 + \varepsilon)p_i^\circ\right] \geq 1 - n_c\delta.$$

Thus, to ensure the estimation accuracy with a high probability of $1 - \delta$, we can simply scale δ by a factor of $1/n_c$ as an input to Algorithm 1.

5.4 Estimation of Optimal Budgeted Prices

Based on the proof of Corollary 1, under the condition that $b \geq \sigma(C)$, p_i^* can be represented as

$$p_i^* = 2(\lambda + c_i) = A_i + \frac{b}{n_c},$$

where A_i can be estimated using the RIS method. When θ RR sets are generated, an RR set R contributes to the estimation \tilde{A}_i by an additive factor of $\Delta(\tilde{A}_i, R)$. According to (18), $\Delta(\tilde{A}_i, R)$ can be computed by

$$\Delta(\tilde{A}_i, R) = \begin{cases} \frac{n}{\theta} \cdot (1 - \frac{n_r}{n_c}) \cdot 2^{1-n_r}, & \text{if } s_i \in R, \\ -\frac{n}{\theta} \cdot \frac{n_r}{n_c} \cdot 2^{1-n_r}, & \text{otherwise,} \end{cases} \quad (30)$$

where $n_r = |C \cap R|$. Note that the random variables $\Delta(\tilde{A}_i, R)$ may be negative. To tackle this issue, we shift the random variables to fall in the range of $[0, 1]$ so that the stopping rule algorithm can be applied.

Inspired by the proof of Corollary 1, $\Delta(\tilde{A}_i, R)$ can be made non-negative by adding a factor of $\frac{n}{\theta} \cdot \frac{\mathbf{1}_{\{C \cap R \neq \emptyset\}}}{n_c}$. Let

$$\Delta(\tilde{A}'_i, R) := \Delta(\tilde{A}_i, R) + \frac{n}{\theta} \cdot \frac{\mathbf{1}_{\{C \cap R \neq \emptyset\}}}{n_c}.$$

Table 2 Datasets.

Dataset	#nodes	#edges	Type	Avg. deg.
Facebook	4.0K	88.2K	Undirected	43.7
Google+	107.6K	13.7M	Directed	254.1
LiveJournal	4.8M	69.0M	Directed	28.5
Orkut	3.1M	117.2M	Undirected	76.3
Twitter	41.7M	1.5G	Directed	70.5

If we aggregate $\Delta(\tilde{A}'_i, R)$ over θ random RR sets, the actual value estimated is

$$A'_i := A_i + \frac{\sigma(C)}{n_c} = p_i^* + \frac{\sigma(C) - b}{n_c}.$$

Thus, we shall first estimate A'_i using the stopping rule algorithm and then compute p_i^* . Analogous to $X_{i,j}$ defined in (19), let $X'_{i,j}$ be a random variable defined as

$$X'_{i,j} = \begin{cases} 0 & \text{if } R_j \cap C = \emptyset, \\ \frac{1}{n_c} + 2^{1-n_r} \cdot (1 - \frac{n_r}{n_c}) & \text{if } s_i \in R_j \cap C, \\ \frac{1}{n_c} - 2^{1-n_r} \cdot \frac{n_r}{n_c} & \text{otherwise,} \end{cases} \quad (31)$$

where $n_c = |C|$ and $n_r = |R_j \cap C|$. It is easy to verify that $0 \leq X'_{i,j} \leq 1$, which meets the requirement of our stopping rule algorithm.

In addition to A_i , we also need to estimate $\frac{b - \sigma(C)}{n_c}$ in order to compute p_i^* as $p_i^* = A_i + \frac{b - \sigma(C)}{n_c}$. According to Corollary 1, to guarantee $p_i^* \geq 0$, the total price b of all the candidate nodes must be set no less than $\sigma(C)$, i.e., $b \geq \sigma(C)$. In general, we can also use the stopping rule algorithm and the RIS method to estimate $b - \sigma(C)$. By the union bound, we can set the failure probability $\delta' = \frac{\delta}{n_c + 1}$ in the stopping rule algorithm so that all the values A'_i and $b - \sigma(C)$ are estimated within a multiplicative factor of ε with probability at least $1 - \delta$, which gives rise to an (ε, δ) -approximation of all the prices p_i^* .

6 Experiments

This section experimentally evaluates the quality and scalability of our proposed algorithms. We implement our algorithms using C++. All experiments are run on a machine with Intel Xeon 2.4GHz CPU and 384GB memory.

6.1 Experimental Setup

Datasets. We evaluate our algorithms by several real datasets including Facebook, Google+, LiveJournal,

Table 3 Running time and number of RR sets generated.

Dataset	Time (s)			#RR sets		
	$n_c = 200$	$n_c = 500$	$n_c = 1000$	$n_c = 200$	$n_c = 500$	$n_c = 1000$
Facebook	10.97	16.95	27.52	2.27E+06	2.97E+06	3.42E+06
Google+	41.04	73.33	140.14	8.11E+06	1.33E+07	2.20E+07
LiveJournal	528.183	1295.84	1928.26	3.24E+07	7.09E+07	1.08E+08
Orkut	1439.66	3234.26	6027.80	9.83E+06	2.08E+07	3.70E+07
Twitter	4273.82	10869.80	26840.70	3.20E+06	7.49E+06	1.40E+07

Orkut and Twitter. The first four datasets are available at <http://snap.stanford.edu/data> and the Twitter dataset is obtained from <http://an.kaist.ac.kr/traces/WWW2010.html> [23]. Table 2 gives the details of these datasets.

Parameter Settings. We adopt the Independent Cascade diffusion model and set the propagation probability $p_{u,v}$ of each edge (u,v) to the reciprocal of v 's in-degree which is a commonly used setting by other studies [34, 35, 37, 38]. We set the number of candidate nodes $n_c = 200, 500$ or 1000 , the failure probability $\delta = \frac{1}{n}$ (n is the number of nodes in the OSN) and the error threshold $\varepsilon = 0.1$ by default. We assume that the candidate node set C includes the top- n_c nodes with the highest out-degrees. These nodes are offered by the OSN provider to advertisers for seed selection.

Algorithms. We compare the price profile calculated by our pricing algorithm (Algorithm 1), referred to as OptPrice, with the following baselines:

- **Uniform:** The prices of all the candidate nodes are set the same.
- **Degree:** The price of each candidate node is set proportional to its out-degree.
- **SingletonInf:** The price of each candidate node is set proportional to the influence spread it can produce when selected as the only seed. We estimate the influence spread using the RIS method and the stopping rule algorithm.
- **IMRank [8]:** A ranking of candidate nodes is generated in decreasing order of their marginal gains in influence spread. We generate the ranking by applying the greedy hill-climbing algorithm for influence maximization [21]. The price of each candidate node is set proportional to its marginal gain.

Comparison of Time Complexity. In the stopping rule algorithm (Algorithm 1), by the analysis in Theorem 3, the number of samples θ_i generated for estimating μ_i satisfies

$$\Pr \left[\frac{\Upsilon}{(1-\varepsilon) \cdot \mu_i} < \theta_i \right] = \Pr[\tilde{\mu}_i < (1-\varepsilon)\mu_i] < \frac{\delta}{2}.$$

To estimate all the node prices in the candidate set, our stopping rule algorithm finishes with $O\left(\frac{\Upsilon}{\min_{s_i \in C} \mu_i}\right)$ samples generated with probability at least $1 - \frac{\delta}{2}$. In addition, the expected number of RR sets generated is $\frac{\Upsilon}{\min_{s_i \in C} \mu_i}$. Let EPT be the expected time complexity to generate an RR set. According to [37], let v^* be a random node chosen from V with probability proportional to its in-degree and we have $\text{EPT} = \mathbb{E}[\sigma(\{v^*\})] \cdot \frac{m}{n}$ where the expectation is over the randomness of v^* , $n = |V|$ is the number of nodes and $m = |E|$ is the number of edges in the network. The expected time complexity of our pricing algorithm is then $O\left(\frac{\Upsilon}{\min_{s_i \in C} \mu_i} \cdot \text{EPT}\right)$. Let α_i be the estimation variable for the singleton influence spread for a seed $s_i \in C$, i.e., $\alpha_i = \frac{\sigma(s_i)}{n}$. Similarly, the expected time complexity of the SingletonInf pricing is $O\left(\frac{\Upsilon}{\min_{s_i \in C} \alpha_i} \cdot \text{EPT}\right)$. When using the RR sets generated by our stopping rule algorithm, the IMRank algorithm has the same time complexity as our pricing algorithm. Since the baselines of Degree and Uniform pricing do not incur computational cost for sampling, their time complexities are $O(|C|)$.

6.2 Experimental Results

6.2.1 Efficiency of Our Algorithm

Table 3 shows the running time of our pricing algorithm for computing the optimal price profile $\langle p_1^o, p_2^o, \dots, p_{n_c}^o \rangle$ as derived in Section 4.2 and the number of RR sets generated for various datasets. As can be seen, our OptPrice algorithm can compute the price profile within hours even for large-scale datasets. This demonstrates the scalability of our OptPrice algorithm.

6.2.2 Evaluation of Divergence Function

A straightforward evaluation is to compare the values of the divergence function (2) produced by the pricing profiles of different algorithms. According to Lemma 2, the divergence function can be divided into three parts: $\frac{1}{2n_c} \sum_{S \subseteq C} \sigma(S)^2$, $\frac{b^2}{4} = \frac{(\sum_{i=1}^{n_c} p_i)^2}{4}$ and $\sum_{i=1}^{n_c} g(p_i)$.

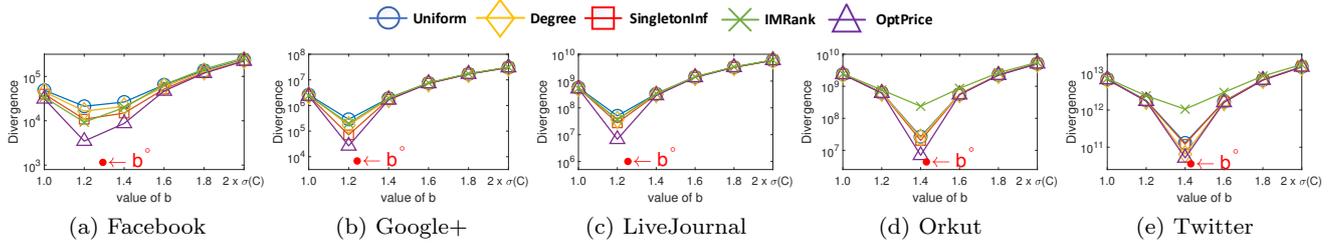


Fig. 1 Value of divergence function under different price budgets ($n_c = 200$).

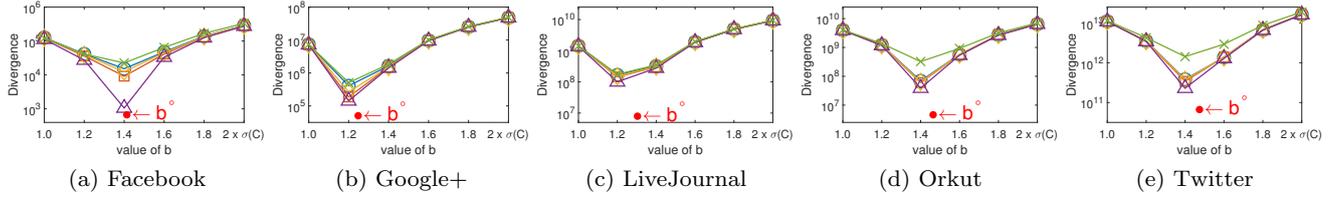


Fig. 2 Value of divergence function under different price budgets ($n_c = 500$).

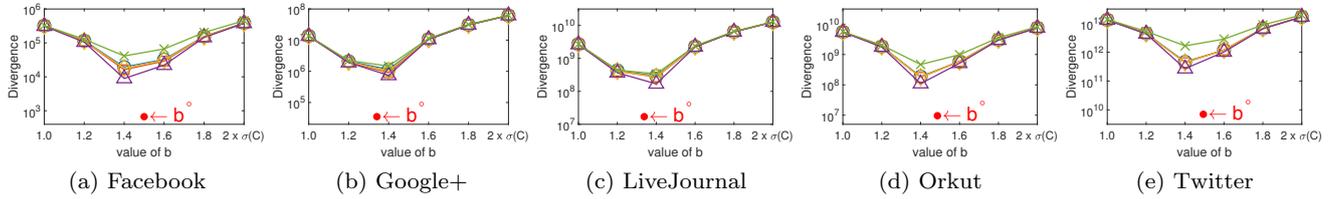


Fig. 3 Value of divergence function under different price budgets ($n_c = 1000$).

Table 4 Comparison between ω and $\beta \cdot \sigma(C)^2$.

Dataset	$n_c = 200$		$n_c = 500$		$n_c = 1000$	
	ω	$\beta \cdot \sigma(C)^2$	ω	$\beta \cdot \sigma(C)^2$	ω	$\beta \cdot \sigma(C)^2$
Facebook	6.82E+05	3.29E+03	1.43E+06	5.86E+03	2.99E+06	1.09E+04
Google+	7.34E+07	3.93E+05	1.50E+08	7.32E+05	2.30E+08	1.06E+06
LiveJournal	1.51E+10	7.94E+07	2.92E+10	1.42E+08	4.62E+10	2.14E+08
Orkut	2.83E+10	1.14E+08	4.24E+10	1.63E+08	5.47E+10	2.05E+08
Twitter	8.10E+13	3.24E+11	1.13E+14	4.27E+11	1.29E+14	4.75E+11

For each $g(p_i)$, the value of $\frac{1}{2^{n_c-1}} \sum_{S \subseteq C \setminus \{s_i\}} \sigma(S \cup \{s_i\})$ can be estimated using the RIS method and the stopping rule algorithm. Then, $\sum_{i=1}^{n_c} g(p_i)$ can be computed based on (3). The challenge lies in evaluating $\frac{1}{2^{n_c}} \sum_{S \subseteq C} \sigma(S)^2$. This part is non-linear with respect to the influence spread. There are an exponential number of seed sets to measure in order to obtain the sum. To make the evaluation tractable, we use the sample average to estimate the value of the sum. Note that this sum is an additive term in the divergence function that is independent from the price profile, which indicates that its estimation accuracy will not affect the relative performance of different algorithms. Lemma 4 shows the theoretical guarantees of the estimation accuracy when a given number of T seed sets are measured.

Let S_1, S_2, \dots, S_T be a sequence of randomly generated subsets of the candidate node set C . Let $\omega =$

$\frac{1}{T} \sum_{j=1}^T \sigma(S_j)^2$. Then, for each $1 \leq j \leq T$, we have $\mathbb{E}[\sigma(S_j)^2] = \mathbb{E}[\omega] = \frac{1}{2^{n_c}} \sum_{S \subseteq C} \sigma(S)^2$. To make $\sigma(S_j)^2$ fall in the range of $[0, 1]$, we normalize the value of $\sigma(S_j)^2$ by $\sigma(C)^2$ so that $0 \leq \frac{\sigma(S_j)^2}{\sigma(C)^2} \leq 1$.

Suppose that we set $T = \frac{(1+\varepsilon/3)^2 \ln(2/\delta)}{2\varepsilon\beta}$. Then, according to Lemma 4, the estimation $\frac{1}{T} \sum_{j=1}^T \sigma(S_j)^2$ is in the range of $[(1-\varepsilon)\mathbb{E}[\omega] - \beta \cdot \sigma(C)^2, (1+\varepsilon)\mathbb{E}[\omega] + \beta \cdot \sigma(C)^2]$ with probability at least $1-\delta$. In the experiments, we set $\varepsilon = 0.01$, $\beta = 0.002$, and $\delta = 0.01$. Then, $T = 133,342$. So, we randomly generate 133,342 subsets of C to estimate the value of $\frac{1}{2^{n_c}} \sum_{S \subseteq C} \sigma(S)^2$. As can be seen from Table 4, under this setting, the additive estimation error of $\beta \cdot \sigma(C)^2$ is negligible compared to the estimated value ω .

Figures 1, 2 and 3 show the value of divergence function (in log scale) produced by different pricing algo-

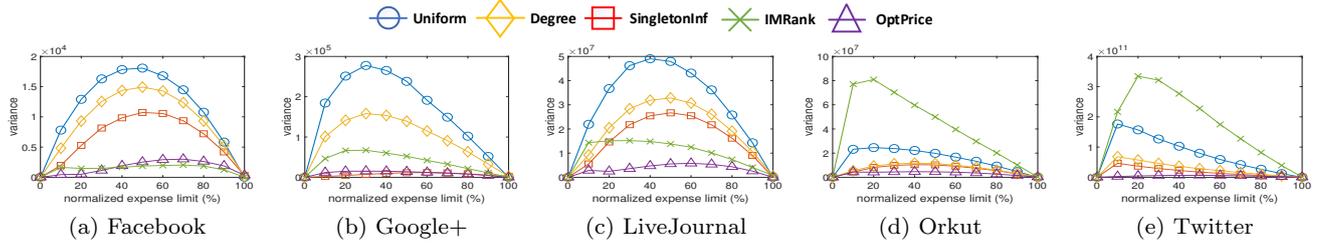


Fig. 4 Variance of 10,000 seed sets for different expense limits ($n_c = 200$).

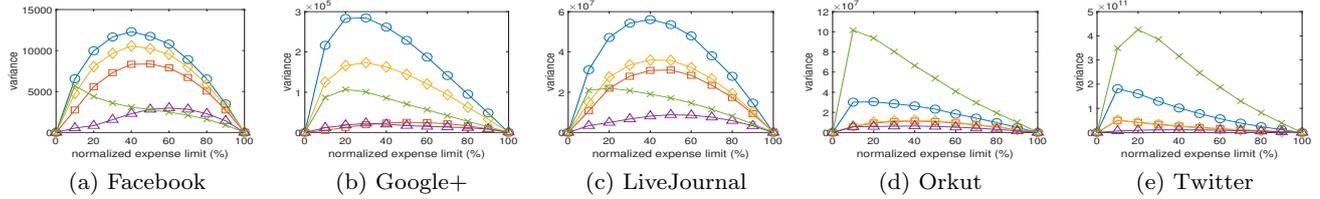


Fig. 5 Variance of 10,000 seed sets for different expense limits ($n_c = 500$).

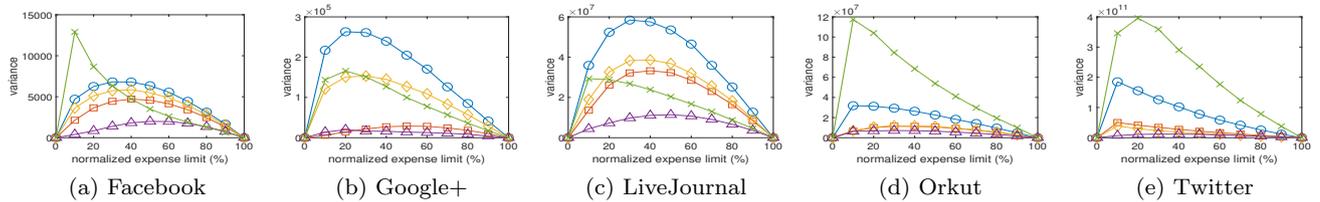


Fig. 6 Variance of 10,000 seed sets for different expense limits ($n_c = 1000$).

gorithms under different price budgets and we also highlight the value delivered by the optimal price profile under $b^\circ = \sum_{i=1}^{n_c} p_i^\circ$ derived in Section 4.2. We can see that the value of b can significantly affect the divergence value and our optimal price profile under b° results in significantly lower divergence values. For example, for Twitter, the divergence value produced by the optimal price profile under b° is around 95% less than that under $b = 1.4 \cdot \sigma(C)$ when $n_c = 1000$. This indicates that the value of b is critical in optimizing the divergence function. In addition, it can also be seen that our OptPrice algorithm can effectively reduce the divergence value and dramatically outperforms other baselines when the price budget b is close to b° . This is because our divergence function is a quadratic function regarding b , i.e., the divergence value is dominated by b^2 or $\sum_{S \subseteq C} \sigma(S)^2$ when b is far from b° , which indicates that optimizing the price profile can do little help to bring the divergence value down in these cases. In summary, our OptPrice algorithm can deliver much lower divergence values and better reflect the influence spread of any seed set by optimizing both the price budget and the price profile.

Discussion. Interestingly, we found that the empirical setting of b in our preliminary work [39] happens to be rather close to the value of b° (but not the same) derived in this paper and thus the divergence value pro-

duced by the earlier set budget is almost the same as the divergence value achieved at b° in this work. The difference is insignificant and hard to distinguish when plotted, so we did not mark the divergence value achieved at the earlier set budget. We think that this observation provides strong justification for the empirical setting of b in our preliminary work [39].

6.2.3 Stability of Influence Spread

Recall that the objective of our pricing profile is to minimize the divergence between the influence profiles and the seed prices for all possible seed sets. When such a divergence is minimized, the influence spread of any seed set would closely match its price. This implies that two seed sets with the same price would also be close in their influence spreads. This property is valuable to the advertisers because the advertisers, holding a given fund to purchase seeds according to their prices, are expecting to achieve a more predictable influence spread from the selected seeds. To verify this property, we design the experiments to construct a collection of seed sets subject to an expense limit and compare their influence spreads. We use the optimal price profile $\langle p_1^\circ, p_2^\circ, \dots, p_{n_c}^\circ \rangle$ derived in Section 4.2. The expense limit is expressed as a percentage of $b^\circ = \sum_{i=1}^{n_c} p_i^\circ$ (the total price of all the candidate nodes). Given an ex-

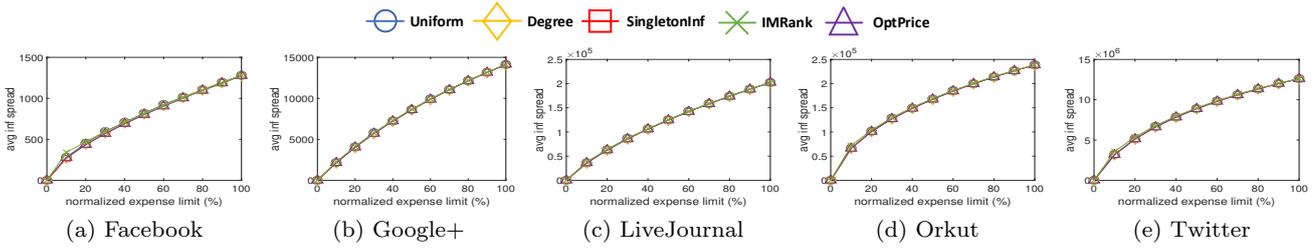


Fig. 7 Average influence spread of 10,000 seed sets for different expense limits ($n_c = 200$).

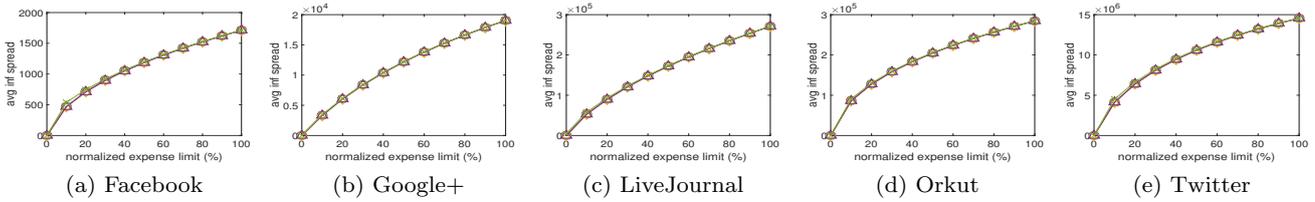


Fig. 8 Average influence spread of 10,000 seed sets for different expense limits ($n_c = 500$).

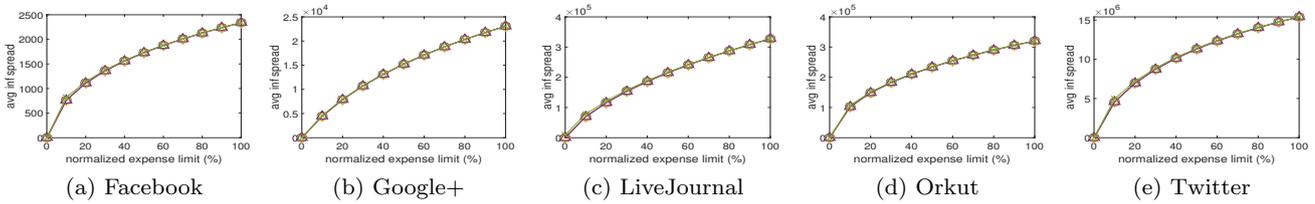


Fig. 9 Average influence spread of 10,000 seed sets for different expense limits ($n_c = 1000$).

pense limit, we randomly choose seeds among the candidate nodes to fill up the limit. We construct 10,000 random seed sets and estimate the influence spreads achieved by these seed sets using the RIS method with 100,000 RR sets. We test different expense limits from 10% to 90%.

Figures 4, 5 and 6 show the variance of the influence spreads of the 10,000 seed sets under different pricing algorithms when there are $n_c = 200$, 500 and 1000 candidate nodes respectively. When the expense limit is very low, only nodes of relatively cheap prices can be selected within the limit. When the expense limit is very high, the majority of the candidate nodes need to be selected to fill up the budget. In both cases, the number of different combinations of the seeds to fill up the expense limit is relatively small, which gives rise to a relatively stable influence spread. As a result, the variance of influence spread is low when the expense limit is very low or very high and it shows a concave shape. It can be seen that different pricing algorithms result in quite different variances. In general, our OptPrice algorithm has significantly lower variance of influence spread than the baselines. This shows that our OptPrice algorithm can capture the influence potentials of the candidate nodes more accurately and give the advertisers a more stable and predictable return (influence spread) for their purchasing (seeding) activities.

Figures 7, 8 and 9 show the average influence spread of the 10,000 seed sets. As can be seen, the average influence spread increases with the given expense limit of seed selection for all the datasets. The seed sets chosen under our proposed pricing algorithm achieve comparable average influence spreads to those under the baselines. This shows that our pricing profile is able to maintain the same expected influence spread as other baselines under a given expense limit.

7 Conclusion

In this work, we build a bridge between OSN providers and advertisers by proposing a pricing mechanism to facilitate the initiator selection of marketing campaigns without the knowledge of OSN structures. In particular, we study the problem of minimizing the pricing divergence from the influence spread and derive an optimal price profile. A scalable estimation algorithm is devised to yield an (ϵ, δ) -approximation of the optimal prices. Through extensive experiments, we demonstrate the performance advantages of our approach over other baselines.

Acknowledgements This work is partially supported by HKUST(GZ) under a Startup Grant, and by Singapore Ministry of Education Academic Research Fund Tier 1 under Grants 2018-T1-002-063 and 2019-T1-002-042.

References

1. Aslay C, Lu W, Bonchi F, Goyal A, Lakshmanan LV (2015) Viral marketing meets social advertising: Ad allocation with minimum regret. *Proc VLDB Endowment* 8(7):814–825
2. Aslay C, Bonchi F, Lakshmanan LV, Lu W (2017) Revenue maximization in incentivized social advertising. *Proc VLDB Endowment* 10(11):1238–1249
3. Borgs C, Brautbar M, Chayes J, Lucier B (2014) Maximizing social influence in nearly optimal time. In: *Proc. SODA*, pp 946–957
4. Chalermsook P, Das Sarma A, Lall A, Nanongkai D (2015) Social network monetization via sponsored viral marketing. *ACM SIGMETRICS Performance Evaluation Review* 43(1):259–270
5. Chen W, Wang Y, Yang S (2009) Efficient influence maximization in social networks. In: *Proc. ACM KDD*, pp 199–208
6. Chen W, Wang C, Wang Y (2010) Scalable influence maximization for prevalent viral marketing in large-scale social networks. In: *Proc. ACM KDD*, pp 1029–1038
7. Chen W, Yuan Y, Zhang L (2010) Scalable influence maximization in social networks under the linear threshold model. In: *Proc. IEEE ICDM*, pp 88–97
8. Cheng S, Shen H, Huang J, Chen W, Cheng X (2014) Imrank: influence maximization via finding self-consistent ranking. In: *Proc. ACM SIGIR*, pp 475–484
9. Chung F, Lu L (2006) Concentration inequalities and martingale inequalities: a survey. *Internet Mathematics* 3(1):79–127
10. Cohen E, Delling D, Pajor T, Werneck RF (2014) Sketch-based influence maximization and computation: Scaling up with guarantees. In: *Proc. ACM CIKM*, pp 629–638
11. Dagum P, Karp R, Luby M, Ross S (2000) An optimal algorithm for monte carlo estimation. *SIAM Journal on Computing* 29(5):1484–1496
12. Domingos P, Richardson M (2001) Mining the network value of customers. In: *Proc. ACM KDD*, pp 57–66
13. eMarketer (2015) Social network ad spending worldwide, by region, 2013–2017. <https://www.emarketer.com/Chart/Social-Network-Ad-Spending-Worldwide-by-Region-2013-2017/168356>
14. eMarketer (2017) Influencer marketing is about data, not celebrity deals. <https://www.emarketer.com/Article/Influencer-Marketing-About-Data-Not-Celebrity-Deals/1016683>
15. eMarketer (2019) How much are brands paying influencers? <https://www.emarketer.com/content/how-much-are-brands-paying-influencers>
16. eMarketer (2019) Is everyone on instagram an influencer? <https://www.emarketer.com/content/is-everyone-on-instagram-an-influencer>
17. Han K, Huang K, Xiao X, Tang J, Sun A, Tang X (2018) Efficient algorithms for adaptive influence maximization. *Proc VLDB Endowment* 11(9):1029–1040
18. Huang K, Tang J, Han K, Xiao X, Chen W, Sun A, Tang X, Lim A (2020) Efficient approximation algorithms for adaptive influence maximization. *The VLDB Journal* 29(6):1385–1406
19. Huang K, Tang J, Xiao X, Sun A, Lim A (2020) Efficient approximation algorithms for adaptive target profit maximization. In: *Proc. IEEE ICDE*, pp 649–660
20. Hub IM (2020) The remarkable rise of influencer marketing. <https://influencermarketinghub.com/the-rise-of-influencer-marketing/>
21. Kempe D, Kleinberg J, Tardos É (2003) Maximizing the spread of influence through a social network. In: *Proc. ACM KDD*, pp 137–146
22. Khan A, Zehnder B, Kossmann D (2016) Revenue maximization by viral marketing: A social network host’s perspective. In: *Proc. IEEE ICDE*, pp 37–48
23. Kwak H, Lee C, Park H, Moon S (2010) What is Twitter, a social network or a news media? In: *Proc. WWW*, pp 591–600
24. Leskovec J, Krause A, Guestrin C, Faloutsos C, VanBriesen J, Glance N (2007) Cost-effective outbreak detection in networks. In: *Proc. ACM KDD*, pp 420–429
25. Lu W, Lakshmanan LV (2012) Profit maximization over social networks. In: *Proc. IEEE ICDM*, pp 479–488
26. Nguyen HT, Thai MT, Dinh TN (2016) Stop-and-stare: Optimal sampling algorithms for viral marketing in billion-scale networks. In: *Proc. ACM SIGMOD*, pp 695–710
27. Nguyen HT, Nguyen TP, Vu TN, Dinh TN (2017) Outward influence and cascade size estimation in billion-scale networks. In: *Proc. ACM SIGMETRICS*, pp 63–63
28. Reuters (2016) Social media ad spending is expected to pass newspapers by 2020. <http://fortune.com/2016/12/05/social-media-ad-spending-newspapers-zenith-2020/>
29. Richardson M, Domingos P (2002) Mining knowledge-sharing sites for viral marketing. In:

- Proc. ACM KDD, pp 61–70
30. Tang J, Tang X, Yuan J (2016) Profit maximization for viral marketing in online social networks. In: Proc. IEEE ICNP, pp 1–10
 31. Tang J, Tang X, Yuan J (2017) Influence maximization meets efficiency and effectiveness: A hop-based approach. In: Proc. IEEE/ACM ASONAM, pp 64–71
 32. Tang J, Tang X, Xiao X, Yuan J (2018) Online processing algorithms for influence maximization. In: Proc. ACM SIGMOD, pp 991–1005
 33. Tang J, Tang X, Yuan J (2018) An efficient and effective hop-based approach for influence maximization in social networks. *Social Network Analysis and Mining* 8(1):10
 34. Tang J, Tang X, Yuan J (2018) Profit maximization for viral marketing in online social networks: Algorithms and analysis. *IEEE Transactions on Knowledge and Data Engineering* 30(6):1095–1108
 35. Tang J, Tang X, Yuan J (2018) Towards profit maximization for online social network providers. In: Proc. IEEE INFOCOM, pp 1178–1186
 36. Tang J, Huang K, Xiao X, Lakshmanan LV, Tang X, Sun A, Lim A (2019) Efficient approximation algorithms for adaptive seed minimization. In: Proc. ACM SIGMOD, pp 1096–1113
 37. Tang Y, Xiao X, Shi Y (2014) Influence maximization: Near-optimal time complexity meets practical efficiency. In: Proc. ACM SIGMOD, pp 75–86
 38. Tang Y, Shi Y, Xiao X (2015) Influence maximization in near-linear time: A martingale approach. In: Proc. ACM SIGMOD, pp 1539–1554
 39. Zhu Y, Tang J, Tang X (2020) Pricing influential nodes in online social networks. *Proc VLDB Endowment* 13(10):1614–1627