

Gradient Method for Continuous Influence Maximization with Budget-Saving Considerations

Wei Chen
Microsoft Research
weic@microsoft.com

Weizhong Zhang
IIIS, Tsinghua University
zwz15@mails.tsinghua.edu.cn

Haoyu Zhao
IIIS, Tsinghua University
zhaohy16@mails.tsinghua.edu.cn

Abstract

Continuous influence maximization (CIM) generalizes the original influence maximization by incorporating general marketing strategies: a marketing strategy mix is a vector $\mathbf{x} = (x_1, \dots, x_d)$ such that for each node v in a social network, v could be activated as a seed of diffusion with probability $h_v(\mathbf{x})$, where h_v is a strategy activation function satisfying DR-submodularity. CIM is the task of selecting a strategy mix \mathbf{x} with constraint $\sum_i x_i \leq k$ where k is a budget constraint, such that the total number of activated nodes after the diffusion process, called influence spread and denoted as $g(\mathbf{x})$, is maximized. In this paper, we extend CIM to consider budget saving, that is, each strategy mix \mathbf{x} has a cost $c(\mathbf{x})$ where c is a convex cost function, and we want to maximize the balanced sum $g(\mathbf{x}) + \lambda(k - c(\mathbf{x}))$ where λ is a balance parameter, subject to the constraint of $c(\mathbf{x}) \leq k$. We denote this problem as CIM-BS. The objective function of CIM-BS is neither monotone, nor DR-submodular or concave, and thus neither the greedy algorithm nor the standard result on gradient method could be directly applied. Our key innovation is the combination of the gradient method with reverse influence sampling to design algorithms that solve CIM-BS: For the general case, we give an algorithm that achieves $(\frac{1}{2} - \varepsilon)$ -approximation, and for the case of independent strategy activations, we present an algorithm that achieves $(1 - \frac{1}{e} - \varepsilon)$ approximation.

1 Introduction

Influence maximization is the task of selecting a small number of seed nodes in a social network such that the influence spread from the seeds when following an influence diffusion model is maximized. It models the viral marketing scenario and has been extensively studied (cf. [12, 5, 13]). Continuous influence maximization (CIM) generalizes the original influence maximization by incorporating general marketing strategies: a marketing strategy mix is a vector $\mathbf{x} = (x_1, \dots, x_d)$ such that for each node v in a social network, v could be activated as a seed of diffusion with probability $h_v(\mathbf{x})$, where h_v is a strategy activation function satisfying monotonicity and DR-submodularity. CIM is the task of selecting a strategy mix \mathbf{x} with constraint $\sum_i x_i \leq k$ where k is a budget constraint, such that the total number of activated nodes after the diffusion process, called influence spread and denoted as $g(\mathbf{x})$, is maximized. CIM is proposed in [12], and recently followed up by a few studies [27, 7].

In this paper, we extend CIM to consider budget saving: each strategy mix \mathbf{x} has a cost $c(\mathbf{x})$ where c is a convex cost function, and we want to maximize the balanced sum $g(\mathbf{x}) + \lambda(k - c(\mathbf{x}))$ where λ is a balance parameter, subject to the constraint of $c(\mathbf{x}) \leq k$. We denote this problem as CIM-BS. The objective reflects the realistic consideration of balancing between increasing influence spread and saving marketing budget. In general we have $g(\mathbf{x})$ monotone (increasing) and DR-submodular (diminishing return property formally defined in Section 2), but $\lambda(k - c(\mathbf{x}))$ is concave and likely to be monotonically decreasing, and thus the objective function $g(\mathbf{x}) + \lambda(k - c(\mathbf{x}))$ is neither monotone, nor DR-submodular or concave, and thus neither

the greedy algorithm nor the standard result on gradient method could be directly applied for a theoretical guarantee.

In this paper, we apply the gradient method [17, 20] to solve CIM-BS with theoretical approximation guarantees. This is the only case we know of that the gradient method is applied to influence maximization with a theoretical guarantee while the greedy method cannot. The gradient method may be applied to the original objective function $g(\mathbf{x}) + \lambda(k - c(\mathbf{x}))$, but $g(\mathbf{x})$ is a complicated combinatorial function and its exact gradient is infeasible to compute. We could use stochastic gradient instead, but it results in large variance and very slow convergence. Instead, we integrate the gradient method with the reverse influence sampling (RIS) approach [2, 24, 23], which is the main technical innovation in our paper. RIS is proposed for improving the efficiency of the influence maximization task, but when integrating with the gradient method, it brings two additional benefits: (a) it allows the efficient computation of the exact gradient of an estimator function of $g(\mathbf{x})$, which avoids slow convergence caused by the large variance in the stochastic gradient, and (b) for a class of independent strategy activation functions h_v where each strategy dimension independently attempts to activate node v , the new objective function is in the form of coverage functions [11], which allows a tight concave upper bound function and leads to a better approximation ratio.

For the general case, we apply the proximal gradient method originally designed for concave functions to work with RIS and achieve an approximation of $(\frac{1}{2} - \epsilon)$ (Theorem 3). This requires an adaptation of the proximal gradient method for the functions of the form $f_1(\mathbf{x}) + f_2(\mathbf{x})$ where f_1 is non-negative, monotone and DR-submodular and f_2 is non-negative and concave, and the result of this adaptation (Theorem 2) may be of independent interest. For the independent strategy activation case, we apply the projected subgradient method on a tight concave upper bound of the objective function and achieve a $(1 - \frac{1}{e} - \epsilon)$ approximation (Theorem 4). We test our algorithms on a real-world dataset and validate its effectiveness comparing with other algorithms.

In summary, our contributions include: (a) we propose the study of CIM-BS problem to balance influence spread with budget saving; and (b) we integrate the gradient method with reverse influence sampling and provide two algorithms with theoretical approximation guarantees, on the objective function that is neither monotone, nor DR-submodular or concave. Our study is one of the first studies that introduce the gradient method to influence maximization, and hopefully it will enrich the scope of the influence maximization research.

For clarity, the detailed proofs and full experiment results are moved to the appendix.

Related works. Influence maximization was first proposed by Kempe et al. [12] as a discrete optimization problem, and has been extensively studied since (cf. [5, 13]). CIM is also proposed in [12]. Yang et al. [27] propose heuristic algorithms to solve CIM more efficiently, while Wu et al. [7] consider discrete version of CIM and apply RIS to solve it efficiently. Profit maximization in [14, 21] introduces linear cost with no budget constraint to influence maximization. Our CIM-BS problem is new and more general than both CIM and profit maximization studied before. The RIS approach is originally proposed in [2], and is further improved in [24, 23, 18].

Recently, a number of studies have applied gradient methods to DR-submodular maximization: Bian et al. [1] apply the Frank-Wolfe algorithm to achieve $1 - 1/e$ approximation for down-closed sets; Hassani et al. [9] apply stochastic projected gradient descent to achieve $1/2$ approximation; Karimi et al. [10] achieve $1 - 1/e$ approximation for coverage functions, which we adopt for the independent strategy activation case; Mokhtari et al. [16] apply more complicated conditional gradient method to achieve $1 - 1/e$ approximation. Our study is not a simple adoption of such methods to CIM-BS, because our objective function is not DR-submodular, and gradient computation cannot be treated as an oracle — we have to provide exact gradient computation and an end-to-end integration with the RIS approach.

2 Preliminary and Model

In this paper, we focus on the triggering model for influence maximization problem. We use a directed graph $G = (V, E)$ to represent a social network, where V is the set of nodes representing individuals, and E is the set of directed edges with edge (u, v) representing that u could directly influence v . Let $n = |V|$ and

$m = |E|$. In the diffusion process, each node is either active or inactive, and a node will stay active if it is activated. In the triggering model, every node $v \in V$ has a distribution D_v on the subsets of v 's in-neighbors $N^-(v) = \{u | (u, v) \in E\}$. Before the diffusion starts, each node $v \in V$ samples a triggering set $T_v \subseteq N^-(v)$ from the distribution D_v , denoted $T_v \sim D_v$. At time $t = 0$, the nodes in a pre-determined *seed set* S are activated. For any time $t = 1, 2, \dots$, the node v is activated if at least one of nodes in its triggering set T_v is activated at time $t - 1$. The whole propagation stops when no new node is activated in a step. An important quantity is the *influence spread* of the seed set S , denoted as $\sigma(S)$, which is defined as the expected number of the final activated nodes with seed set S . The classical *influence maximization* problem is to maximize $\sigma(S)$ such that $|S| \leq k$ for some given budget k .

A generalization of the classical influence maximization problem is the *continuous influence maximization (CIM)* problem with general marketing strategies [12]. A mix of marketing strategies is represented by a d -dimensional vector $\mathbf{x} = (x_1, x_2, \dots, x_d) \in \mathbb{R}_+^d$, where \mathbb{R}_+ is the set of non-negative real numbers. In the general case, we consider strategy mix \mathbf{x} in a general convex set $\mathcal{D} \subseteq \mathbb{R}_+^d$, but most commonly we consider $\mathcal{D} = \mathbb{R}_+^d$ or \mathcal{D} has an upper bound in each dimension, e.g. $\mathcal{D} = [0, 1]^d$. Given the strategy \mathbf{x} , each node $v \in V$ is independently activated as a seed with probability $h_v(\mathbf{x})$, where $h_v : \mathbb{R}_+^d \rightarrow [0, 1]$ is referred to as a *strategy activation function*. Once a set of seeds S is activated by a marketing strategy mix \mathbf{x} , the influence propagates from seeds in S following the triggering model. Then we define the influence spread of strategy mix \mathbf{x} , $g(\mathbf{x})$, as the expected number of nodes activated by \mathbf{x} , and formally,

$$\begin{aligned} g(\mathbf{x}) &= \mathbb{E}_S[\sigma(S)] \\ &= \sum_{S \subseteq V} \sigma(S) \cdot \prod_{u \in S} h_u(\mathbf{x}) \cdot \prod_{v \notin S} (1 - h_v(\mathbf{x})). \end{aligned} \quad (1)$$

The above formula means that we enumerate through all possible seed sets S , and due to independent seed activation by \mathbf{x} the probability of S being the seed set is $\prod_{u \in S} h_u(\mathbf{x}) \cdot \prod_{v \notin S} (1 - h_v(\mathbf{x}))$ and its influence spread is $\sigma(S)$.

In many situations, each strategy dimension in \mathbf{x} activates each node independently. That is, for each node v and each strategy $j \in [d]$, there is a function $q_{v,j}$ such that strategy j with amount x_j activates node v with probability $q_{v,j}(x_j)$. Then we have $h_v(\mathbf{x}) = 1 - \prod_{j \in [d]} (1 - q_{v,j}(x_j))$. We call this case *independent strategy activation*. Independent strategy activation models many scenarios such as personalized marketing and event marketing [7].

In this paper, we focus on an extension of the continuous influence maximization problem — *continuous influence maximization with budget saving (CIM-BS)*. We have a total budget k , and for every strategy mix \mathbf{x} , there is a cost $c(\mathbf{x})$. We do not want the cost to exceed the budget, and we want to maximize the budget balanced influence spread: a combination of the expected influence spread and the remaining budget. More formally, we have the following definition.

Definition 1 (Continuous Influence Maximization with Budget Saving). *The continuous influence maximization with budget saving (CIM-BS) is the problem of given (a) a social network $G = (V, E)$ and the triggering model $\{D_v\}_{v \in V}$ on G , (b) strategy activation functions $\{h_v\}_{v \in V}$, (c) cost function c , total budget k , and a balance parameter $\lambda \geq 0$, finding a strategy mix $\mathbf{x}^* \in \mathbb{R}_+^d$ to maximize its balanced sum of influence spread and budget savings, i.e., find \mathbf{x}^* such that $\mathbf{x}^* \in \arg\max_{\mathbf{x} \in \mathcal{D}, c(\mathbf{x}) \leq k} (g(\mathbf{x}) + \lambda(k - c(\mathbf{x})))$.*

Note that when $\lambda = 0$ and $c(\mathbf{x}) = \sum_{i \in [d]} x_i$, the CIM-BS problem falls back to the CIM problem defined in [12], and also appears in [27, 7]. When $c(\mathbf{x})$ is a linear function and the budget constraint $c(\mathbf{x}) \leq k$ is dropped, the problem resembles the profit maximization studied in [14, 21]. However, the general version of the problem as defined here with $\lambda > 0$, a general cost function $c(\mathbf{x})$ and constraint $c(\mathbf{x}) \leq k$ together is new. Henceforth, let $s(\mathbf{x}) = \lambda(k - c(\mathbf{x}))$. Intuitively, $s(\mathbf{x})$ is the budget-saving part of the objective. The overall objective of $g(\mathbf{x}) + s(\mathbf{x})$ is trying to find the balance between maximizing influence and saving budget. We call $g(\mathbf{x}) + s(\mathbf{x})$ the *budget-balanced influence spread*. For convenience, we denote $\mathcal{P} = \{\mathbf{x} \mid \mathbf{x} \in \mathcal{D}, c(\mathbf{x}) \leq k\}$.

We say that a vector function $f : \mathcal{D} \rightarrow \mathbb{R}$ is DR-submodular if for any $\mathbf{x}, \mathbf{y} \in \mathcal{D}$ with $\mathbf{x} \leq \mathbf{y}$ (coordinate-wise), for any unit vector \mathbf{e}_i with the i -th dimension 1 and all other dimensions 0, and for any $\delta > 0$, we have

$f(\mathbf{x} + \delta \cdot \mathbf{e}_i) - b(\mathbf{x}) \geq f(\mathbf{y} + \delta \cdot \mathbf{e}_i) - f(\mathbf{y})$. DR-submodularity characterizes the diminishing marginal return on function f as vector \mathbf{x} increases, hence the name. We also say that f is monotone if for any $\mathbf{x}, \mathbf{y} \in \mathcal{D}$ with $\mathbf{x} \leq \mathbf{y}$, $f(\mathbf{x}) \leq f(\mathbf{y})$; f is convex if for any $\mathbf{x}, \mathbf{y} \in \mathcal{D}$, any $\lambda \in [0, 1]$, $f(\lambda \mathbf{x} + (1 - \lambda)\mathbf{y}) \leq \lambda f(\mathbf{x}) + (1 - \lambda)f(\mathbf{y})$; f is L -Lipschitz if for any $\mathbf{x}, \mathbf{y} \in \mathcal{D}$, $|f(\mathbf{x}) - f(\mathbf{y})| \leq L \cdot \|\mathbf{x} - \mathbf{y}\|_2$, where $\|\cdot\|_2$ is the vector 2-norm; f is β -smooth if it has gradients everywhere and for any $\mathbf{x}, \mathbf{y} \in \mathcal{D}$, $\|\nabla f(\mathbf{x}) - \nabla f(\mathbf{y})\|_2 \leq \beta \|\mathbf{x} - \mathbf{y}\|_2$. Note that when the gradients exist, the L -Lipschitz condition is equivalent to $\|\nabla f(\mathbf{x})\|_2 \leq L$ for all $\mathbf{x} \in \mathcal{D}$.

In this paper, we assume that the strategy activation function h_v is monotone and DR-submodular, which implies that the influence spread function g is also monotone and DR-submodular, same as assumed in [12, 7]. It is reasonable in that, with more marketing effort, the probability of seed activation would increase (monotonicity) but the marginal effect may be decreasing. For the case of independent strategy activation ($h_v(\mathbf{x}) = 1 - \prod_{j \in [d]} (1 - q_{v,j}(x_j))$), we assume $q_{v,j}$ is non-decreasing and concave, which implies that h_v is monotone and DR-submodular [7]. In term of the cost function c , we assume that c is convex and L_c -Lipschitz. The most common function is the simple summation (or 1-norm of \mathbf{x}): $c(\mathbf{x}) = \sum_{i \in [d]} x_i$, but more general convex functions are also common in the economics literature (e.g. [15]), for example $c(x) = \|\mathbf{x}\|_2$.

An important remark is now in order. When h_v 's are monotone and DR-submodular and c is convex, g is monotone and DR-submodular and s is concave, and as a result $g + s$ may be neither monotone nor DR-submodular. This means the greedy hill-climbing algorithm of [12, 7] no longer has theoretical approximation guarantee for the CIM-BS problem. This motivates us to apply the gradient method to solve CIM-BS.

3 Gradient Method with Reverse Influence Sampling

Gradient method has been applied to many continuous optimization problems. For our CIM-BS problem, a natural option is to apply the gradient method directly on the objective function $g + s$. However, the influence spread function g is a complicated combinatorial function, such that its gradient ∇g is too complex to compute in practice. We could apply stochastic gradient on g (see Appendix D) but it has very large variance due to the significant amount of randomness from both strategies activating seeds and influence propagation from seeds, which leads to very slow convergence of the method. Instead, in this section, we propose a novel integration of the gradient method with the reverse influence sampling (RIS) approach [2] for CIM-BS. The key insight is that RIS allows the efficient computation of the exact gradient of an alternative objective function $\hat{g}_R + s$ while maintaining an approximation guarantee of $1/2 - \varepsilon$. Moreover, when independent strategy activation is satisfied by the model, the alternative objective enables a tight concave upper bound, which leads to a $1 - 1/e - \varepsilon$ approximation.

In Section 3.1, we first review existing results on RIS with the continuous domain. Then in Sections 3.2 and 3.3, we present the gradient method, its integration with RIS, and its theoretical analysis, which are our main technical contribution.

3.1 Properties of the Reverse Reachable Sets

The central concept in the RIS approach is the *reverse reachable set*, as defined below.

Definition 2 (Reverse Reachable Set). *Under the triggering model, a reverse reachable (RR) set with a root node v , denoted R_v , is the random set of nodes that v reaches in one reverse propagation: sample all triggering sets $T_u, u \in V$, such that edges $\{(w, u) | u \in V, w \in T_u\}$ together with nodes V form a live-edge graph, and R_v is the set of nodes that can reach v (or v can reach reversely) in this live-edge graph. An RR set R without specifying a root is one with root v selected uniformly at random from V .*

An RR set R_v includes nodes that would activate v in one sample propagation. Then, the key insight is that for a collection of RR sets, if some node u appears in many of these RR sets, it means u is likely to activate many nodes, and thus has high influence. Technically, RR sets connect with the influence spread of a seed set S with the following equation [2, 23]: $\sigma(S) = \mathbb{E}_R[\Pr\{S \cap R \neq \emptyset\}]$. For CIM-BS, we have the following connection as given in [7]:

Algorithm 1 Grad-RIS: Gradient-RIS Meta-Algorithm for CIM-BS.

Input: Directed graph G , triggering model $\{D_v\}_{v \in V}$, domain \mathcal{P} , strategy activation functions $\{h_v\}_{v \in V}$, cost function c , budget k , balance parameter λ , Lipschitz constants L_1, L_2 for the functions $g + s$ and $\hat{g}_{\mathcal{R}} + s$, approximation parameter ε , confidence parameter ℓ , gradient algorithm \mathcal{A} with the approximation guarantee α

Output: A strategy mix \mathbf{x}

- 1: $\mathcal{R}, LB \leftarrow \text{Sampling}(\text{all parameters received})$
 - 2: $\mathbf{x} \leftarrow \mathcal{A}(\mathcal{R}, \hat{g}_{\mathcal{R}} + s, \varepsilon LB) / * \hat{g}_{\mathcal{R}}$ defined in Eq. (2), $s(\mathbf{x}) = \lambda(k - c(\mathbf{x})) *$
 - 3: **return** \mathbf{x}
-

Lemma 1 ([7]). *For any strategy $\mathbf{x} \in \mathcal{P}$, we have $g(\mathbf{x}) = n \cdot \mathbb{E}_R [1 - \prod_{u \in R} (1 - h_u(\mathbf{x}))]$.*

Intuitively, the above lemma means that a node $u \in R$ would activate R 's root if u itself is activated, which happens with probability $h_u(\mathbf{x})$, and thus strategy mix successfully activate R 's root with probability $1 - \prod_{u \in R} (1 - h_u(\mathbf{x}))$. We can generate θ independent RR-sets $\mathcal{R} = \{R_1, \dots, R_\theta\}$, and take the average among them as defined below:

$$\hat{g}_{\mathcal{R}}(\mathbf{x}) = \frac{n}{\theta} \sum_{R \in \mathcal{R}} \left(1 - \prod_{v \in R} (1 - h_v(\mathbf{x})) \right). \quad (2)$$

We can see that $\hat{g}_{\mathcal{R}}(\mathbf{x})$ is an unbiased estimator of $g(\mathbf{x})$. If θ is large enough, then $\hat{g}_{\mathcal{R}}(\mathbf{x})$ should be close to $g(\mathbf{x})$ at every $\mathbf{x} \in \mathcal{P}$. We also have the following lemma from [7].

Lemma 2 ([7]). *If h_v is monotone and DR-submodular for all $v \in V$, then g and $\hat{g}_{\mathcal{R}}$ are also monotone and DR-submodular.*

3.2 Algorithmic Framework Integrating Gradient Method with RIS

We first introduce the general algorithmic framework that integrates any gradient algorithm with the RIS approach. We assume a generic gradient algorithm \mathcal{A} that takes a set of RR sets $\mathcal{R} = \{R_1, \dots, R_\theta\}$, an objective function $\hat{g}_{\mathcal{R}} + s$, and an additive error ε as input, and return a solution $\hat{\mathbf{x}}$ that guarantees $(\hat{g}_{\mathcal{R}} + s)(\hat{\mathbf{x}}) \geq \alpha \cdot \max_{\mathbf{y} \in \mathcal{P}} (\hat{g}_{\mathcal{R}} + s)(\mathbf{y}) - \varepsilon$, in time $T = \text{poly}(\frac{1}{\varepsilon})$, where α is some constant approximation ratio. We call such an \mathcal{A} an (α, ε) -approximate gradient algorithm.

Algorithm 1 gives the meta-algorithm. It first calls the **Sampling** procedure to sample enough RR sets \mathcal{R} , together with an estimated lower bound LB of the optimal solution. Then it calls the gradient algorithm \mathcal{A} with \mathcal{R} , using $\hat{g}_{\mathcal{R}} + s$ as the objective function and εLB as the additive error.

The **Sampling** procedure is to sample enough RR sets for the theoretical guarantee. We adapt the sampling procedure of the IMM algorithm [23], as shown in Algorithm 2. We use the IMM sampling procedure mainly because of its clarity in analysis and theoretical guarantee, while other sampling procedures (e.g. [18]) could be adapted too. The main structure of the sampling procedure is the same as in IMM, where we repeatedly halving the guess x_i of the lower bound LB of the optimal budget-balanced influence spread to find a good lower bound estimate, and then use LB to estimate the final number of RR sets needed and regenerate these RR sets (the regeneration is the workaround 1 proposed in [4] to fix a bug in the original IMM). There are two important differences worth to mention. First, in line 8, we call the gradient algorithm \mathcal{A} to find an approximate solution \mathbf{y}_i , which replaces the original greedy algorithm in IMM. Second, and more importantly, the original IMM algorithm works on a finite solution space — at most $\binom{n}{k}$ feasible seed sets of size k , and $\binom{n}{k}$ is used to bound the number of RR sets needed. However, in CIM-BS, we are working on an infinite solution space, and thus we cannot directly have such a bound. To tackle this problem, we utilize the concept of ε -net and covering number to turn the infinite solution space into a finite space:

Definition 3 (ε -Net and Covering Number). *A finite set N is called an ε -net for \mathcal{P} if for every $\mathbf{x} \in \mathcal{P}$, there exists $\pi(\mathbf{x}) \in N$ such that $\|\mathbf{x} - \pi(\mathbf{x})\|_2 \leq \varepsilon$. The smallest cardinality of an ε -net for \mathcal{P} is called the covering number: $\mathcal{N}(\mathcal{P}, \varepsilon) = \inf\{|N| : N \text{ is an } \varepsilon\text{-net of } \mathcal{P}\}$.*

Algorithm 2 Sampling Procedure.

Input: Same as in Algorithm 1

Output: The RR-sets \mathcal{R} and an estimated lower bound LB

- 1: $LB \leftarrow 1, \mathcal{R}_0 \leftarrow \phi, \theta_0 = 0, \varepsilon' \leftarrow \sqrt{2}\varepsilon/3$
 - 2: **for** $i = 1, 2, \dots, \lfloor \log_2(n + \lambda k) - 1 \rfloor$ **do**
 - 3: $x_i \leftarrow (n + \lambda k)/2^i$
 - 4: $\theta_i \leftarrow \lceil n(2 + \frac{2}{3}\varepsilon') \rceil$
 - 5:
$$\left\lceil \frac{\ln \mathcal{N}(\mathcal{P}, \frac{\varepsilon/3}{L_2} x_i) + \ell \ln n + \ln 2 + \ln \log_2(n + \lambda k)}{\varepsilon'^2 x_i} \right\rceil$$
 - 6: Generate $\theta_i - \theta_{i-1}$ independent RR-sets \mathcal{R}'
 - 7: $\mathcal{R}_i \leftarrow \mathcal{R}' \cup \mathcal{R}_{i-1}$
 - 8: $\mathbf{y}_i \leftarrow \mathcal{A}(\mathcal{R}_i, \hat{g}_{\mathcal{R}_i} + s, \varepsilon x_i/3)$
 - 9: **if** $(\hat{g}_{\mathcal{R}_i} + s)(\mathbf{y}_i) \geq (1 + \varepsilon' + \varepsilon/3) \cdot x_i$ **then**
 - 10: $LB \leftarrow \frac{(\hat{g}_{\mathcal{R}_i} + s)(\mathbf{y}_i)}{1 + \varepsilon' + \varepsilon/3}; \theta \leftarrow \theta_i$
 - 11: **break**
 - 12: **end if**
 - 13: **end for**
 - 14: $\theta^{(1)} = \frac{8n \cdot \ln(4n^\ell)}{LB \cdot (\alpha - \varepsilon/3)^2 \varepsilon^2/9}$
 - 15: $\theta^{(2)} = \frac{2\alpha' \cdot n \cdot \ln(4n^\ell \mathcal{N}(\mathcal{P}, \frac{\varepsilon/3}{L_1 + L_2} LB))}{(\varepsilon/3 - \frac{1}{4}(\alpha - \varepsilon/3)^2 \varepsilon/3)^2 LB}$
 - 16: $\tilde{\theta} \leftarrow \max\{\theta^{(1)}, \theta^{(2)}\}$.
 - 17: Generate $\tilde{\theta}$ independent RR-sets $R_1, \dots, R_{\tilde{\theta}}, \mathcal{R} \leftarrow \{R_1, \dots, R_{\tilde{\theta}}\}$
 - 18: **return** (\mathcal{R}, LB)
-

As a concrete example, suppose we have 1-norm or 2-norm cost function $c(\mathbf{x}) = \|\mathbf{x}\|_1$ or $\|\mathbf{x}\|_2$. With budget k , we know that \mathcal{P} is bounded by the ball $\mathbb{B}_1(k)$ or $\mathbb{B}_2(k)$ with radius k . Then, As shown in [25], the covering number satisfies $\mathcal{N}(\mathcal{P}, \varepsilon) \leq \mathcal{N}(\mathbb{B}_1(k), \varepsilon) \leq \mathcal{N}(\mathbb{B}_2(k), \varepsilon) \leq (3k/\varepsilon)^d$.

Besides the ε -net, we also need to have the upper bounds L_1 and L_2 on the Lipschitz constants of functions $g + s$ and $\hat{g}_{\mathcal{R}} + s$. We defer the discussion on L_1 and L_2 to the next subsection. Covering number and L_1, L_2 together are used to bound the number of RR sets needed, as used in lines 4 and 15 of Algorithm 2. We denote Algorithms 1 and 2 together as Grad-RIS, and we show that Grad-RIS achieves the following approximation guarantee:

Theorem 1. *For any $\varepsilon, \ell, \alpha > 0$, for any $(\alpha, \varepsilon/3)$ -approximate gradient algorithm \mathcal{A} for $\hat{g}_{\mathcal{R}} + s$, with probability at least $1 - \frac{1}{n^\ell}$, Grad-RIS outputs a solution \mathbf{x} that is an $(\alpha - \varepsilon)$ -approximation of the optimal solution OPT_{g+s} of CIM-BS, i.e. $(g + s)(\mathbf{x}) \geq (\alpha - \varepsilon)OPT_{g+s}$.*

The proof of the above theorem follows the proof structure of IMM [23], where the number of the RR sets needed is carefully adapted to accommodate the covering number of the ε -net, and the Lipschitz constants of the objective functions.

3.3 Gradient Algorithms

In this subsection, we will show two gradient algorithms that approximately maximize the function $\hat{g}_{\mathcal{R}}(\mathbf{x}) + s(\mathbf{x})$. They are the instantiations of the generic algorithm \mathcal{A} in Section 3.2: the first one works on the general model and uses proximal gradient to achieve $(\frac{1}{2}, \varepsilon)$ -approximation, while the second one works on the special case of independent strategy activation and uses gradient on a concave upper bound to achieve $(1 - \frac{1}{\varepsilon}, \varepsilon)$ -approximation.

General Case: ProxGrad-RIS. We first consider the general case where the strategy activation functions h_v 's are monotone, DR-submodular, L_h -Lipschitz and β_h -smooth, and the cost function c is convex and

L_c -Lipschitz. In this case, we have that g and $\hat{g}_{\mathcal{R}}$ are monotone and DR-submodular (Lemma 2), and budget-saving function s is concave. To solve this problem, we adapt the (stochastic) proximal gradient algorithm [20, 19] to provide a $\frac{1}{2}$ -approximate solution to the following problem: given a convex set \mathcal{P} , a β -smooth, non-negative, monotone, and DR-submodular function $f_1(\mathbf{x})$ on \mathcal{P} and a non-negative and concave function $f_2(\mathbf{x})$ on \mathcal{P} , find a solution in \mathcal{P} maximizing $f_1(\mathbf{x}) + f_2(\mathbf{x})$. The original proximal gradient is for the case when both f_1 and f_2 are concave, and we adapt it to the case when f_1 is monotone and DR-submodular to provide an approximate solution. The reason we use proximal gradient is that our budget-saving function s may not be smooth (e.g. when the cost function is the 2-norm function). We present the general solution first, since it may be of independent interest. The following is the iteration procedure for the stochastic proximal gradient algorithm.

$$\begin{cases} \mathbf{x}^{(t+1)} = \text{prox}_{-\eta_t f_2}(\mathbf{x}^{(t)} + \eta_t \mathbf{v}^{(t)}), \\ \text{where } \mathbb{E}[\mathbf{v}^{(t)}] = \nabla f_1(\mathbf{x}^{(t)}), \\ \text{prox}_{\phi}(\mathbf{x}) := \text{argmin}_{\mathbf{y} \in \mathcal{P}} (\phi(\mathbf{y}) + \frac{1}{2} \|\mathbf{x} - \mathbf{y}\|_2^2), \\ \text{for any convex function } \phi, \end{cases} \quad (3)$$

where $\eta_t \leq \frac{1}{\beta}$ is the step size and $\mathbf{v}^{(t)}$ is the stochastic gradient at $\mathbf{x}^{(t)}$. We use Δ to denote an upper bound of the diameter of \mathcal{P} , i.e. $\Delta \geq \max_{\mathbf{x}, \mathbf{y} \in \mathcal{P}} \|\mathbf{x} - \mathbf{y}\|_2$. The following is the main result for the above stochastic proximal gradient algorithm, with its proof adapted from the original proof.

Theorem 2. *Suppose that \mathcal{P} is a convex set, function $f_1(\mathbf{x})$ is β -smooth, non-negative, monotone, and DR-submodular on \mathcal{P} , $f_2(\mathbf{x})$ is non-negative and concave on \mathcal{P} . Let \mathbf{x}^* be the point that maximizes $f_1(\mathbf{x}) + f_2(\mathbf{x})$. Suppose that for some $\sigma > 0$, the stochastic gradient $\mathbf{v}^{(t)}$ satisfies $\mathbb{E}\|\mathbf{v}^{(t)} - \nabla f_1(\mathbf{x}^{(t)})\|_2^2 \leq \sigma^2$ for all t , then for all $T > 0$, if we set $\eta_t = \eta = 1/(\beta + \frac{\sigma}{\Delta} \sqrt{2T})$, and iterate as shown in (3) starting from $\mathbf{x}^{(0)} \in \mathcal{P}$, we have*

$$\begin{aligned} & \mathbb{E} \left[\max_{t=0,1,2,\dots,T} (f_1 + f_2)(\mathbf{x}^{(t)}) \right] \\ & \geq \frac{1}{2} (f_1 + f_2)(\mathbf{x}^*) - \frac{\beta \Delta^2}{4T} - \frac{\sigma \Delta}{\sqrt{2T}}. \end{aligned}$$

Note that if we use exact gradient instead of the stochastic gradient, we simply set $\sigma = 0$ in the above theorem. To apply the proximal gradient algorithm and Theorem 2 to maximize $\hat{g}_{\mathcal{R}} + s$, we compute the exact gradient of $\hat{g}_{\mathcal{R}}$ and also derive the Lipschitz and smoothness constants, as shown below. For RR set sequence $\mathcal{R} = \{R_1, \dots, R_\theta\}$, let $\nu^{(1)}(\mathcal{R}) = \sum_{R \in \mathcal{R}} |R|/\theta$ be the average RR set size in \mathcal{R} , $\nu^{(2)}(\mathcal{R}) = \sum_{R \in \mathcal{R}} |R|^2/\theta$ be the average squared size in \mathcal{R} , and $\nu^{(3)}(\mathcal{R}) = \sum_{R \in \mathcal{R}} |R|^3/\theta$ be the average cubed size.

Lemma 3. *If functions $h_v(\mathbf{x})$'s are L_h -Lipschitz, then function $\hat{g}_{\mathcal{R}}(\mathbf{x})$ is $(\nu^{(1)}(\mathcal{R})nL_h)$ -Lipschitz, and function $g(\mathbf{x})$ is (n^2L_h) -Lipschitz. If functions $h_v(\mathbf{x})$'s are β_h -smooth, then function $\hat{g}_{\mathcal{R}}(\mathbf{x})$ is $(\nu^{(1)}(\mathcal{R})n\beta_h + \nu^{(2)}(\mathcal{R})nL_h^2)$ -smooth. The gradient of function $\hat{g}_{\mathcal{R}}(\mathbf{x})$ is*

$$\nabla \hat{g}_{\mathcal{R}}(\mathbf{x}) = \frac{n}{\theta} \sum_{R \in \mathcal{R}, v' \in R} \nabla h_{v'}(\mathbf{x}) \prod_{v \in R, v \neq v'} (1 - h_v(\mathbf{x})). \quad (4)$$

With Lemma 3 and Theorem 2, we can conclude the gradient algorithm \mathcal{A} working with Grad-RIS with the following settings: (a) we use the proximal gradient iteration given in Eq. (3), with stochastic gradient $\mathbf{v}^{(t)}$ replaced with the exact gradient $\nabla \hat{g}_{\mathcal{R}}(\mathbf{x}^{(t)})$ as given in Eq. (4); (b) we set step size $\eta_t = 1/(\nu^{(1)}(\mathcal{R})n\beta_h + \nu^{(2)}(\mathcal{R})nL_h^2)$ when calling the algorithm with \mathcal{R} ; (c) we set number of steps $T = 3(\nu^{(1)}(\mathcal{R})n\beta_h + \nu^{(2)}(\mathcal{R})nL_h^2) \cdot \Delta^2/4\epsilon$; (d) we use $L_1 = L_2 = n^2L_h + \lambda L_c$ (since $n \geq \nu^{(1)}(\mathcal{R})$) as parameters in Grad-RIS. We refer to the full algorithm with the above setting as ProxGrad-RIS. The following theorem summarizes the approximation guarantee of ProxGrad-RIS.

Theorem 3. *For any $\epsilon, \ell > 0$, with probability at least $1 - \frac{1}{n^\ell}$, ProxGrad-RIS outputs a solution \mathbf{x} that is a $(\frac{1}{2} - \epsilon)$ -approximation of the optimal solution OPT_{g+s} of CIM-BS, i.e. $(g + s)(\mathbf{x}) \geq (\frac{1}{2} - \epsilon)OPT_{g+s}$.*

We remark that the actual computation of the proximal step $\text{prox}_{-\eta_t f_2}(\cdot)$ in Eq.(3) depends on domain \mathcal{D} and cost function c . When $\mathcal{D} = \mathbb{R}_+^d$ and c is 1-norm or 2-norm function, we can derive efficient algorithm for the proximal step, as summarized below.

Lemma 4. *When $c(\mathbf{x}) = \|\mathbf{x}\|_1$ and $\mathcal{D} = \mathbb{R}_+^d$, the proximal step can be done in time $O(d \log d)$. When $c(\mathbf{x}) = \|\mathbf{x}\|_2$ and $\mathcal{D} = \mathbb{R}_+^d$, the proximal step can be done in $O(d)$.*

Independent Strategy Activation Case: UpperGrad-RIS. Next, we introduce an $(1 - \frac{1}{e})$ -approximation for maximizing $\hat{g}_{\mathcal{R}} + s$ under the case of independent strategy activation. Recall that in the independent strategy activation case, each function $h_v(\mathbf{x}) = 1 - \prod_{j \in [d]} (1 - q_{v,j}(x_j))$, where $q_{v,j}(x_j)$ is monotone and concave in x_j . In this case, we can write $\hat{g}(\mathbf{x})$ into the following form.

$$\begin{aligned} \hat{g}_{\mathcal{R}}(\mathbf{x}) &= \frac{n}{\theta} \sum_{R \in \mathcal{R}} \left(1 - \prod_{v \in R} (1 - h_v(\mathbf{x})) \right) \\ &= \frac{n}{\theta} \sum_{R \in \mathcal{R}} \left(1 - \prod_{v \in R} \left(\prod_{j \in [d]} (1 - q_{v,j}(\mathbf{x})) \right) \right). \end{aligned}$$

The above form makes $\hat{g}_{\mathcal{R}}(\mathbf{x})$ belong to coverage functions, which has the following concave upper and lower bounds ([10]):

Proposition 1 ([10]). *For any $\mathbf{x} \in [0, 1]^l$, let $\alpha(\mathbf{x}) = 1 - \prod_{i=1}^l (1 - x_i)$ and $\beta(\mathbf{x}) = \min\{1, \sum_{i=1}^l x_i\}$, we have $(1 - 1/e)\beta(\mathbf{x}) \leq \alpha(\mathbf{x}) \leq \beta(\mathbf{x})$.*

By Proposition 1, we can optimize $\bar{g}_{\mathcal{R}} + s$ where $\bar{g}_{\mathcal{R}}(\mathbf{x})$ is the upper bound of $\hat{g}_{\mathcal{R}}(\mathbf{x})$ defined as:

$$\bar{g}_{\mathcal{R}}(\mathbf{x}) := \frac{n}{\theta} \sum_{R \in \mathcal{R}} \min\{1, \sum_{j \in [d], v \in R} q_{v,j}(\mathbf{x})\}. \quad (5)$$

Since the function $g_{\mathcal{R}}(\mathbf{x})$ is non-smooth, we use the projected subgradient method to maximize the function $\bar{g}_{\mathcal{R}} + s$ [17], as summarized by the following lemma.

Lemma 5. *In the case of independent strategy activation, suppose that $\bar{g}_{\mathcal{R}} + s$ is $L_{\bar{g}+s}$ -Lipschitz. If we use projected subgradient descent to optimize the function $(\bar{g}_{\mathcal{R}} + s)(\mathbf{x})$ with step size $\eta_t = \frac{\Delta}{L_{\bar{g}+s}\sqrt{t}}$ and let \mathbf{y} denote the output where $(\bar{g}_{\mathcal{R}} + s)(\mathbf{y}) \geq \max_{\mathbf{x} \in \mathcal{P}} (\bar{g}_{\mathcal{R}} + s)(\mathbf{x}) - \varepsilon$. Then $(\hat{g}_{\mathcal{R}} + s)(\mathbf{y}) \geq (1 - \frac{1}{e}) \max_{\mathbf{x} \in \mathcal{P}} (\hat{g}_{\mathcal{R}} + s)(\mathbf{x}) - \varepsilon$, and \mathbf{y} can be solved in $\frac{(\Delta L_{\bar{g}+s})^2}{\varepsilon^2}$ iterations.*

The following lemma presents the Lipschitz constants and subgradients needed in Lemma 5.

Lemma 6. *Suppose that functions $q_{v,j}(x_j)$'s are L_q -Lipschitz, then function $g(\mathbf{x})$ is $n^2 \sqrt{d} L_q$ -Lipschitz, and functions $\hat{g}_{\mathcal{R}}(\mathbf{x})$ and $\bar{g}_{\mathcal{R}}(\mathbf{x})$ are $\nu^{(1)}(\mathcal{R}) n \sqrt{d} L_q$ -Lipschitz. The subgradient of the function $\bar{g}_{\mathcal{R}}(\mathbf{x})$ is*

$$\frac{n}{\theta} \sum_{R \in \mathcal{R}} \begin{cases} 0, & \text{if } \sum_{j \in [d], v \in R} q_{v,j}(x_j) \geq 1, \\ \sum_{v \in R, j \in [d]} \nabla q_{v,j}(x_j), & \text{if } \sum_{j \in [d], v \in R} q_{v,j}(x_j) < 1. \end{cases} \quad (6)$$

Combining Lemma 6 with Lemma 5, we can conclude our subgradient algorithm based on the upper bound function $\bar{g}_{\mathcal{R}} + s$: (a) we use the projected subgradient algorithm with the subgradient of $\bar{g}_{\mathcal{R}}$ given in Eq.(6); (b) we set step size $\eta_t = \Delta / (\nu^{(1)}(\mathcal{R}) n \sqrt{d} L_q \sqrt{t})$; (c) we use $T = 9(\Delta \nu^{(1)}(\mathcal{R}) n \sqrt{d} L_q + \lambda L_c)^2 / \varepsilon^2$ iterations to get ε accuracy; and (d) we set $L_1 = L_2 = n^2 \sqrt{d} L_q + \lambda L_c$ in Grad-RIS. We refer to the full algorithm with the above setting as UpperGrad-RIS. The following theorem summarizes the approximation guarantee of UpperGrad-RIS.

Theorem 4. For any $\varepsilon, \ell > 0$, with probability at least $1 - \frac{1}{n^\ell}$, UpperGrad-RIS outputs a solution \mathbf{x} that is an $(1 - 1/e - \varepsilon)$ -approximation of the optimal solution OPT_{g+s} of CIM-BS, i.e. $(g+s)(\mathbf{x}) \geq (1 - 1/e - \varepsilon)OPT_{g+s}$.

Total Time Complexity. For the time complexity of ProxGrad-RIS and UpperGrad-RIS, we make the following reasonable assumptions: (1) the time for sampling a trigger set $T_v \sim D_v$ is proportional to the in-degree of v ; (2) the optimal influence spread $\max_{\mathbf{x} \in \mathcal{P}} g(\mathbf{x})$ among strategy mixes is at least the optimal single node influence spread $\max_{v \in V} \sigma(\{v\})$; and (3) $\lambda k \leq n$, otherwise the budget saving is more important than influencing the entire network, and CIM-BS problem no longer makes much sense. The following theorem summarizes the time complexity result when $\mathcal{D} = \mathbb{R}_+^d$ and $c(\mathbf{x}) = \|\mathbf{x}\|_1$ or $c(\mathbf{x}) = \|\mathbf{x}\|_2$. The more general result is given in Appendix C. Notation $\tilde{O}(\cdot)$ ignores poly-logarithmic factors.

Theorem 5. Suppose that $\mathcal{D} = \mathbb{R}_+^d$ and $c(\mathbf{x}) = \|\mathbf{x}\|_1$ or $\|\mathbf{x}\|_2$, $h_v(\mathbf{x})$'s are L_h -Lipschitz and β_h -smooth. If $\nabla h_v(\mathbf{x})$ can be computed in time T_h , the expected running time of ProxGrad-RIS is bounded by $\tilde{O}\left(\frac{\beta_h n^2 + L_h^2 n^3}{\varepsilon} \cdot \frac{T_h(m+n) \cdot (d+\ell)}{\varepsilon^2}\right)$. Under independent strategy activation, if $q_{v,j}(x_j)$'s are L_q -Lipschitz and the gradient and function value of $q_{v,l}(x_j)$ can be computed in time T_q , the expected running time of UpperGrad-RIS is bounded by $\tilde{O}\left(\frac{n^4 d L_q^2}{\varepsilon^2} \cdot \frac{T_q(m+n) \cdot (d+\ell)}{\varepsilon^2}\right)$.

From the time complexity result, we can see that the two gradient algorithms still have high-order dependency on the graph size. This is mainly because we need the conservative bounds on the number of gradient algorithm iterations for the theoretical guarantee (terms $\frac{\beta_h n^2 + L_h^2 n^3}{\varepsilon}$ and $\frac{n^4 d L_q^2}{\varepsilon^2}$). In our actual algorithms, we already use $\nu^{(1)}(\mathcal{R})$ and $\nu^{(2)}(\mathcal{R})$ instead of n and n^2 in the upper bound of the gradient decent steps for $\hat{g}_{\mathcal{R}} + s$, so our actual performance would be reduced by corresponding factors. For details, please see Appendix C.2 for the results and discussions on using the moments of RR set size in the time complexity bounds.

4 Experiments

Experiment setup. We test on two network dataset. The first dataset is the DM network, which is a network of data mining researchers extracted from the ArnetMiner archive (arnetminer.org), with 679 nodes and 3,374 edges, and edge weights are learned from a topic affinity model and obtained from the authors [22]. The second dataset is NetHEPT, a popular dataset used in many influence maximization studies (e.g. [6, 26, 23]). It is an academic collaboration network from the ‘‘High Energy Physics Theory’’ section of arXiv from 1991 to 2003, where nodes represent the authors and each edge represents one paper co-authored by two nodes. After removing duplicated edges, we have 15,233 nodes and 62,774 directed edges. The influence probabilities on edges are assigned according to the weighted cascade setting [12]: the influence probability of edge (u, v) is $1/d_v$, where d_v is the in-degree of v .

Besides our ProxGrad-RIS and UpperGrad-RIS algorithms, we test two more algorithms: (a) ProxGrad-Org: stochastic proximal gradient algorithm on the original objective function, and the stochastic gradient computation as well as step size and step count settings are given in Appendix D. By Theorem 2, ProxGrad-Org would achieve 1/2 approximation in expectation. (b) Greedy-RIS: simply replace the gradient algorithm \mathcal{A} in Grad-RIS with the greedy algorithm for the objective $\hat{g}_{\mathcal{R}}(\mathbf{x}) + s(\mathbf{x})$ on generated RR sets \mathcal{R} , and the greedy algorithm stops either when the budget is exhausted or the marginal gain is negative. This is similar to the algorithm in [7], but since $\hat{g}_{\mathcal{R}}(\mathbf{x}) + s(\mathbf{x})$ is neither monotone nor DR-submodular, Greedy-RIS has no theoretical guarantee and it is only a heuristic algorithm for our tests. For the three gradient-based algorithms ProxGrad-RIS, UpperGrad-RIS, and ProxGrad-Org, we further test their heuristic versions that may lead to faster running time: instead of using a conservative number of iteration steps for theoretical guarantees, we heuristically terminate the gradient iteration if the difference in the objective function values for two consecutive iterations is within a small value of 0.3 (we will justify the choice of this parameter in our tests). We put suffix HEU for the three versions of the heuristic gradient termination algorithms.

For parameter settings, we set $\varepsilon = 0.3$ and $\ell = 1$ for all algorithms. For Greedy-RIS, we set the greedy step size to be 0.1 on each dimension. For 1-norm cost function ($c(\mathbf{x}) = \|\mathbf{x}\|_1$), we test (a) vary k from 5 to 50 while

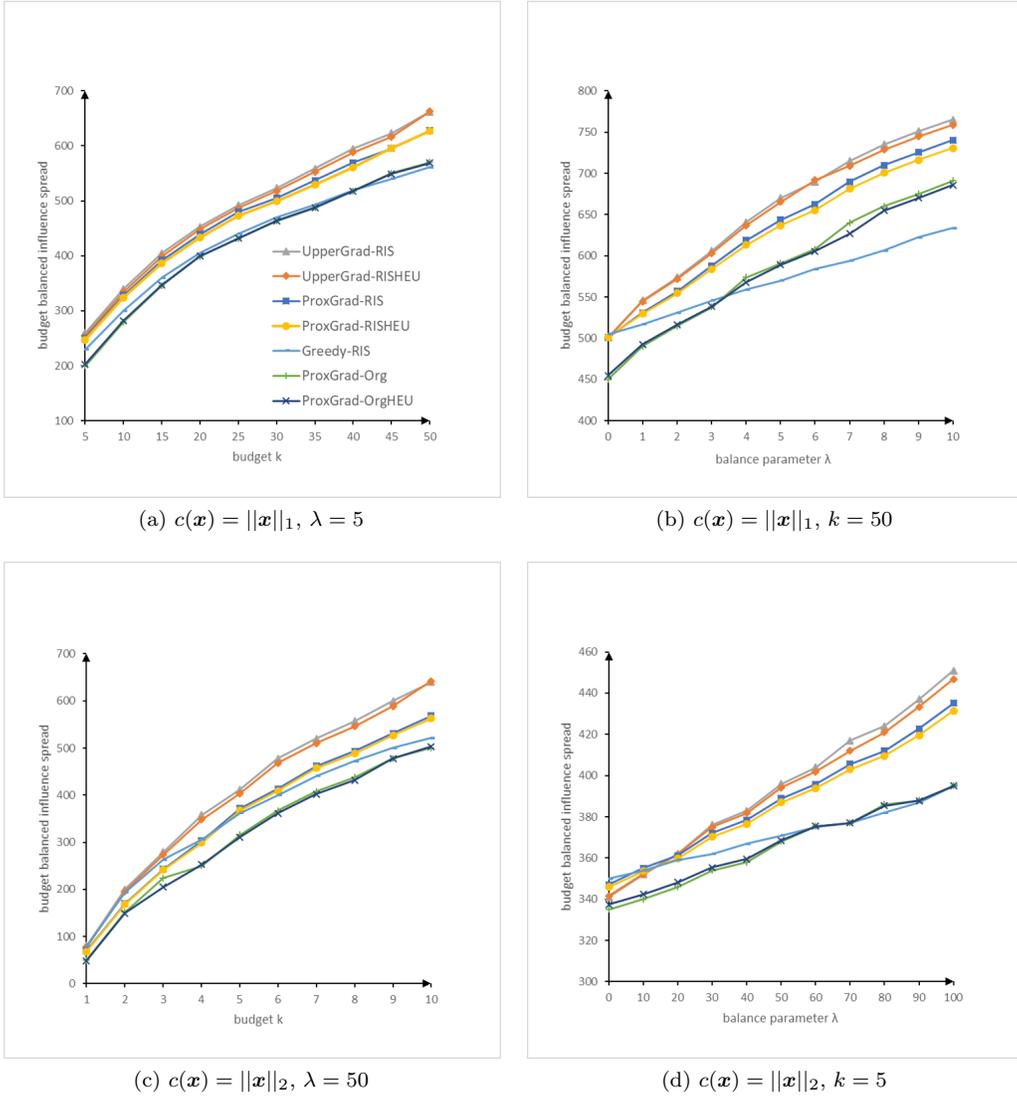


Figure 1: Budget balanced influence spread results for the personalized marketing scenario on the DM dataset. The legends shown in (a) apply to all other figures.

keeping $\lambda = 5$, and (b) vary λ from 0 to 10 while keeping $k = 50$. For 2-norm cost function ($c(\mathbf{x}) = \|\mathbf{x}\|_2$), we test (c) varying vary k from 1 to 10 while keeping $\lambda = 50$, and (d) vary λ from 0 to 100 while keeping $k = 5$. The reason we use a smaller budget k for 2-norm is because $\|\mathbf{x}\|_2 \leq \|\mathbf{x}\|_1/\sqrt{d}$, and thus we need a significantly small budget for 2-norm in order to have a similar feasible region. Parameter λ is adjusted accordingly so that $\lambda \cdot k$ is at the same scale as the influence spread, otherwise either influence spread or budget saving is dominant, and the problem is degenerated. For functions $h_v(\mathbf{x})$, we test two cases: the personalized marketing case and the segment marketing case [27, 7]. In the personalized marketing case, each node v receives a separate discount $x_v \in [0, 1]$. This corresponds to the independent strategy activation case with $d = n$, and $q_{v,j}(x_j) > 0$ only when $j = v$. We set $q_{v,v}(x_v) = 2x_v - x_v^2$ as in [27, 7]. In the segment marketing case, we have 10 strategies in total, i.e $d = 10$, each strategy targets to a disjoint segment of users, and each user has exactly one corresponding strategy. Each user is randomly put into one of the 10 user segments with equal probability, and if one segment in the end has less than 50 or larger than 80 users, we

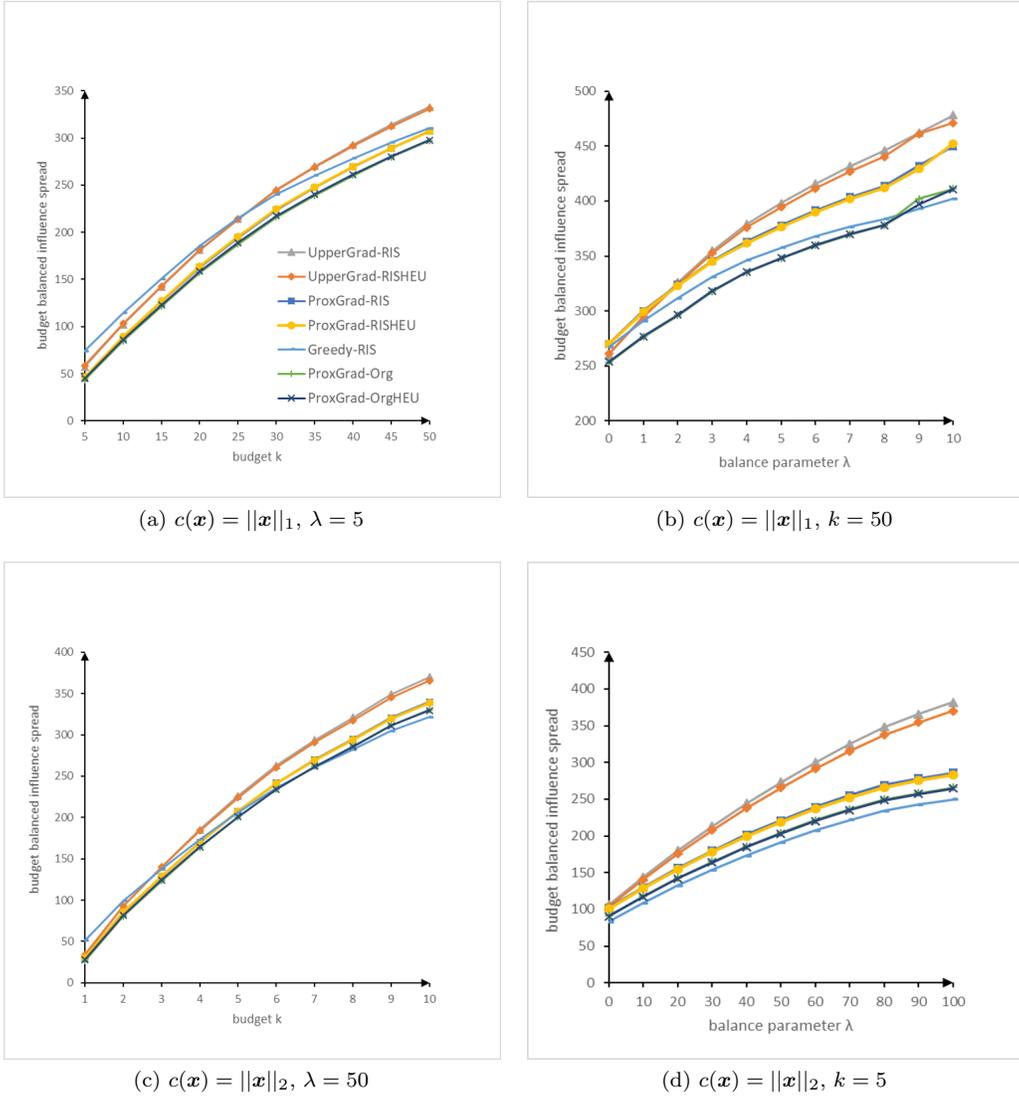


Figure 2: Budget balanced influence spread results for the segment marketing scenario on the DM dataset. The legends shown in (a) apply to all other figures.

regenerate the user segments, so that in the end all segments have sizes within 50 to 80.

The experiments are run on a Ubuntu 17.04 server machine with 2.9GHz and 128GB memory. The code is written in C++ and compiled by g++.

Experimental results. We first show the results on the DM dataset. Figure 1 shows the influence spread results of the personalized marketing scenario, and Figure 2 shows the influence results of the segment marketing scenario, both on the DM dataset. Each data point on an influence spread curve is the average of five solutions found by five runs of the same algorithm, and the influence spread of each solution is an average of 1000 simulation runs. In all cases, UpperGrad-RIS/UpperGrad-RISHEU has the best performance, followed by ProxGrad-RIS/ProxGrad-RISHEU, which coincides with our theoretical analysis that UpperGrad-RIS has a better theoretical guarantee. Both algorithms outperform two baselines in most cases, especially when λ is getting large. Large λ indicates that we need to pay more attention to budget saving, and thus the result suggests that our algorithm handles much better in the balance between influence spread and budget saving.

Comparing the heuristic termination version of each gradient-based algorithm with their corresponding theory-guided termination, the heuristic versions almost match the theoretical version in influences spread in all cases, showing that the heuristic termination seems to perform well in practice.

Table 1: Running time results for the personalized marketing scenario on the DM dataset (in seconds).

| | $c(\mathbf{x}) = \ \mathbf{x}\ _1, k = 50$ and $\lambda = 5$ | $c(\mathbf{x}) = \ \mathbf{x}\ _2, k = 5$ and $\lambda = 50$ |
|------------------|--|--|
| ProxGrad-RIS | 33.2 | 27.3 |
| ProxGrad-RISHEU | 6.5 | 5.5 |
| UpperGrad-RIS | 81.8 | 72.3 |
| UpperGrad-RISHEU | 13.2 | 13.9 |
| Greedy-RIS | 10.2 | 8.7 |
| ProxGrad-Org | 1043.9 | 1021.6 |
| ProxGrad-OrgHEU | 243.4 | 187.8 |

Table 2: Running time results for the segment marketing scenario on the DM dataset (in seconds).

| | $c(\mathbf{x}) = \ \mathbf{x}\ _1, k = 50$ and $\lambda = 5$ | $c(\mathbf{x}) = \ \mathbf{x}\ _2, k = 5$ and $\lambda = 50$ |
|------------------|--|--|
| ProxGrad-RIS | 0.58 | 0.46 |
| ProxGrad-RISHEU | 0.12 | 0.11 |
| UpperGrad-RIS | 1.9 | 2.4 |
| UpperGrad-RISHEU | 0.32 | 0.77 |
| Greedy-RIS | 0.12 | 0.13 |
| ProxGrad-Org | 28.6 | 19.2 |
| ProxGrad-OrgHEU | 4.8 | 4.3 |

Table 1 shows the running time of the personalized marketing scenario, and Table 2 shows the running time of the segment marketing scenario, Each running time number is the average of five runs. The result shows that ProxGrad-RIS and UpperGrad-RIS are 9 to 49 times faster than ProxGrad-Org. This is mainly due to the high variance in the stochastic gradient for the original objective function, as we discussed before. Moreover, ProxGrad-RIS and UpperGrad-RIS is slower than Greedy-RIS. This is mainly because our conservative bounds on the number of gradient iterations make ProxGrad-RIS and UpperGrad-RIS slow, while Greedy-RIS only use the heuristic greedy approach with step size 0.1 without any theoretical guarantee. Indeed Greedy-RIS is inferior to ProxGrad-RIS and UpperGrad-RIS in terms of the influence spread achieved.

The heuristic termination significantly improves the running time. Comparing against their respective theory-guided termination counterparts, we can see that heuristic termination in general improves the running time for 5 — 7 times. Comparing against the Greedy-RIS algorithm, we can see that ProxGrad-RISHEU is now faster than Greedy-RIS and UpperGrad-RISHEU is close to Greedy-RIS in running time. Therefore, this means that our gradient-based algorithms could achieve faster running time with heuristic termination while still providing better influence spread quality than the greedy heuristic, and if we want a theoretical guarantee, we could use more conservative theory-guided termination, which runs a few times slower but provides both theoretical guarantee and best empirical performance on influence spread.

Next, we test on the larger dataset NetHEPT. On this larger dataset, the gradient algorithms with theoretical guarantee is too slow to run, so we only run the heuristic versions of the algorithms and comparing them with the heuristic greedy algorithm. Figure 3 show the result of the budget balanced influence spread versus k and λ respectively, for the case of the personalized marketing scenario. We use basically the same parameter settings as in the DM dataset, except that we try large λ values (e.g. $\lambda = 10$ instead of $\lambda = 5$ as in the DM dataset when varying k), because NetHEPT dataset has larger influence spread, and we need a larger value of λ to balance that. From the result, we can see that the UpperGrad-RISHEU and ProxGrad-RISHEU still perform better than the greedy heuristic. Moreover, the advantage is larger when the balance parameter λ is getting large, similar to the results we see on the DM dataset. Table 3 reports the running time of the algorithms on the NetHEPT dataset. We see that UpperGrad-RISHEU and ProxGrad-RISHEU are a few times

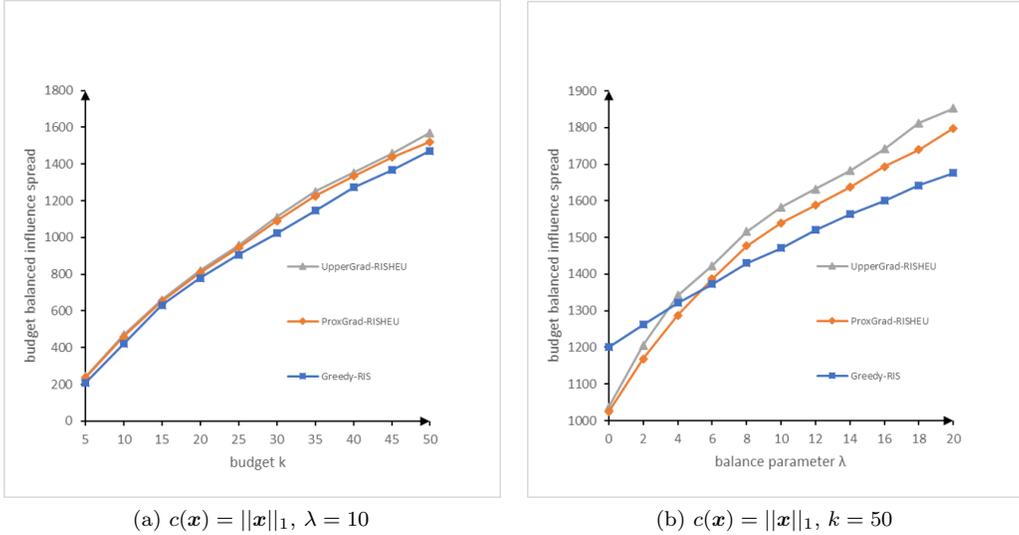


Figure 3: Budget balanced influence spread results for the personalized marketing scenario on the NetHEPT dataset.

slower than the greedy heuristic.

Table 3: Running time results for the personalized marketing scenario on the NetHEPT dataset (in seconds).

| | $c(\mathbf{x}) = \ \mathbf{x}\ _1, k = 50$ and $\lambda = 10$ |
|------------------|---|
| ProxGrad-RISHEU | 964.1 |
| UpperGrad-RISHEU | 1522.7 |
| Greedy-RIS | 334.0 |

Besides looking at the budget balanced influence spread $g(\mathbf{x}) + s(\mathbf{x})$ as a whole, we would also like to decompose this overall objective into the influence spread $g(\mathbf{x})$ and budget saving $s(\mathbf{x})$ and see how each of them behaves, especially when λ changes. Figure 4 (a) shows this test result on the DM dataset with $k = 50$ and $c(\mathbf{x}) = \|\mathbf{x}\|_1$, focusing on the UpperGrad-RIS and Greedy-RIS algorithms. The result shows that when λ increases, the influence spread objective $g(\mathbf{x})$ in general decreases while the budget saving objective $s(\mathbf{x})$ increases, indicating that both algorithms lean towards budget saving when more weight is put on budget saving. Comparing the two algorithms, we clearly see that UpperGrad-RIS put much more emphasis on budget saving than Greedy-RIS, with budget saving objective $s(\mathbf{x})$ more than doubled.

For the heuristic version of our gradient algorithms, we further verify the stopping criteria parameter 0.3 that we use. To do so, we vary this parameter from 0.1 to 1 and see the result comparing to the theory-guided version of the algorithms. Figure 4 (b) shows this test result on the DM dataset with $k = 50$, $\lambda = 5$ and $c(\mathbf{x}) = \|\mathbf{x}\|_1$. The result shows that in general before 0.3 or 0.4, the performance of our heuristic algorithms match very closely with the theory-guided versions, but when the parameter increases to 0.5 or above, the performance of the heuristic algorithms starts to drop significantly. Therefore, in our main experiments, we set this parameter to 0.3.

Finally, we collect the statistics for the first three moments of the average RR set sizes, which are closely related to the running time of the gradient-based algorithms, as discussed at the end of Section 3 and shown in Theorems 8 and 9. In particular, by random sampling 10,000 RR sets and taking the average, we obtain $\nu^{(1)} = \mathbb{E}[|R|] = 7.2$, $\nu^{(2)} = \mathbb{E}[|R|^2] = 62.9$, $\nu^{(3)} = \mathbb{E}[|R|^3] = 501.4$. Following the remark after Theorems 8, we can see that without using these moments in the time complexity bound, we would have relaxed the bound for a large factor. In particular, according to the remark after Theorems 8, the relaxation factors for various

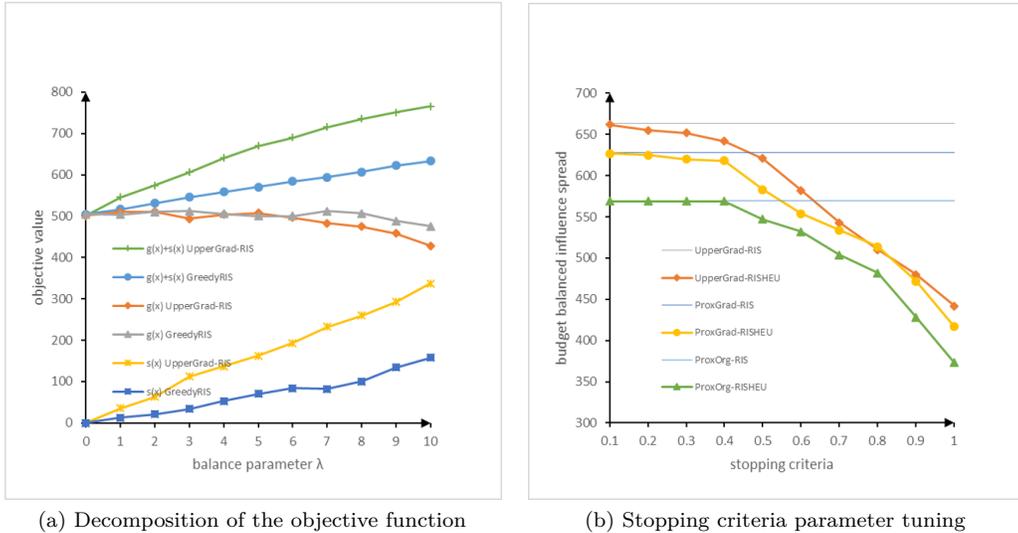


Figure 4: Decomposition of the objective function, and stopping criteria parameter tuning for the heuristic version of algorithms, on the DM dataset. We use $c(\mathbf{x}) = \|\mathbf{x}\|_1$, $k = 50$, and for (b) $\lambda = 5$.

components of computations are: $n/\nu^{(1)} = 94$, $n^2/\nu^{(2)} = 7,330$, $\nu^{(1)} \cdot n/\nu^{(2)} = 78$, and $\nu^{(1)} \cdot n^2/\nu^{(3)} = 6,620$. This suggests that using moments of mean RR set size would significantly reduce the time complexity bound.

5 Conclusion and Further Work

In this paper, we tackle the new problem of continuous influence maximization with budget saving (CIM-BS), whose objective function is neither monotone, nor DR-submodular or concave. We use the gradient method to solve CIM-BS, and provide innovative integration with the reverse influence sampling method to achieve theoretical approximation guarantees. One important direction of future study is to make the gradient method more scalable, which requires more detailed study of convergence behavior and properties of the gradient method in the influence maximization domain. Another direction is to investigate if the gradient method can be applied to other influence maximization settings such as competitive influence maximization. Gradient method is a rich and powerful approach that has been already applied to many application domains, and thus we hope our work could inspire more studies incorporating the gradient method into the influence maximization research.

References

- [1] Andrew An Bian, Baharan Mirzasoleiman, Joachim M. Buhmann, and Andreas Krause. Guaranteed non-convex optimization: Submodular maximization over continuous domains. In *AISTATS*, 2017.
- [2] Christian Borgs, Michael Brautbar, Jennifer Chayes, and Brendan Lucier. Maximizing social influence in nearly optimal time. In *SODA*, 2014.
- [3] Chandra Chekuri, Jan Vondrák, and Rico Zenklusen. Submodular function maximization via the multilinear relaxation and contention resolution schemes. *SIAM Journal on Computing*, 43(6):1831–1879, 2014.
- [4] Wei Chen. An issue in the martingale analysis of the influence maximization algorithm imm. In *CSoNet*, 2019.

- [5] Wei Chen, Laks VS Lakshmanan, and Carlos Castillo. *Information and Influence Propagation in Social Networks*. Morgan & Claypool Publishers, 2013.
- [6] Wei Chen, Yajun Wang, and Siyu Yang. Efficient influence maximization in social networks. In *KDD*, pages 199–208, 2009.
- [7] Wei Chen, Ruihan Wu, and Zheng Yu. Scalable lattice influence maximization. Technical Report arXiv:1802.04555v2, arXiv, 2019.
- [8] Patrick L Combettes and Valérie R Wajs. Signal recovery by proximal forward-backward splitting. *Multiscale Modeling & Simulation*, 4(4):1168–1200, 2005.
- [9] S. Hamed Hassani, Mahdi Soltanolkotabi, and Amin Karbasi. Gradient methods for submodular maximization. In *NIPS*, 2017.
- [10] Mohammad Karimi, Mario Lucic, Hamed Hassani, and Andreas Krause. Stochastic submodular maximization: The case of coverage functions. In *NIPS*, 2017.
- [11] Mohammad Reza Karimi, Mario Lucic, S. Hamed Hassani, and Andreas Krause. Stochastic submodular maximization: The case of coverage functions. In *NIPS*, pages 6856–6866, 2017.
- [12] David Kempe, Jon M. Kleinberg, and Éva Tardos. Maximizing the spread of influence through a social network. *Theory of Computing*, 11(4):105–147, 2015. First appeared in KDD’03.
- [13] Yuchen Li, Ju Fan, Yanhao Wang, and Kian-Lee Tan. Influence maximization on social graphs: A survey. *IEEE Trans. Knowl. Data Eng.*, 30(10):1852–1872, 2018.
- [14] Wei Lu and Laks V S Lakshmanan. Profit maximization over social networks. In *ICDM*, pages 479–488, 2012.
- [15] N Gregory Mankiw. *Principles of economics*. Cengage Learning, 2014.
- [16] Aryan Mokhtari, Hamed Hassani, and Amin Karbasi. Conditional gradient method for stochastic submodular maximization: Closing the gap. In *AISTATS*, 2018.
- [17] Yurii Nesterov. *Introductory lectures on convex optimization: A basic course*, volume 87. Springer Science & Business Media, 2013.
- [18] Hung T. Nguyen, My T. Thai, and Thang N. Dinh. Stop-and-stare: Optimal sampling algorithms for viral marketing in billion-scale networks. In *SIGMOD*, pages 695–710, 2016.
- [19] Atsushi Nitanda. Stochastic proximal gradient descent with acceleration techniques. In *NIPS*, 2014.
- [20] Neal Parikh, Stephen Boyd, et al. Proximal algorithms. *Foundations and Trends® in Optimization*, 1(3):127–239, 2014.
- [21] J Tang, X Tang, and J Yuan. Profit maximization for viral marketing in online social networks. In *ICNP*, 2016.
- [22] Jie Tang, Jimeng Sun, Chi Wang, and Zi Yang. Social influence analysis in large-scale networks. In *KDD*, 2009.
- [23] Youze Tang, Yanchen Shi, and Xiaokui Xiao. Influence maximization in near-linear time: a martingale approach. In *SIGMOD*, 2015.
- [24] Youze Tang, Xiaokui Xiao, and Yanchen Shi. Influence maximization: near-optimal time complexity meets practical efficiency. In *SIGMOD*, 2014.
- [25] Ramon van Handel. Probability in high dimension. Technical report, Princeton University, 2014.

- [26] Chi Wang, Wei Chen, and Yajun Wang. Scalable influence maximization for independent cascade model in large-scale social networks. *Data Mining and Knowledge Discovery*, 25(3):545–576, 2012.
- [27] Yu Yang, Xiangbo Mao, Jian Pei, and Xiaofei He. Continuous influence maximization: What discounts should we offer to social network users? In *SIGMOD*, pages 727–741, 2016.

Appendix

A Proof of Theorem 1

In this section, we present the detailed proof of Theorem 1. The proof follows the structure of the proof of IMM [23].

Our analysis is based on a version of Chernoff bound, which is shown as follow. For convenience, we will use the notation OPT_g to denote the maximum value of g and OPT_{g+s} to denote the maximum value of $g + s$ in the set \mathcal{P} , and we use \mathbf{x}_g^* to denote the solution when maximizing the function g in the set \mathcal{P} , i.e. $g(\mathbf{x}_g^*) = \text{OPT}_g$. We will also use \mathbf{x}_{g+s}^* to denote the point maximizing $g + s$ in the set \mathcal{P} , i.e. $g(\mathbf{x}_{g+s}^*) = \text{OPT}_{g+s}$.

Proposition 2 (Chernoff Bound [23]). *Let X_1, X_2, \dots, X_t be t independent random variable with support $[0, 1]$, and let $\mathbb{E}[X_i] = \mu$ for all $i \in [t]$. Let $Y = \sum_{i=1}^t X_i$, we have for any $\gamma > 0$,*

$$\Pr\{Y - t\mu \geq \gamma \cdot t\mu\} \leq \exp\left(-\frac{\gamma^2}{2 + \frac{2}{3}\gamma}t\mu\right).$$

For any $0 < \gamma < 1$, we have

$$\Pr\{Y - t\mu < -\gamma \cdot t\mu\} \leq \exp\left(-\frac{\gamma^2}{2}t\mu\right).$$

Recall that we want to optimize the function $(g + s)(\mathbf{x})$ in the set \mathcal{P} , where $h(\mathbf{x}) \geq 0$ for all $\mathbf{x} \in \mathcal{P}$, and $g(\mathbf{x})$ is the influence of the network with strategy \mathbf{x} and $g(\mathbf{x})$, and we use

$$\hat{g}_{\mathcal{R}}(\mathbf{x}) = \frac{n}{\theta} \sum_{R \in \mathcal{R}} \left(1 - \prod_{v \in R} (1 - h_v(\mathbf{x}))\right)$$

to approximate $g(\mathbf{x})$.

Then given the Chernoff bound, our proof comes as follow. We first fix the number of generated independent RR-sets $\theta = |\mathcal{R}|$ and we assume that we have an lower bound LB for the optimal value OPT_{g+s} . We first show that with the randomness of the generated RR-sets, with high probability, optimizing the function $(\hat{g}_{\mathcal{R}} + s)$ will lead to a guaranteed approximation of the function $(g + s)$. (Lemma 7,8,9) Then we show that with high probability, the Sampling Procedure(Algorithm 2) will return a lower bound LB for the optimal value OPT_{g+s} .

Lemma 7. *Given a constant $0 < \alpha' < 1$. For any $\varepsilon > 0$, any $0 < \alpha'\varepsilon_1 < \varepsilon/3$, and any $\delta_1, \delta_2 > 0$. If we have an lower bound LB for the optimal value OPT_{g+s} , let*

$$\theta^{(1)} = \frac{2n \cdot \ln\left(\frac{1}{\delta_1}\right)}{\varepsilon_1^2 \cdot LB}, \quad \theta^{(2)}(\mathcal{N}) = \frac{2\alpha' \cdot n \cdot \ln\left(\frac{\mathcal{N}}{\delta_2}\right)}{(\varepsilon/3 - \alpha'\varepsilon_1)^2 LB},$$

where \mathcal{N} is a variable. Recall that $\mathbf{x}_{g+s}^* = \text{argmax}_{\mathbf{x} \in \mathcal{P}}(g(\mathbf{x}) + s(\mathbf{x}))$, then for any fixed $|\mathcal{R}| = \theta \geq \theta^{(1)}$, we have

$$\Pr\{\hat{g}_{\mathcal{R}}(\mathbf{x}_{g+s}^*) + h(\mathbf{x}_{g+s}^*) < (1 - \varepsilon_1) \cdot \text{OPT}_{g+s}\} \leq \delta_1.$$

For any fixed $|\mathcal{R}| = \theta \geq \theta^{(2)}(\mathcal{N})$ and any fixed possible solution $\mathbf{x} \in \mathcal{P}$ that satisfies $g(\mathbf{x}) + h(\mathbf{x}) < (\alpha' - \varepsilon/3)\text{OPT}_{g+s} - T$ where $T \geq 0$ is a constant, we have

$$\Pr\{\hat{g}_{\mathcal{R}}(\mathbf{x}) + s(\mathbf{x}) \geq \alpha'(1 - \varepsilon_1) \cdot \text{OPT}_{g+s} - T\} \leq \frac{\delta_2}{\mathcal{N}}.$$

Proof. First recall that

$$\hat{g}_{\mathcal{R}}(\mathbf{x}) = \frac{n}{\theta} \sum_{R \in \mathcal{R}} \left(1 - \prod_{v \in R} (1 - h_v(\mathbf{x})) \right).$$

Let $X_i^{\mathcal{R}}(\mathbf{x}) = 1 - \prod_{v \in R_i} (1 - h_v(\mathbf{x}))$ where $\mathcal{R} = \{R_1, R_2, \dots, R_\theta\}$, then $X_i^{\mathcal{R}}(\mathbf{x}) \in [0, 1]$ and $\hat{g}_{\mathcal{R}}(\mathbf{x}) = \frac{n}{\theta} \sum_{i=1}^{\theta} X_i^{\mathcal{R}}(\mathbf{x})$. We also know that $X_i^{\mathcal{R}}(\mathbf{x})$ are independent. Then from the Chernoff bound (Proposition 2), we have

$$\begin{aligned} & \Pr\{\hat{g}_{\mathcal{R}}(\mathbf{x}_{g+s}^*) + s(\mathbf{x}_{g+s}^*) < (1 - \varepsilon_1) \cdot \text{OPT}_{g+s}\} \\ &= \Pr\left\{\frac{n}{\theta} \sum_{i=1}^{\theta} X_i^{\mathcal{R}}(\mathbf{x}^*) + s(\mathbf{x}_{g+s}^*) < (1 - \varepsilon_1) \cdot (g(\mathbf{x}_{g+s}^*) + s(\mathbf{x}_{g+s}^*))\right\} \\ &= \Pr\left\{\sum_{i=1}^{\theta} X_i^{\mathcal{R}}(\mathbf{x}_{g+s}^*) - \frac{\theta}{n} g(\mathbf{x}_{g+s}^*) < -\varepsilon_1 \cdot \frac{\theta}{n} \text{OPT}_{g+s}\right\} \\ &\leq \exp\left(-\frac{\left(\varepsilon_1 \cdot \frac{\text{OPT}_{g+s}}{g(\mathbf{x}_{g+s}^*)}\right)^2 \theta}{2} \frac{\theta}{n} g(\mathbf{x}_{g+s}^*)\right) \\ &\leq \exp\left(-\frac{\varepsilon_1^2 \theta}{2} \frac{\text{OPT}_{g+s}}{n}\right) \\ &\leq \exp\left(-\frac{\varepsilon_1^2}{2} \frac{2n \cdot \ln\left(\frac{1}{\delta_1}\right)}{n \text{LB} \cdot \varepsilon_1^2} \text{OPT}_{g+s}\right) \leq \delta_1. \end{aligned}$$

Let $\varepsilon_2 = \varepsilon - \alpha' \varepsilon_1$. Let $\mathbf{x} \in \mathcal{P}$ be a point such that $g(\mathbf{x}) + s(\mathbf{x}) < (\alpha' - \varepsilon/3) \text{OPT}_{g+s} - T$, and we have

$$\begin{aligned} & \Pr\{\hat{g}_{\mathcal{R}}(\mathbf{x}) + s(\mathbf{x}) \geq \alpha'(1 - \varepsilon_1) \cdot \text{OPT}_{g+s} - T\} \\ &= \Pr\left\{\frac{n}{\theta} \sum_{i=1}^{\theta} X_i^{\mathcal{R}}(\mathbf{x}) + s(\mathbf{x}) - g(\mathbf{x}) - s(\mathbf{x}) \geq \alpha'(1 - \varepsilon_1) \cdot \text{OPT}_{g+s} - g(\mathbf{x}) - s(\mathbf{x}) - T\right\} \\ &\leq \Pr\left\{\sum_{i=1}^{\theta} X_i^{\mathcal{R}}(\mathbf{x}) - \frac{\theta}{n} g(\mathbf{x}) \geq \frac{\theta}{n} (\alpha'(1 - \varepsilon_1) \cdot \text{OPT}_{g+s} - (\alpha' - \varepsilon/3) \text{OPT}_{g+s})\right\} \\ &= \Pr\left\{\sum_{i=1}^{\theta} X_i^{\mathcal{R}}(\mathbf{x}) - \frac{\theta}{n} g(\mathbf{x}) \geq \frac{\theta}{n} (\varepsilon_2 \text{OPT}_{g+s})\right\} \\ &= \Pr\left\{\sum_{i=1}^{\theta} X_i^{\mathcal{R}}(\mathbf{x}) - \frac{\theta}{n} g(\mathbf{x}) \geq \left(\frac{\varepsilon_2 \text{OPT}_{g+s}}{g(\mathbf{x})}\right) \frac{\theta}{n} g(\mathbf{x})\right\} \\ &\leq \exp\left(-\frac{\left(\frac{\varepsilon_2 \text{OPT}_{g+s}}{g(\mathbf{x})}\right)^2 \theta}{2 + \frac{2}{3} \frac{\varepsilon_2 \text{OPT}_{g+s}}{g(\mathbf{x})}} \frac{\theta}{n} g(\mathbf{x})\right) \\ &\leq \exp\left(-\frac{\varepsilon_2^2 \text{OPT}_{g+s}^2}{2g(\mathbf{x}) + \frac{2}{3} \varepsilon_2 \text{OPT}_{g+s}} \frac{\theta}{n}\right) \\ &\leq \exp\left(-\frac{\varepsilon_2^2 \text{OPT}_{g+s}^2}{2(\alpha' - \varepsilon/3) \text{OPT}_{g+s} + \frac{2}{3} \varepsilon_2 \text{OPT}_{g+s}} \frac{\theta}{n}\right) \\ &\leq \exp\left(-\frac{(\varepsilon/3 - \alpha' \varepsilon_1)^2 \text{OPT}_{g+s}}{2\alpha'} \frac{\theta}{n}\right) \end{aligned}$$

$$\leq \exp \left(-\frac{(\varepsilon/3 - \alpha'\varepsilon_1)^2 \text{OPT}_{g+s}}{2\alpha'} \frac{1}{n} \cdot \frac{2\alpha' \cdot n \cdot \ln \left(\frac{\mathcal{N}}{\delta_2} \right)}{(\varepsilon/3 - \alpha'\varepsilon_1)^2 \text{LB}} \right) \leq \frac{\delta_2}{\mathcal{N}}.$$

□

Lemma 8. *Given a constant $0 < \alpha < 1, L_1, L_2$. For any $\varepsilon > 0$, any $\varepsilon_2 > 0, 0 < \varepsilon_3 < \alpha$, any $0 < \varepsilon_1 < (\alpha - \varepsilon_3)\varepsilon/3$, and any $\delta_1, \delta_2 > 0$. Suppose that we have an lower bound LB for the optimal value OPT_{g+s} , if*

(a) $\Pr\{\hat{g}_{\mathcal{R}}(\mathbf{x}_{g+s}^*) + s(\mathbf{x}_{g+s}^*) < (1 - \varepsilon_1) \cdot \text{OPT}_{g+s}\} \leq \delta_1$;

(b) for any $\mathbf{x} \in \mathcal{P}$ that satisfies $g(\mathbf{x}) + s(\mathbf{x}) < (\alpha - \varepsilon_3 - \varepsilon/3)\text{OPT}_{g+s} - L_2\varepsilon_2\text{LB}$, we have

$$\Pr\{\hat{g}_{\mathcal{R}}(\mathbf{x}) + s(\mathbf{x}) \geq ((\alpha - \varepsilon_3)(1 - \varepsilon_1) \cdot \text{OPT}_{g+s} - L_2\varepsilon_2\text{LB})\} \leq \frac{\delta_2}{\mathcal{N}(\mathcal{P}, \varepsilon_2\text{LB})};$$

(c) the algorithm \mathcal{A} will output \mathbf{x}_{out} such that $(\hat{g}_{\mathcal{R}} + s)(\mathbf{x}_{out}) \geq (\alpha - \varepsilon_3) \cdot \max_{\mathbf{x} \in \mathcal{P}}(\hat{g}_{\mathcal{R}} + s)(\mathbf{x})$;

(d) $(g + s)(\mathbf{x})$ is L_1 -Lipschitz, and $(\hat{g}_{\mathcal{R}} + s)(\mathbf{x})$ is L_2 -Lipschitz;

then with probability at least $1 - \delta_1 - \delta_2$,

$$(g + s)(\mathbf{x}_{out}) \geq (\alpha - \varepsilon/3 - \varepsilon_3 - (L_1 + L_2)\varepsilon_2)\text{OPT}_{g+s}.$$

Proof. First we fix an $\varepsilon_2\text{LB}$ -net E for the set \mathcal{P} with number of points $\mathcal{N}(\mathcal{P}, \varepsilon_2\text{LB})$. Let $\pi(\mathbf{x}) : \mathcal{P} \rightarrow E$ denote the mapping from \mathcal{P} to the $\varepsilon_2\text{LB}$ -net E such that $\|\pi(\mathbf{x}) - \mathbf{x}\|_2 \leq \varepsilon_2\text{LB}$. From assumption (a), we know that with probability at least $1 - \delta_1$, we have

$$\begin{aligned} (\hat{g}_{\mathcal{R}} + s)(\mathbf{x}_{out}) &\geq (\alpha - \varepsilon_3) \max_{\mathbf{x} \in \mathcal{P}}(\hat{g}_{\mathcal{R}} + s)(\mathbf{x}) \\ &\geq (\alpha - \varepsilon_3)(\hat{g}_{\mathcal{R}} + s)(\mathbf{x}_{g+s}^*) \\ &\geq (\alpha - \varepsilon_3)(1 - \varepsilon_1) \cdot \text{OPT}_{g+s}. \end{aligned}$$

Since $(\hat{g}_{\mathcal{R}} + s)(\mathbf{x})$ is L_2 -Lipschitz, we have

$$\begin{aligned} (\hat{g}_{\mathcal{R}} + s)(\pi(\mathbf{x}_{out})) &= (\hat{g}_{\mathcal{R}} + s)(\mathbf{x}_{out}) + ((\hat{g}_{\mathcal{R}} + s)(\pi(\mathbf{x}_{out})) - (\hat{g}_{\mathcal{R}} + s)(\mathbf{x}_{out})) \\ &\geq (\hat{g}_{\mathcal{R}} + s)(\mathbf{x}_{out}) - |(\hat{g}_{\mathcal{R}} + s)(\pi(\mathbf{x}_{out})) - (\hat{g}_{\mathcal{R}} + s)(\mathbf{x}_{out})| \\ &\geq (\alpha - \varepsilon_3)(1 - \varepsilon_1) \cdot \text{OPT}_{g+s} - L_2\varepsilon_2\text{LB}. \end{aligned}$$

Then from (b) and the union bound, we know that with probability at least $1 - \delta_2$, for every $\mathbf{x} \in E$, if $g(\mathbf{x}) + s(\mathbf{x}) < (\alpha - \varepsilon_3 - \varepsilon/3)\text{OPT}_{g+s} - L_2\varepsilon_2\text{LB}$, then

$$\hat{g}_{\mathcal{R}}(\mathbf{x}) + s(\mathbf{x}) < (\alpha - \varepsilon_3)(1 - \varepsilon_1) \cdot \text{OPT}_{g+s} - L_2\varepsilon_2\text{LB}.$$

Then by the union bound, we know that with probability at least $1 - \delta_1 - \delta_2$,

$$(g + s)(\pi(\mathbf{x}_{out})) \geq (\alpha - \varepsilon_3 - \varepsilon/3)\text{OPT}_{g+s} - L_2\varepsilon_2\text{LB}.$$

Since from (d), $(g + s)(\mathbf{x})$ is L_1 -Lipschitz, then with probability at least $1 - \delta_1 - \delta_2$, we have

$$(g + s)(\mathbf{x}_{out}) \geq (g + s)(\pi(\mathbf{x}_{out})) - L_1\varepsilon_2\text{LB} \geq (\alpha - \varepsilon_3 - \varepsilon/3 - (L_1 + L_2)\varepsilon_2)\text{OPT}_{g+s}.$$

□

Combining Lemma 7 and Lemma 8 together, we have the following lemma,

Lemma 9. If $(g + s)(\mathbf{x})$ is L_1 -Lipschitz, and $(\hat{g}_{\mathcal{R}} + s)(\mathbf{x})$ is L_2 -Lipschitz. Suppose that we have an lower bound LB for the optimal value OPT_{g+s} , and an oracle \mathcal{A} that can get an $(\alpha - \varepsilon/3)$ -approximation for OPT_{g+s} , where $\alpha - \varepsilon/3 > 0$. Let

$$\theta^{(1)} = \frac{8n \cdot \ln(4n^\ell)}{LB \cdot (\alpha - \varepsilon/3)^2 \varepsilon^2 / 9}, \quad \theta^{(2)} = \frac{2\alpha' \cdot n \cdot \ln\left(4n^\ell \mathcal{N}\left(\mathcal{P}, \frac{\varepsilon/3}{L_1+L_2} LB\right)\right)}{(\varepsilon/3 - \frac{1}{4}(\alpha - \varepsilon/3)^2 \varepsilon/3)^2 LB}.$$

If $|\mathcal{R}| = \theta \geq \max\{\theta^{(1)}, \theta^{(2)}\}$, and $\mathbf{x}_{out} = \mathcal{A}(\hat{g}_{\mathcal{R}} + h, \varepsilon LB)$, then with probability at least $1 - \frac{1}{2n^\ell}$,

$$(g + s)(\mathbf{x}_{out}) \geq (\alpha - \varepsilon) OPT_{g+s}.$$

Proof. The proof of this lemma is a direct combination of Lemma 7 and Lemma 8. We choose the parameters $\delta_1 = \delta_2 = \frac{1}{4n^\ell}$ in Lemma 7 and Lemma 8. We choose $\varepsilon_1 = \frac{1}{2}(\alpha - \varepsilon/3)$, $\varepsilon_2 = \frac{\varepsilon/3}{L_1+L_2}$, $\varepsilon_3 = \varepsilon/3$ in Lemma 8, and $\alpha' = \alpha - \varepsilon/3$ in Lemma 7.

Now since $|\mathcal{R}| = \theta \geq \max\{\theta^{(1)}, \theta^{(2)}\}$, then the assumption of Lemma 7 is satisfied, and then the assumption (a),(b) of Lemma 8 is satisfied.

Then based on our assumption on the oracle \mathcal{A} , we know that

$$(\hat{g}_{\mathcal{R}} + s)(\mathbf{x}_{out}) \geq \alpha \max_{\mathbf{x} \in \mathcal{P}} (\hat{g}_{\mathcal{R}} + s)(\mathbf{x}) - \varepsilon/3 \cdot LB \geq (\alpha - \varepsilon/3) \max_{\mathbf{x} \in \mathcal{P}} (\hat{g}_{\mathcal{R}} + s)(\mathbf{x}),$$

then the assumption (c) of Lemma 8 is satisfied.

We also assume that $(g + s)(\mathbf{x})$ is L_1 -Lipschitz, and $(\hat{g}_{\mathcal{R}} + s)(\mathbf{x})$ is L_2 -Lipschitz, so assumption (d) is also satisfied. Then we know that with probability at least $1 - \frac{1}{2n^\ell}$,

$$(g + s)(\mathbf{x}_{out}) \geq (\alpha - \varepsilon) OPT_{g+s}.$$

□

Now we show that with high probability, the sampling procedure will return a lower bound $LB \leq OPT_{g+s}$.

Lemma 10. For every $i = 1, 2, \dots, \lfloor \log_2(n + \lambda k) - 1 \rfloor$, suppose that $g_{\mathcal{R}_i} + h$ is L_2 -Lipschitz where $|\mathcal{R}_i| =$

$$\theta_i = \left\lceil \frac{n \cdot (2 + \frac{2}{3}\varepsilon') \cdot (\ln \mathcal{N}(\mathcal{P}, \frac{\varepsilon/3}{L_2} x_i) + \ell \ln n + \ln 2 + \ln \log_2(n + \lambda k))}{\varepsilon'^2 x_i} \right\rceil,$$

(a) if $x_i = \frac{n+k}{2^i} > OPT_{g+s}$, then with probability at most $\frac{1}{2n^\ell \log_2(n + \lambda k)}$,

$$(\hat{g}_{\mathcal{R}_i} + s)(\mathbf{y}_i) \geq (1 + \varepsilon' + \varepsilon/3)x_i;$$

(b) if $x_i = \frac{n+k}{2^i} \leq OPT_{g+s}$, then with probability at most $\frac{1}{2n^\ell \log_2(n + \lambda k)}$,

$$(\hat{g}_{\mathcal{R}_i} + s)(\mathbf{y}_i) \geq (1 + \varepsilon' + \varepsilon/3) OPT_{g+s}.$$

Proof. Let E_i be the $\frac{\varepsilon/3}{L_2} x_i$ -net such that $|E_i| = \mathcal{N}(\mathcal{P}, \frac{\varepsilon/3}{L_2} x_i)$ and let $\pi_i(\mathbf{x}) : \mathcal{P} \rightarrow E_i$ denote the mapping such that $\|\pi_i(\mathbf{x}) - \mathbf{x}\|_2 \leq \frac{\varepsilon/3}{L_2} x_i$. Since $g_{\mathcal{R}_i} + h$ is L_2 -Lipschitz, we only have to prove that

(a') if $x_i = \frac{n+k}{2^i} > OPT_{g+s}$, then with probability at most $\frac{1}{2n^\ell \log_2(n + \lambda k)}$,

$$(\hat{g}_{\mathcal{R}_i} + s)(\pi_i(\mathbf{y}_i)) \geq (1 + \varepsilon')x_i;$$

(b') if $x_i = \frac{n+k}{2^i} \leq OPT_{g+s}$, then with probability at most $\frac{1}{2n^\ell \log_2(n + \lambda k)}$,

$$(\hat{g}_{\mathcal{R}_i} + s)(\pi_i(\mathbf{y}_i)) \geq (1 + \varepsilon') OPT_{g+s}.$$

Let $X_j^{\mathcal{R}_i}(\mathbf{x}) = 1 - \prod_{v \in \mathcal{R}_j} (1 - h_v(\mathbf{x}))$ where $\mathcal{R}_i = \{R_1, R_2, \dots, R_{\theta_i}\}$, then $X_j^{\mathcal{R}_i}(\mathbf{x}) \in [0, 1]$ and $\hat{g}_{\mathcal{R}_i}(\mathbf{x}) = \frac{n}{\theta} \sum_{j=1}^{\theta} X_h^{\mathcal{R}_i}(\mathbf{x})$.

We first prove (a'). We first fix a $\mathbf{y} \in \mathcal{P}$, then if $x_i = \frac{n+k}{2^i} > \text{OPT}_{g+s}$, we have

$$\begin{aligned}
& \Pr \{ (\hat{g}_{\mathcal{R}_i} + s)(\mathbf{y}) \geq (1 + \varepsilon')x_i \} \\
&= \Pr \left\{ \frac{n}{\theta} \sum_{j=1}^{\theta_i} X_j^{\mathcal{R}_i}(\mathbf{y}) + s(\mathbf{y}) \geq (1 + \varepsilon')x_i \right\} \\
&= \Pr \left\{ \frac{n}{\theta} \sum_{j=1}^{\theta_i} X_j^{\mathcal{R}_i}(\mathbf{y}) - g(\mathbf{y}) \geq (1 + \varepsilon')x_i - s(\mathbf{y}) - g(\mathbf{y}) \right\} \\
&\leq \Pr \left\{ \sum_{j=1}^{\theta_i} X_j^{\mathcal{R}_i}(\mathbf{y}) - \frac{\theta_i}{n} g(\mathbf{y}) \geq \frac{\theta_i}{n} \varepsilon' x_i \right\} \\
&\leq \exp \left(- \frac{\left(\frac{\varepsilon' x_i}{g(\mathbf{y})} \right)^2 \theta_i}{2 + \frac{2}{3} \frac{\varepsilon' x_i}{g(\mathbf{y})} n} g(\mathbf{y}) \right) \\
&= \exp \left(- \frac{(\varepsilon' x_i)^2 \theta_i}{2g(\mathbf{y}) + \frac{2}{3} \varepsilon' x_i n} \right) \\
&\leq \exp \left(- \frac{\varepsilon'^2 x_i \theta_i}{2 + \frac{2}{3} \varepsilon' n} \right) \\
&\leq \exp \left(- \frac{\varepsilon'^2 x_i \frac{1}{n} \cdot \left(2 + \frac{2}{3} \varepsilon' \right) \cdot \left(\ln \mathcal{N}(\mathcal{P}, \frac{\varepsilon/3}{L_2} x_i) + \ell \ln n + \ln 2 + \ln \log_2(n + \lambda k) \right)}{2 + \frac{2}{3} \varepsilon' n} \frac{1}{\varepsilon'^2 x_i} \right) \\
&= \frac{1}{2n^\ell \cdot \mathcal{N}(\mathcal{P}, \frac{\varepsilon/3}{L_2} x_i) \cdot \log_2(n + \lambda k)}.
\end{aligned}$$

Then note that there are at most $\mathcal{N}(\mathcal{P}, \frac{\varepsilon}{L_2} x_i)$ possibilities for $\pi_i(\mathbf{y}_i)$, so applying the union bound, we have

$$\Pr \{ (\hat{g}_{\mathcal{R}_i} + s)(\pi_i(\mathbf{y}_i)) \geq (1 + \varepsilon')x_i \} \leq \frac{1}{2n^\ell \log_2(n + \lambda k)}.$$

Then we prove (b'). If $x_i = \frac{n+k}{2^i} \leq \text{OPT}_{g+s}$, then for any fixed $\mathbf{y} \in \mathcal{P}$, we have

$$\begin{aligned}
& \Pr \{ (\hat{g}_{\mathcal{R}_i} + s)(\mathbf{y}) \geq (1 + \varepsilon')\text{OPT}_{g+s} \} \\
&= \Pr \left\{ \frac{n}{\theta} \sum_{j=1}^{\theta_i} X_j^{\mathcal{R}_i}(\mathbf{y}) + s(\mathbf{y}) \geq (1 + \varepsilon')\text{OPT}_{g+s} \right\} \\
&= \Pr \left\{ \frac{n}{\theta} \sum_{j=1}^{\theta_i} X_j^{\mathcal{R}_i}(\mathbf{y}) - g(\mathbf{y}) \geq (1 + \varepsilon')\text{OPT}_{g+s} - s(\mathbf{y}) - g(\mathbf{y}) \right\} \\
&\leq \Pr \left\{ \sum_{j=1}^{\theta_i} X_j^{\mathcal{R}_i}(\mathbf{y}) - \frac{\theta_i}{n} g(\mathbf{y}) \geq \frac{\theta_i}{n} \varepsilon' \text{OPT}_{g+s} \right\} \\
&\leq \exp \left(- \frac{\left(\frac{\varepsilon' \text{OPT}_{g+s}}{g(\mathbf{y})} \right)^2 \theta_i}{2 + \frac{2}{3} \frac{\varepsilon' \text{OPT}_{g+s}}{g(\mathbf{y})} n} g(\mathbf{y}) \right) \\
&= \exp \left(- \frac{(\varepsilon' \text{OPT}_{g+s})^2 \theta_i}{2g(\mathbf{y}) + \frac{2}{3} \varepsilon' \text{OPT}_{g+s} n} \right)
\end{aligned}$$

$$\begin{aligned}
&\leq \exp\left(-\frac{\varepsilon'^2 \text{OPT}_{g+s} \theta_i}{2 + \frac{2}{3}\varepsilon'} \frac{1}{n}\right) \\
&\leq \exp\left(-\frac{\varepsilon'^2 \text{OPT}_{g+s} \frac{1}{n} \cdot \left(2 + \frac{2}{3}\varepsilon'\right) \cdot \left(\ln \mathcal{N}(\mathcal{P}, \frac{\varepsilon/3}{L_2} x_i) + \ell \ln n + \ln 2 + \ln \log_2(n + \lambda k)\right)}{2 + \frac{2}{3}\varepsilon'} \frac{1}{n} \frac{1}{\varepsilon'^2 x_i}}\right) \\
&\leq \frac{1}{2n^\ell \cdot \mathcal{N}(\mathcal{P}, \frac{\varepsilon/3}{L_2} x_i) \cdot \log_2(n + \lambda k)}.
\end{aligned}$$

Then similar to the proof of (a'), applying the union bound will conclude the proof. \square

Then with the help of Lemma 10, we can prove Theorem 1.

Theorem 1. *For any $\varepsilon, \ell, \alpha > 0$, for any $(\alpha, \varepsilon/3)$ -approximate gradient algorithm \mathcal{A} for $\hat{g}_{\mathcal{R}} + s$, with probability at least $1 - \frac{1}{n^\ell}$, Grad-RIS outputs a solution \mathbf{x} that is an $(\alpha - \varepsilon)$ -approximation of the optimal solution OPT_{g+s} of CIM-BS, i.e. $(g + s)(\mathbf{x}) \geq (\alpha - \varepsilon)\text{OPT}_{g+s}$.*

Proof of Theorem 1. First we show that with probability at least $1 - \frac{1}{2n^\ell}$, the output lower bound $\text{LB} \leq \text{OPT}_{g+s}$. We first prove the case when $\text{OPT}_{g+s} \geq x_{\lfloor \log_2(n + \lambda k) \rfloor - 1}$. Let k denote the smallest index such that $\text{OPT}_{g+s} \geq x_k$. Then for any $i \leq k - 1$, we have $\text{OPT}_{g+s} < x_i$ and for any $j \geq k$, we have $\text{OPT}_{g+s} \geq x_j$. Then from Lemma 10 and union bound, we know that with probability at least $\frac{1}{2n^\ell}$, for every $i \leq k - 1$, $(\hat{g}_{\mathcal{R}_i} + s)(\pi_i(\mathbf{y}_i)) \geq (1 + \varepsilon')x_i$, and for every $j \geq k$, we have $(\hat{g}_{\mathcal{R}_j} + s)(\mathbf{y}_j) \geq (1 + \varepsilon' + \varepsilon/3)\text{OPT}_{g+s}$. Then from the definition of the algorithm, we know that $\text{LB} \leq \text{OPT}_{g+s}$.

Then for the case when $\text{OPT}_{g+s} < x_{\lfloor \log_2(n + \lambda k) \rfloor - 1}$, from the union bound, we know that with probability at least $1 - \frac{1}{2n^\ell}$ the 'break' statement will not be executed. So $\text{LB} = 1 \leq \text{OPT}_{g+s}$.

Then we bound the probability that the algorithm does not return an $\alpha - \varepsilon$ approximation. Let \mathcal{A} denote the event that the algorithm does not return an $\alpha - \varepsilon$ approximation, and \mathcal{B} denote the event that the output of the Sampling procedure $\text{LB} > \text{OPT}_{g+s}$. We want to show that $\Pr\{\mathcal{A}\} \leq \frac{1}{n^\ell}$. We have

$$\begin{aligned}
\Pr\{\mathcal{A}\} &= \Pr\{\mathcal{A} \wedge \mathcal{B}\} + \Pr\{\mathcal{A} \wedge \neg \mathcal{B}\} \\
&\leq \Pr\{\mathcal{B}\} + \Pr\{\mathcal{A} | \neg \mathcal{B}\}.
\end{aligned}$$

From Lemma 10, we have $\Pr\{\mathcal{B}\} \leq \frac{1}{2n^\ell}$. Since we generate new RR-sets before using the oracle to get the solution, so LB can be viewed as fixed, and from Lemma 9, we know that $\Pr\{\mathcal{A} | \neg \mathcal{B}\} \leq \frac{1}{2n^\ell}$. Combined them together, we complete the proof. \square

B Omitted Proofs in Subsection 3.3

B.1 Proof of Theorem 2

In this section, we give the formal proof of Theorem 2. First we slightly review the iteration step and the notations. For the proximal gradient descent, the iteration is shown as follows:

$$\begin{cases} \mathbf{x}^{(t+1)} = \text{prox}_{-\eta_t f_2}(\mathbf{x}^{(t)} + \eta_t \mathbf{v}^{(t)}), \text{ where } \mathbb{E}[\mathbf{v}^{(t)}] = \nabla f_1(\mathbf{x}^{(t)}), \\ \text{prox}_\phi(\mathbf{x}) := \text{argmin}_{\mathbf{y} \in \mathcal{P}} (\phi(\mathbf{y}) + \frac{1}{2} \|\mathbf{x} - \mathbf{y}\|_2^2), \text{ for any convex function } \phi, \end{cases} \quad (7)$$

where $\eta_t \leq \frac{1}{\beta}$ is the step size and $\mathbf{v}^{(t)}$ is the stochastic gradient at $\mathbf{x}^{(t)}$. Note that without loss of generality, we can assume that $f_2(\mathbf{x}) = -\infty$ for all $\mathbf{x} \notin \mathcal{P}$, and we have $\text{argmin}_{\mathbf{y} \in \mathcal{P}} (\phi(\mathbf{y}) + \frac{1}{2} \|\mathbf{x} - \mathbf{y}\|_2^2) = \text{argmin}_{\mathbf{y}} (\phi(\mathbf{y}) + \frac{1}{2} \|\mathbf{x} - \mathbf{y}\|_2^2)$. Let $G_\eta(\mathbf{x}) = \frac{1}{\eta}(\mathbf{x} - \text{prox}_{-\eta f_2}(\mathbf{x} + \eta \mathbf{v}))$ where \mathbf{v} is the stochastic gradient of f_1 at point \mathbf{x} , then we have $\mathbf{x}^{(t+1)} = \mathbf{x}^{(t)} - \eta_t G_{\eta_t}(\mathbf{x}^{(t)})$.

To analyze the convergence of the proximal gradient descent, we have the following proposition from [3].

Proposition 3. *If the function $f(\mathbf{x})$ is DR-submodular and monotone, we have*

$$\langle \nabla f(\mathbf{x}), \mathbf{x} - \mathbf{y} \rangle \leq 2f(\mathbf{x}) - f(\mathbf{x} \wedge \mathbf{y}) - f(\mathbf{x} \vee \mathbf{y}).$$

Lemma 11. *Let $\mathbf{u} = G_\eta(\mathbf{x}) + \mathbf{v}$, then $\mathbf{u} \in -\partial f_2(\mathbf{x} - \eta G_\eta(\mathbf{x}))$, i.e. \mathbf{u} is the subgradient of $-f_2$ at point $\mathbf{x} - \eta G_\eta(\mathbf{x})$. Here \mathbf{v} is the stochastic gradient of f_1 at point \mathbf{x} , and $G_\eta(\mathbf{x}) = \frac{1}{\eta}(\mathbf{x} - \text{prox}_{-\eta f_2}(\mathbf{x} + \eta \mathbf{v}))$.*

Proof. First it is easy to show that if g is convex and $\mathbf{u} = \text{prox}_g(\mathbf{x})$, then we have $\mathbf{x} - \mathbf{u} \in \partial g(\mathbf{u})$. This is due to the fact that \mathbf{u} minimize the function $g_1(\mathbf{y}) = g(\mathbf{y}) + \frac{1}{2}\|\mathbf{y} - \mathbf{x}\|_2^2$, and $g_1(\mathbf{y})$ is convex in \mathbf{y} . Then we have

$$0 \in \partial g(\mathbf{u}) + \mathbf{u} - \mathbf{x} \Rightarrow \mathbf{x} - \mathbf{u} \in \partial g(\mathbf{u}).$$

Then note that $\mathbf{x} - \eta G_\eta(\mathbf{x}) = \text{prox}_{-\eta f_2}(\mathbf{x} + \eta \mathbf{v})$, then we have

$$\mathbf{x} + \eta \mathbf{v} - \text{prox}_{-\eta f_2}(\mathbf{x} + \eta \nabla f(\mathbf{x})) \in -\eta \partial f_2(\mathbf{x} - \eta G_\eta(\mathbf{x})),$$

and rearranging the terms we have

$$\eta(\mathbf{v} + G_\eta(\mathbf{x})) \in -\eta \partial f_2(\mathbf{x} - \eta G_\eta(\mathbf{x})),$$

which concludes the proof. \square

Lemma 12. *Suppose $f_1(\mathbf{x})$ is a monotone DR-submodular function on convex set \mathcal{P} and $f_2(\mathbf{x})$ is concave on set \mathcal{P} . Note that we assume $f_2(\mathbf{x}) = -\infty$ for all $\mathbf{x} \notin \mathcal{P}$. If $f_1(\mathbf{x})$ is β -smooth, and $\eta \leq \frac{1}{\beta}$ is the step size, then for any $\mathbf{x}, \mathbf{z} \in \mathcal{P}$, we have*

$$\begin{aligned} (f_1 + f_2)(\mathbf{x} - \eta G_\eta(\mathbf{x})) + (f_1 + f_2)(\mathbf{x}) &\geq (f_1 + f_2)(\mathbf{z}) + \frac{\eta}{2} \|G_\eta(\mathbf{x})\|_2^2 - G_\eta(\mathbf{x})^T(\mathbf{x} - \mathbf{z}) \\ &\quad + f_2(\mathbf{x}) - (\mathbf{v} - \nabla f_1(\mathbf{x}))^T(\mathbf{x} - \eta G_\eta(\mathbf{x}) - \mathbf{z}). \end{aligned}$$

Proof. First note that $\mathbf{x} - \eta G_\eta(\mathbf{x}) \in \mathcal{P}$, since $\mathbf{x} - \eta G_\eta(\mathbf{x}) = \text{prox}_{-\eta f_2}(\mathbf{x}^{(k)} + \eta \mathbf{v}^{(k)})$ where $\mathbf{v}^{(k)}$ is the stochastic gradient of f_1 at point $\mathbf{x}^{(k)}$, and we know that $-\eta f_2(\mathbf{x}) = +\infty$ for all $\mathbf{x} \notin \mathcal{P}$.

From the smoothness of function f_1 and the convexity of $-f_2$ and the previous lemma(Lemma 11), we have

$$\begin{aligned} &-(f_1 + f_2)(\mathbf{x} - \eta G_\eta(\mathbf{x})) \\ &= -f_1(\mathbf{x} - \eta G_\eta(\mathbf{x})) - f_2(\mathbf{x} - \eta G_\eta(\mathbf{x})) \\ &\leq -f_1(\mathbf{x}) + \langle -\nabla f_1(\mathbf{x}), -\eta G_\eta(\mathbf{x}) \rangle + \frac{\beta}{2} \|\eta G_\eta(\mathbf{x})\|_2^2 - f_2(\mathbf{x} - \eta G_\eta(\mathbf{x})) \\ &\leq -f_1(\mathbf{x}) + \eta \nabla f_1(\mathbf{x})^T G_\eta(\mathbf{x}) + \frac{\beta}{2} \|\eta G_\eta(\mathbf{x})\|_2^2 \\ &\quad - f_2(\mathbf{z}) + (\mathbf{v} + G_\eta(\mathbf{x}))^T(\mathbf{x} - \eta G_\eta(\mathbf{x}) - \mathbf{z}) \\ &= -f_1(\mathbf{x}) + \eta \nabla f_1(\mathbf{x})^T G_\eta(\mathbf{x}) + \frac{\beta}{2} \|\eta G_\eta(\mathbf{x})\|_2^2 \\ &\quad - f_2(\mathbf{z}) + (\nabla f_1(\mathbf{x}) + G_\eta(\mathbf{x}))^T(\mathbf{x} - \eta G_\eta(\mathbf{x}) - \mathbf{z}) + (\mathbf{v} - \nabla f_1(\mathbf{x}))^T(\mathbf{x} - \eta G_\eta(\mathbf{x}) - \mathbf{z}) \\ &= -f_1(\mathbf{x}) + \nabla f_1(\mathbf{x})^T(\mathbf{x} - \mathbf{z}) - \frac{\beta}{2} \|\eta G_\eta(\mathbf{x})\|_2^2 - f_2(\mathbf{z}) + G_\eta(\mathbf{x})^T(\mathbf{x} - \mathbf{z}) \\ &\quad + (\mathbf{v} - \nabla f_1(\mathbf{x}))^T(\mathbf{x} - \eta G_\eta(\mathbf{x}) - \mathbf{z}). \end{aligned}$$

Then from the proposition(Proposition 3)

$$\langle \nabla f_1(\mathbf{x}), \mathbf{x} - \mathbf{z} \rangle \leq 2f_1(\mathbf{x}) - f_1(\mathbf{x} \wedge \mathbf{z}) - f_1(\mathbf{x} \vee \mathbf{z}) \leq 2f_1(\mathbf{x}) - f_1(\mathbf{z}),$$

we have

$$\begin{aligned}
-(f_1 + f_2)(\mathbf{x} - \eta G_\eta(\mathbf{x})) &\leq -f_1(\mathbf{x}) + \nabla f_1(\mathbf{x})^T(\mathbf{x} - \mathbf{z}) - \frac{\beta}{2} \|\eta G_\eta(\mathbf{x})\|_2^2 - f_2(\mathbf{z}) \\
&\quad + G_\eta(\mathbf{x})^T(\mathbf{x} - \mathbf{z}) + (\mathbf{v} - \nabla f_1(\mathbf{x}))^T(\mathbf{x} - \eta G_\eta(\mathbf{x}) - \mathbf{z}) \\
&\leq -f_1(\mathbf{x}) + 2f_1(\mathbf{x}) - f_1(\mathbf{z}) - \frac{\eta}{2} \|G_\eta(\mathbf{x})\|_2^2 - f_2(\mathbf{z}) \\
&\quad + G_\eta(\mathbf{x})^T(\mathbf{x} - \mathbf{z}) + (\mathbf{v} - \nabla f_1(\mathbf{x}))^T(\mathbf{x} - \eta G_\eta(\mathbf{x}) - \mathbf{z}) \\
&= f_1(\mathbf{x}) - (f_1 + f_2)(\mathbf{z}) + G_\eta(\mathbf{x})^T(\mathbf{x} - \mathbf{z}) - \frac{\eta}{2} \|G_\eta(\mathbf{x})\|_2^2 \\
&\quad + (\mathbf{v} - \nabla f_1(\mathbf{x}))^T(\mathbf{x} - \eta G_\eta(\mathbf{x}) - \mathbf{z}).
\end{aligned}$$

Rearranging the terms, we have

$$\begin{aligned}
(f_1 + f_2)(\mathbf{x} - \eta G_\eta(\mathbf{x})) + (f_1 + f_2)(\mathbf{x}) &\geq (f_1 + f_2)(\mathbf{z}) + \frac{\eta}{2} \|G_\eta(\mathbf{x})\|_2^2 - G_\eta(\mathbf{x})^T(\mathbf{x} - \mathbf{z}) \\
&\quad + f_2(\mathbf{x}) - (\mathbf{v} - \nabla f_1(\mathbf{x}))^T(\mathbf{x} - \eta G_\eta(\mathbf{x}) - \mathbf{z}).
\end{aligned}$$

□

Theorem 2. Suppose that \mathcal{P} is a convex set, function $f_1(\mathbf{x})$ is β -smooth, non-negative, monotone, and DR-submodular on \mathcal{P} , $f_2(\mathbf{x})$ is non-negative and concave on \mathcal{P} . Let \mathbf{x}^* be the point that maximizes $f_1(\mathbf{x}) + f_2(\mathbf{x})$. Suppose that for some $\sigma > 0$, the stochastic gradient $\mathbf{v}^{(t)}$ satisfies $\mathbb{E}\|\mathbf{v}^{(t)} - \nabla f_1(\mathbf{x}^{(t)})\|_2^2 \leq \sigma^2$ for all t , then for all $T > 0$, if we set $\eta_t = \eta = 1/(\beta + \frac{\sigma}{\Delta}\sqrt{2T})$, and iterate as shown in (3) starting from $\mathbf{x}^{(0)} \in \mathcal{P}$, we have

$$\begin{aligned}
&\mathbb{E} \left[\max_{t=0,1,2,\dots,T} (f_1 + f_2)(\mathbf{x}^{(t)}) \right] \\
&\geq \frac{1}{2} (f_1 + f_2)(\mathbf{x}^*) - \frac{\beta\Delta^2}{4T} - \frac{\sigma\Delta}{\sqrt{2T}}.
\end{aligned}$$

Proof of Theorem 2. Note that we assume the function f_2 to be the concave extension of the original function d , i.e. $f_2(\mathbf{x}) = -\infty$ for all $\mathbf{x} \notin \mathcal{P}$. Then we know that

$$\mathbf{x}^{(t+1)} = \text{prox}_{-\eta f_2}(\mathbf{x}^{(t)} + \eta_t \mathbf{v}^{(t)}) \in \mathcal{P},$$

which means that $f_2(\mathbf{x}^{(t+1)}) \geq 0$, for all $t = 1, 2, \dots, T$. Then in the previous lemma, let $\mathbf{z} = \mathbf{x}^*$ where \mathbf{x}^* maximize $f_1 + f_2$ in the set \mathcal{P} . We first consider the case when we have exact gradient $\mathbf{v}^{(t)} = \nabla f_1(\mathbf{x}^{(t)})$ and we set $\eta_t = \eta = \frac{1}{\beta}$, and we have

$$\begin{aligned}
&(f_1 + f_2)(\mathbf{x}^{(t+1)}) + (f_1 + f_2)(\mathbf{x}^{(t)}) \\
&\geq (f_1 + f_2)(\mathbf{x}^*) + \frac{\eta}{2} \|G_\eta(\mathbf{x}^{(t)})\|_2^2 - G_\eta(\mathbf{x}^{(t)})^T(\mathbf{x}^{(t)} - \mathbf{x}^*) + f_2(\mathbf{x}^{(t)}) \\
&\geq (f_1 + f_2)(\mathbf{x}^*) + \frac{\eta}{2} \|G_\eta(\mathbf{x}^{(t)})\|_2^2 - G_\eta(\mathbf{x}^{(t)})^T(\mathbf{x}^{(t)} - \mathbf{x}^*) \\
&= (f_1 + f_2)(\mathbf{x}^*) + \frac{1}{2\eta} (\eta G_\eta(\mathbf{x}^{(t)}))^T (\eta G_\eta(\mathbf{x}^{(t)}) - 2\mathbf{x}^{(t)} + 2\mathbf{x}^*) \\
&= (f_1 + f_2)(\mathbf{x}^*) + \frac{1}{2\eta} \left(\|\eta G_\eta(\mathbf{x}^{(t)}) - \mathbf{x}^{(t)} + \mathbf{x}^*\|_2^2 - \|\mathbf{x}^{(t)} - \mathbf{x}^*\|_2^2 \right) \\
&= (f_1 + f_2)(\mathbf{x}^*) + \frac{1}{2\eta} \left(\|\mathbf{x}^{(t+1)} - \mathbf{x}^*\|_2^2 - \|\mathbf{x}^{(t)} - \mathbf{x}^*\|_2^2 \right).
\end{aligned}$$

Then we sum up the above inequalities and we get

$$\begin{aligned}
& \sum_{t=0}^{T-1} ((f_1 + f_2)(\mathbf{x}^{(t+1)}) + (f_1 + f_2)(\mathbf{x}^{(t)})) \\
& \geq \sum_{t=0}^{T-1} \left((f_1 + f_2)(\mathbf{x}^*) + \frac{1}{2\eta} \left(\|\mathbf{x}^{(t+1)} - \mathbf{x}^*\|_2^2 - \|\mathbf{x}^{(t)} - \mathbf{x}^*\|_2^2 \right) \right) \\
& = T(f_1 + f_2)(\mathbf{x}^*) + \frac{1}{2\eta} \|\mathbf{x}^{(T)} - \mathbf{x}^*\|_2^2 - \frac{1}{2\eta} \|\mathbf{x}^{(0)} - \mathbf{x}^*\|_2^2 \\
& \geq T(f_1 + f_2)(\mathbf{x}^*) - \frac{1}{2\eta} \|\mathbf{x}^{(0)} - \mathbf{x}^*\|_2^2 \\
& \geq T(f_1 + f_2)(\mathbf{x}^*) - \frac{1}{2\eta} \Delta^2.
\end{aligned}$$

Then, we complete the proof by the fact that

$$\max_{t=0,1,2,\dots,T} h(\mathbf{x}^{(t)}) \geq \frac{1}{2T} \sum_{t=0}^{T-1} (h(\mathbf{x}^{(t+1)}) + h(\mathbf{x}^{(t)})),$$

and we have

$$\max_{t=0,1,2,\dots,T} h(\mathbf{x}^{(t)}) \geq \frac{1}{2} (f_1 + f_2)(\mathbf{x}^*) - \frac{\beta \Delta^2}{4T}.$$

Then we prove the stochastic gradient case. We first have the following property [8]: For convex function ϕ and any \mathbf{x}, \mathbf{y} , we have

$$\|\text{prox}_\phi(\mathbf{x}) - \text{prox}_\phi(\mathbf{y})\|_2 \leq \|\mathbf{x} - \mathbf{y}\|_2.$$

From the previous lemma and the previous property, we take the expectation of $\mathbf{v}^{(t)}$ and we can get

$$\begin{aligned}
& \mathbb{E}_{\mathbf{v}^{(t)}} (f_1 + f_2)(\mathbf{x}^{(t+1)}) + (f_1 + f_2)(\mathbf{x}^{(t)}) \\
& \geq (f_1 + f_2)(\mathbf{x}^*) + \frac{\eta}{2} \mathbb{E}_{\mathbf{v}^{(t)}} \|G_\eta(\mathbf{x}^{(t)})\|_2^2 - \mathbb{E}_{\mathbf{v}^{(t)}} G_\eta(\mathbf{x}^{(t)})^T (\mathbf{x}^{(t)} - \mathbf{x}^*) + f_2(\mathbf{x}^{(t)}) \\
& \quad - \mathbb{E}_{\mathbf{v}^{(t)}} (\mathbf{v}^{(t)} - \nabla f_1(\mathbf{x}^{(t)}))^T (\mathbf{x}^{(t)} - \eta G_\eta(\mathbf{x}^{(t)}) - \mathbf{x}^*) \\
& = (f_1 + f_2)(\mathbf{x}^*) + \frac{\eta}{2} \mathbb{E}_{\mathbf{v}^{(t)}} \|G_\eta(\mathbf{x}^{(t)})\|_2^2 - \mathbb{E}_{\mathbf{v}^{(t)}} G_\eta(\mathbf{x}^{(t)})^T (\mathbf{x}^{(t)} - \mathbf{x}^*) + f_2(\mathbf{x}^{(t)}) \\
& \quad - \mathbb{E}_{\mathbf{v}^{(t)}} (\mathbf{v}^{(t)} - \nabla f_1(\mathbf{x}^{(t)}))^T (\text{prox}_{-\eta f_2}(\mathbf{x}^{(t)} + \eta \mathbf{v}^{(t)}) - \text{prox}_{-\eta f_2}(\mathbf{x}^{(t)} + \eta \nabla f_1(\mathbf{x}^{(t)}))) \\
& \geq (f_1 + f_2)(\mathbf{x}^*) + \frac{\eta}{2} \mathbb{E}_{\mathbf{v}^{(t)}} \|G_\eta(\mathbf{x}^{(t)})\|_2^2 - \mathbb{E}_{\mathbf{v}^{(t)}} G_\eta(\mathbf{x}^{(t)})^T (\mathbf{x}^{(t)} - \mathbf{x}^*) + f_2(\mathbf{x}^{(t)}) \\
& \quad - \eta \mathbb{E}_{\mathbf{v}^{(t)}} \|\mathbf{v}^{(t)} - \nabla f_1(\mathbf{x}^{(t)})\|_2^2 \\
& \geq (f_1 + f_2)(\mathbf{x}^*) + \frac{\eta}{2} \mathbb{E}_{\mathbf{v}^{(t)}} \|G_\eta(\mathbf{x}^{(t)})\|_2^2 - \mathbb{E}_{\mathbf{v}^{(t)}} G_\eta(\mathbf{x}^{(t)})^T (\mathbf{x}^{(t)} - \mathbf{x}^*) - \sigma^2 \eta \\
& = (f_1 + f_2)(\mathbf{x}^*) + \frac{1}{2\eta} \mathbb{E}_{\mathbf{v}^{(t)}} (\eta G_\eta(\mathbf{x}^{(t)}))^T (\eta G_\eta(\eta \mathbf{x}^{(t)}) - 2\mathbf{x}^{(t)} + 2\mathbf{x}^*) - \sigma^2 \eta \\
& = (f_1 + f_2)(\mathbf{x}^*) + \frac{1}{2\eta} \mathbb{E}_{\mathbf{v}^{(t)}} \left(\|\eta G_\eta(\mathbf{x}^{(t)}) - \mathbf{x}^{(t)} + \mathbf{x}^*\|_2^2 - \|\mathbf{x}^{(t)} - \mathbf{x}^*\|_2^2 \right) - \sigma^2 \eta \\
& = (f_1 + f_2)(\mathbf{x}^*) + \frac{1}{2\eta} \mathbb{E}_{\mathbf{v}^{(t)}} \left(\|\mathbf{x}^{(t+1)} - \mathbf{x}^*\|_2^2 - \|\mathbf{x}^{(t)} - \mathbf{x}^*\|_2^2 \right) - \sigma^2 \eta.
\end{aligned}$$

Then we take expectation through all the randomness and sum them up, we have

$$\sum_{t=0}^{T-1} (\mathbb{E}(f_1 + f_2)(\mathbf{x}^{(t+1)}) + \mathbb{E}(f_1 + f_2)(\mathbf{x}^{(t)}))$$

$$\begin{aligned}
&\geq \sum_{t=0}^{T-1} \mathbb{E} \left((f_1 + f_2)(\mathbf{x}^*) + \frac{1}{2\eta} \left(\|\mathbf{x}^{(t+1)} - \mathbf{x}^*\|_2^2 - \|\mathbf{x}^{(t)} - \mathbf{x}^*\|_2^2 \right) \right) - T\sigma^2\eta \\
&= T(f_1 + f_2)(\mathbf{x}^*) + \frac{1}{2\eta} \mathbb{E} \|\mathbf{x}^{(T)} - \mathbf{x}^*\|_2^2 - \frac{1}{2\eta} \|\mathbf{x}^{(0)} - \mathbf{x}^*\|_2^2 - T\sigma^2\eta \\
&\geq T(f_1 + f_2)(\mathbf{x}^*) - \frac{1}{2\eta} \|\mathbf{x}^{(0)} - \mathbf{x}^*\|_2^2 - T\sigma^2\eta \\
&\geq T(f_1 + f_2)(\mathbf{x}^*) - \frac{1}{2\eta} \Delta^2 - T\sigma^2\eta.
\end{aligned}$$

Then we plug in $\eta = 1/(\beta + \frac{\sigma}{\Delta}\sqrt{2T})$, we have

$$\begin{aligned}
&\mathbb{E} \max_{t=0,1,2,\dots,T} (f_1 + f_2)(\mathbf{x}^{(t)}) \\
&\geq \max_{t=0,1,2,\dots,T} \mathbb{E} (f_1 + f_2)(\mathbf{x}^{(t)}) \\
&\geq \frac{1}{2T} \sum_{t=0}^{T-1} (\mathbb{E} (f_1 + f_2)(\mathbf{x}^{(t+1)}) + \mathbb{E} (f_1 + f_2)(\mathbf{x}^{(t)})) \\
&\geq \frac{1}{2} (f_1 + f_2)(\mathbf{x}^*) - \frac{1}{2T} \left(\frac{1}{2\eta} \Delta^2 + T\sigma^2\eta \right) \\
&\geq \frac{1}{2} (f_1 + f_2)(\mathbf{x}^*) - \frac{\beta\Delta^2}{4T} - \frac{\sigma\Delta}{\sqrt{2T}}.
\end{aligned}$$

□

B.2 Proofs of Lemma 3

The following lemma is the more detailed version of Lemma 3.

Lemma 13 (Detailed version of Lemma 3). *If functions $h_v(\mathbf{x})$'s are L_h -Lipschitz, then function $\hat{g}_{\mathcal{R}}(\mathbf{x})$ is $(\nu^{(1)}(\mathcal{R})nL_h)$ -Lipschitz, and function $g(\mathbf{x})$ is (n^2L_h) -Lipschitz. If functions $h_v(\mathbf{x})$'s are β_h -smooth, then function $\hat{g}_{\mathcal{R}}(\mathbf{x})$ is $(\nu^{(1)}(\mathcal{R})n\beta_h + \nu^{(2)}(\mathcal{R})nL_h^2)$ -smooth. The gradient of function $\hat{g}_{\mathcal{R}}(\mathbf{x})$ is*

$$\nabla \hat{g}_{\mathcal{R}}(\mathbf{x}) = \frac{n}{\theta} \sum_{R \in \mathcal{R}} \sum_{v' \in R} \nabla h_{v'}(\mathbf{x}) \left(\prod_{v \in R, v \neq v'} (1 - h_v(\mathbf{x})) \right),$$

and can be computed in time $O(\sum_{R \in \mathcal{R}} |R|(1 + T_h))$ if we assume that the gradient of $h_v(\mathbf{x})$ can be generated in time $O(T_h)$.

Proof of Lemma 13. First we have the following formula for the gradient of $\hat{g}_{\mathcal{R}}(\mathbf{x})$.

$$\begin{aligned}
\nabla \hat{g}_{\mathcal{R}}(\mathbf{x}) &= \nabla \frac{n}{\theta} \sum_{R \in \mathcal{R}} \left(1 - \prod_{v \in R} (1 - h_v(\mathbf{x})) \right) \\
&= \frac{n}{\theta} \sum_{R \in \mathcal{R}} \nabla \left(1 - \prod_{v \in R} (1 - h_v(\mathbf{x})) \right) \\
&= \frac{n}{\theta} \sum_{R \in \mathcal{R}} \sum_{v' \in R} \nabla h_{v'}(\mathbf{x}) \left(\prod_{v \in R, v \neq v'} (1 - h_v(\mathbf{x})) \right).
\end{aligned}$$

Next we show that we can generate the exact gradient of $\hat{g}_{\mathcal{R}}(\mathbf{x})$ in $O(\sum_{R \in \mathcal{R}} |R|)$ time (assuming that generating the gradient of h_v needs $O(1)$ time). First without loss of generality, we can assume that $(1 - h_v(\mathbf{x})) \neq 0$.

Otherwise, if $1 - h_u(\mathbf{x}) = 0$ and $u \in R$, then $\prod_{v \in R, v \neq u} (1 - h_u(\mathbf{x})) = 0$, and the problem is simpler. Then we can compute $\prod_{v \in R} (1 - h_v(\mathbf{x}))$ in $O(|R|)$ time and then compute $\prod_{v \in R, v \neq v'} (1 - h_v(\mathbf{x}))$ in $O(1)$ time. Then computing $\sum_{v' \in R} \nabla h_{v'}(\mathbf{x}) \left(\prod_{v \in R, v \neq v'} (1 - h_v(\mathbf{x})) \right)$ needs another $O(|R|T_h)$ time, and the total time complexity to compute the gradient $\nabla \hat{g}_{\mathcal{R}}(\mathbf{x})$ is $O(\sum_{R \in \mathcal{R}} |R|(1 + T_h))$. Then we compute the gradient of $g(\mathbf{x})$.

$$\begin{aligned}
\nabla g(\mathbf{x}) &= \sum_{S \subseteq V} \sigma(S) \left(\sum_{u' \in S} \nabla_x h_{u'}(\mathbf{x}) \left(\prod_{u \in S, u \neq u'} h_u(\mathbf{x}) \right) \left(\prod_{v \notin S} (1 - h_v(\mathbf{x})) \right) \right. \\
&\quad \left. - \sum_{v' \notin S} \nabla_x h_{v'}(\mathbf{x}) \left(\prod_{u \in S} h_u(\mathbf{x}) \right) \left(\prod_{v \notin S, v \neq v'} (1 - h_v(\mathbf{x})) \right) \right) \\
&= \sum_{u' \in V} \left[\sum_{S: u' \in S} \sigma(S) \left(\prod_{u \in S, u \neq u'} h_u(\mathbf{x}) \right) \left(\prod_{v \notin S} (1 - h_v(\mathbf{x})) \right) \cdot \nabla h_{u'}(\mathbf{x}) \right. \\
&\quad \left. - \sum_{T: u' \notin T} \sigma(T) \left(\prod_{u \in T} h_u(\mathbf{x}) \right) \left(\prod_{v \notin T, v \neq u'} (1 - h_v(\mathbf{x})) \right) \cdot \nabla h_{u'}(\mathbf{x}) \right] \\
&= \sum_{u' \in V} f_{u'}(\mathbf{x}) \nabla h_{u'}(\mathbf{x}),
\end{aligned}$$

where $f_{u'}(\mathbf{x})$ is defined as

$$\begin{aligned}
f_{u'}(\mathbf{x}) &:= \sum_{S: u' \in S} \sigma(S) \left(\prod_{u \in S, u \neq u'} h_u(\mathbf{x}) \right) \left(\prod_{v \notin S} (1 - h_v(\mathbf{x})) \right) \\
&\quad - \sum_{T: u' \notin T} \sigma(T) \left(\prod_{u \in T} h_u(\mathbf{x}) \right) \left(\prod_{v \notin T, v \neq u'} (1 - h_v(\mathbf{x})) \right).
\end{aligned}$$

We know that

$$\begin{aligned}
&|f_{u'}(\mathbf{x})| \\
&= \left| \sum_{S: u' \in S} \sigma(S) \left(\prod_{u \in S, u \neq u'} h_u(\mathbf{x}) \right) \left(\prod_{v \notin S} (1 - h_v(\mathbf{x})) \right) - \sum_{T: u' \notin T} \sigma(T) \left(\prod_{u \in T} h_u(\mathbf{x}) \right) \left(\prod_{v \notin T, v \neq u'} (1 - h_v(\mathbf{x})) \right) \right| \\
&\leq \left| \sum_{S: u' \in S} \sigma(S) \left(\prod_{u \in S, u \neq u'} h_u(\mathbf{x}) \right) \left(\prod_{v \notin S} (1 - h_v(\mathbf{x})) \right) \right| + \left| \sum_{T: u' \notin T} \sigma(T) \left(\prod_{u \in T} h_u(\mathbf{x}) \right) \left(\prod_{v \notin T, v \neq u'} (1 - h_v(\mathbf{x})) \right) \right| \\
&= \sum_{S: u' \in S} \sigma(S) \left(\prod_{u \in S, u \neq u'} h_u(\mathbf{x}) \right) \left(\prod_{v \notin S} (1 - h_v(\mathbf{x})) \right) + \sum_{T: u' \notin T} \sigma(T) \left(\prod_{u \in T} h_u(\mathbf{x}) \right) \left(\prod_{v \notin T, v \neq u'} (1 - h_v(\mathbf{x})) \right) \\
&\leq n \sum_{S: u' \in S} \left(\prod_{u \in S, u \neq u'} h_u(\mathbf{x}) \right) \left(\prod_{v \notin S} (1 - h_v(\mathbf{x})) \right) + n \sum_{T: u' \notin T} \left(\prod_{u \in T} h_u(\mathbf{x}) \right) \left(\prod_{v \notin T, v \neq u'} (1 - h_v(\mathbf{x})) \right) \\
&= n h_{u'}(\mathbf{x}) + n(1 - h_{u'}(\mathbf{x})) \\
&= n.
\end{aligned}$$

Then we have

$$\|\nabla g(\mathbf{x})\|_2 = \left\| \sum_{u' \in V} f_{u'}(\mathbf{x}) \nabla h_{u'}(\mathbf{x}) \right\|_2$$

$$\begin{aligned}
&\leq \sum_{u' \in \mathcal{V}} |f_{u'}(\mathbf{x})| \cdot \|\nabla h_{u'}(\mathbf{x})\|_2 \\
&\leq \sum_{u' \in \mathcal{V}} nL_h \\
&\leq n^2 L_h.
\end{aligned}$$

Then we know that $g(\mathbf{x})$ is $2n^2 L_h$ -Lipschitz. We also have

$$\begin{aligned}
\|\nabla \hat{g}_{\mathcal{R}}(\mathbf{x})\|_2 &= \left\| \frac{n}{\theta} \sum_{R \in \mathcal{R}} \sum_{v' \in R} \nabla h_{v'}(\mathbf{x}) \left(\prod_{v \in R, v \neq v'} (1 - h_v(\mathbf{x})) \right) \right\|_2 \\
&\leq \frac{n}{\theta} \sum_{R \in \mathcal{R}} \sum_{v' \in R} \|\nabla h_{v'}(\mathbf{x})\|_2 \left(\prod_{v \in R, v \neq v'} (1 - h_v(\mathbf{x})) \right) \\
&\leq \nu^{(1)}(\mathcal{R}) n L_h.
\end{aligned}$$

So the function $\hat{g}_{\mathcal{R}}(\mathbf{x})$ is also $n^2 L_h$ -Lipschitz. Then we show the smoothness of the function $\hat{g}_{\mathcal{R}}(\mathbf{x})$. We have

$$\begin{aligned}
&\|\nabla \hat{g}_{\mathcal{R}}(\mathbf{x}) - \nabla \hat{g}_{\mathcal{R}}(\mathbf{y})\|_2 \\
&= \left\| \frac{n}{\theta} \sum_{R \in \mathcal{R}} \sum_{v' \in R} \nabla h_{v'}(\mathbf{x}) \left(\prod_{v \in R, v \neq v'} (1 - h_v(\mathbf{x})) \right) - \frac{n}{\theta} \sum_{R \in \mathcal{R}} \sum_{v' \in R} \nabla h_{v'}(\mathbf{y}) \left(\prod_{v \in R, v \neq v'} (1 - h_v(\mathbf{y})) \right) \right\|_2 \\
&\leq \frac{n}{\theta} \sum_{R \in \mathcal{R}} \sum_{v' \in R} \left\| \nabla h_{v'}(\mathbf{x}) \left(\prod_{v \in R, v \neq v'} (1 - h_v(\mathbf{x})) \right) - \nabla h_{v'}(\mathbf{y}) \left(\prod_{v \in R, v \neq v'} (1 - h_v(\mathbf{y})) \right) \right\|_2 \\
&= \frac{n}{\theta} \sum_{R \in \mathcal{R}} \sum_{v' \in R} \left\| \nabla h_{v'}(\mathbf{x}) \left(\prod_{v \in R, v \neq v'} (1 - h_v(\mathbf{x})) \right) - \nabla h_{v'}(\mathbf{y}) \left(\prod_{v \in R, v \neq v'} (1 - h_v(\mathbf{x})) \right) \right. \\
&\quad \left. + \nabla h_{v'}(\mathbf{y}) \left(\prod_{v \in R, v \neq v'} (1 - h_v(\mathbf{x})) \right) - \nabla h_{v'}(\mathbf{y}) \left(\prod_{v \in R, v \neq v'} (1 - h_v(\mathbf{y})) \right) \right\|_2 \\
&\leq \frac{n}{\theta} \sum_{R \in \mathcal{R}} \sum_{v' \in R} \underbrace{\left\| \nabla h_{v'}(\mathbf{x}) \left(\prod_{v \in R, v \neq v'} (1 - h_v(\mathbf{x})) \right) - \nabla h_{v'}(\mathbf{y}) \left(\prod_{v \in R, v \neq v'} (1 - h_v(\mathbf{x})) \right) \right\|_2}_{\mathbb{A}} \\
&\quad + \frac{n}{\theta} \sum_{R \in \mathcal{R}} \sum_{v' \in R} \underbrace{\left\| \nabla h_{v'}(\mathbf{y}) \left(\prod_{v \in R, v \neq v'} (1 - h_v(\mathbf{x})) \right) - \nabla h_{v'}(\mathbf{y}) \left(\prod_{v \in R, v \neq v'} (1 - h_v(\mathbf{y})) \right) \right\|_2}_{\mathbb{B}}.
\end{aligned}$$

For term \mathbb{A} , we have

$$\begin{aligned}
\mathbb{A} &= \frac{n}{\theta} \sum_{R \in \mathcal{R}} \sum_{v' \in R} \left| \prod_{v \in R, v \neq v'} (1 - h_v(\mathbf{x})) \right| \cdot \|\nabla h_{v'}(\mathbf{x}) - \nabla h_{v'}(\mathbf{y})\|_2 \\
&\leq \frac{n}{\theta} \sum_{R \in \mathcal{R}} \sum_{v' \in R} \cdot 1 \cdot \beta_h \|\mathbf{x} - \mathbf{y}\|_2 \\
&\leq \beta_h \nu^{(1)}(\mathcal{R}) n \|\mathbf{x} - \mathbf{y}\|_2,
\end{aligned}$$

where we use the assumption that $h_v(\mathbf{x})$ is β_h -smooth. As for the term \mathbb{B} , we first have

$$\begin{aligned}
& \left| \prod_{v \in R \setminus \{v'\} = \{v_1, \dots, v_{|R|-1}\}} (1 - h_v(\mathbf{x})) - \prod_{v \in R \setminus \{v'\}} (1 - h_v(\mathbf{y})) \right| \\
&= \left| \prod_{v \in R \setminus \{v'\}} (1 - h_v(\mathbf{x})) - \prod_{i=1}^{|R|-2} (1 - h_{v_i}(\mathbf{x})) \cdot (1 - h_{v_{|R|-1}}(\mathbf{y})) \right. \\
&\quad + \prod_{i=1}^{|R|-2} (1 - h_{v_i}(\mathbf{x})) \cdot (1 - h_{v_{|R|-1}}(\mathbf{y})) - \prod_{i=1}^{|R|-3} (1 - h_{v_i}(\mathbf{x})) \cdot \prod_{j=|R|-2}^{|R|-1} (1 - h_{v_j}(\mathbf{y})) \\
&\quad \left. + \dots + (1 - h_{v_1}(\mathbf{x})) \cdot \prod_{j=2}^{|R|-1} (1 - h_{v_j}(\mathbf{y})) - \prod_{v \in R \setminus \{v'\}} (1 - h_v(\mathbf{y})) \right| \\
&\leq \sum_{i=1}^{|R|-1} |1 - h_{v_i}(\mathbf{x}) - 1 + h_{v_i}(\mathbf{y})| \\
&\leq |R|L_h \|\mathbf{x} - \mathbf{y}\|_2,
\end{aligned}$$

where the last inequality comes from the L_h -lipschitz property of the function $h_v(\mathbf{x})$. Then we have

$$\begin{aligned}
\mathbb{B} &\leq \frac{n}{\theta} \sum_{R \in \mathcal{R}} \sum_{v' \in R} \left| \prod_{v \in R \setminus \{v'\} = \{v_1, \dots, v_{|R|-1}\}} (1 - h_v(\mathbf{x})) - \prod_{v \in R \setminus \{v'\}} (1 - h_v(\mathbf{y})) \right| \cdot \|\nabla h_{v'}(\mathbf{y})\|_2 \\
&\leq \frac{n}{\theta} \sum_{R \in \mathcal{R}} \sum_{v' \in R} |R|L_h \|\mathbf{x} - \mathbf{y}\|_2 \cdot \|\nabla h_{v'}(\mathbf{y})\|_2 \\
&\leq \nu^{(2)}(\mathcal{R})nL_h^2 \|\mathbf{x} - \mathbf{y}\|_2.
\end{aligned}$$

Then we know that the function is $(\nu^{(1)}(\mathcal{R})n\beta_h + \nu^{(2)}(\mathcal{R})nL_h^2)$ -smooth. \square

B.3 Proof of Lemma 4

Lemma 4. *When $c(\mathbf{x}) = \|\mathbf{x}\|_1$ and $\mathcal{D} = \mathbb{R}_+^d$, the proximal step can be done in time $O(d \log d)$. When $c(\mathbf{x}) = \|\mathbf{x}\|_2$ and $\mathcal{D} = \mathbb{R}_+^d$, the proximal step can be done in $O(d)$.*

Proof of Lemma 4. First, it is easy to know that if $c(\mathbf{x}) = \|\mathbf{x}\|_2$ and $\mathcal{D} = \mathbb{R}_+^d$, then the proximal step can be finished in time $O(d)$. In this case, the set \mathcal{P} is defined as $\mathcal{P} = \{\mathbf{x} \mid \|\mathbf{x}\|_2 \leq k, \mathbf{x} \succeq 0\}$, and proximal step is defined as

$$\text{prox}_{-\eta s}(\mathbf{x}) := \underset{\mathbf{y} \in \mathcal{P}}{\text{argmin}} -\eta s(\mathbf{y}) + \frac{1}{2} \|\mathbf{x} - \mathbf{y}\|^2 = \underset{\mathbf{y} \in \mathcal{P}}{\text{argmin}} \eta \lambda \|\mathbf{y}\|_2 + \frac{1}{2} \|\mathbf{x} - \mathbf{y}\|^2.$$

It is obvious that \mathbf{y} should lie in the line generated by 0 and \mathbf{x} , and we can solve for \mathbf{y} by using the basic technique for solving optimal value of a uni-variate quadratic function.

Then we show that how to do the proximal step when $c(\mathbf{x}) = \|\mathbf{x}\|_1$ and $\mathcal{D} = \mathbb{R}_+^d$. In this case, we know that $\mathcal{P} = \{\mathbf{x} \mid \|\mathbf{x}\|_1 \leq k, \mathbf{x} \succeq 0\}$, and we want to find $\mathbf{y} \in \mathcal{P}$ to minimize

$$\text{prox}_{-\eta s}(\mathbf{x}) := \underset{\mathbf{y} \in \mathcal{P}}{\text{argmin}} -\eta s(\mathbf{y}) + \frac{1}{2} \|\mathbf{x} - \mathbf{y}\|^2 = \underset{\mathbf{y} \in \mathcal{P}}{\text{argmin}} \eta \lambda \|\mathbf{y}\|_1 + \frac{1}{2} \|\mathbf{x} - \mathbf{y}\|^2.$$

We use $\eta' = \eta \cdot \lambda$ for convenience. The proximal step can also be written as minimizing

$$\min_{y_i \geq 0, \sum_i y_i \leq k} \eta' \sum_i y_i + \frac{1}{2} \sum_i (x_i - y_i)^2.$$

First, we have $y_i \leq \max\{x_i - \eta', 0\}$ for all i . If not, suppose $y_i > \max\{x_i - \eta', 0\}$, then we let $y'_i = \max\{x_i - \eta', 0\}$. The new solution has smaller value and is also feasible.

Then we also have the following property: if \mathbf{y} is optimal, then for any $i \neq j$, if $y_i, y_j \neq 0$, then we have $x_i - y_i = x_j - y_j$. Suppose that $x_i - y_i > x_j - y_j$, then let ε be sufficiently small. Let $y'_i = y_i + \varepsilon$ and $y'_j = y_j - \varepsilon$, the new solution also lies in set \mathcal{P} , but the function value of $\eta' \sum_i y_i + \frac{1}{2} \sum_i (x_i - y_i)^2$ become smaller.

We also have: suppose \mathbf{y} is optimal, then if $x_i \geq x_j$. If $y_i = 0$, then $y_j = 0$. Otherwise, suppose $y_i = 0$ but $y_j > 0$. We can pick sufficiently small ε and let $y'_i = y_i + \varepsilon$ and $y'_j = y_j - \varepsilon$. The function value will decrease but the new solution is also feasible.

We also have: suppose \mathbf{y} is optimal, then if $y_i = 0$, then for any j such that $y_j \neq 0$, $x_j - y_j \geq x_i - y_i$. Otherwise, we can find sufficiently small ε and let $y'_i = y_i + \varepsilon$, $y'_j = y_j - \varepsilon$. The function value will decrease and the new solution is feasible.

Given the previous properties, we have the following structure of optimal value \mathbf{y} . Suppose that $\{(i)\}$ is a permutation of $[d]$ such that $x_{(1)} \leq x_{(2)} \leq \dots \leq x_{(d)}$, then there exists $i_0 \in [d]$ such that $y_{(i)} = 0$ for all $i \leq i_0$, and $0 < y_{(i)} \leq x_{(i)} - \lambda'$ for all $i > i_0$. Besides, for all $i, j > i_0$, we have $x_{(i)} - y_{(i)} = x_{(j)} - y_{(j)} > x_{(i_0)} - y_{(i_0)}$. We also have one of the following: For all i , $y_i = \max\{x_i - \eta', 0\}$; otherwise, $\sum_i y_i = k$. Then, there is only one i_0 that satisfy this the previous structure, and we show that we can use $O(d \log d)$ time to find i_0 . We first use $O(d \log d)$ time to sort x_i and get $x_{(i)}$. We first let $y_i = \max\{x_i - \eta', 0\}$ and we test if $\sum_i y_i \leq k$. If yes, then \mathbf{y} is the optimal solution. Otherwise, we know that $\sum_i y_i = k$. We binary search for (i_0) , and we use i_1 to denote the binary search variable. Each time we set $y_{(j)} = 0$ for all $j \leq i_1$, and we set $y_{(j)}$ such that $x_{(j)} - y_{(j)}$ are the same for all $j > i_1$ and $\sum_i y_i = k$. Then we test if $x_{(d)} - y_{(d)} \geq x_{(i_1)} - y_{(i_1)}$, and $y_{(j)} \geq 0$ for all $j > i_1$. If both are yes, then i_0 is i_1 and \mathbf{y} is the optimal solution. If there exists $j > i_1$ such that $y_{(j)} < 0$, then i_1 should be larger. If $x_{(d)} - y_{(d)} < x_{(i_1)} - y_{(i_1)}$, then i_1 should be smaller. Each test needs time $O(d)$, and there are at most $O(\log d)$ binary search step, so the total complexity is $O(d \log d)$. \square

B.4 Proof of Lemma 5

Lemma 5. *In the case of independent strategy activation, suppose that $\bar{g}_{\mathcal{R}} + s$ is $L_{\bar{g}+s}$ -Lipschitz. If we use projected subgradient descent to optimize the function $(\bar{g}_{\mathcal{R}} + s)(\mathbf{x})$ with step size $\eta_t = \frac{\Delta}{L_{\bar{g}+s}\sqrt{t}}$ and let \mathbf{y} denote the output where $(\bar{g}_{\mathcal{R}} + s)(\mathbf{y}) \geq \max_{\mathbf{x} \in \mathcal{P}}(\bar{g}_{\mathcal{R}} + s)(\mathbf{x}) - \varepsilon$. Then $(\hat{g}_{\mathcal{R}} + s)(\mathbf{y}) \geq (1 - \frac{1}{e}) \max_{\mathbf{x} \in \mathcal{P}}(\hat{g}_{\mathcal{R}} + s)(\mathbf{x}) - \varepsilon$, and \mathbf{y} can be solved in $\frac{(\Delta L_{\bar{g}+s})^2}{\varepsilon^2}$ iterations.*

Proof of Lemma 5. First, we optimize the function $(\bar{g}_{\mathcal{R}} + s)(\mathbf{x})$ in the set \mathcal{P} . First, it is easy to know that the function $(\bar{g}_{\mathcal{R}} + s)(\mathbf{x})$ is concave, since first we know that $s(\mathbf{x})$ is concave, $q_{v,j}(\mathbf{x})$ are concave for all v, j , the constant 1 is also concave. Then because addition of 2 concave function is also concave, and the point-wise minimum of 2 concave function is also concave, so the function $(\bar{g}_{\mathcal{R}} + s)(\mathbf{x})$ is concave. Then since we assume that the function $(\bar{g}_{\mathcal{R}} + s)(\mathbf{x})$ is $L_{\bar{g}+s}$ -Lipschitz, by the projected subgradient descent, in $O\left(\frac{L_{\bar{g}+s}}{\varepsilon^2}\right)$ iteration, we can get a solution \mathbf{y} such that $(\bar{g}_{\mathcal{R}} + s)(\mathbf{y}) \geq \max_{\mathbf{x} \in \mathcal{P}}(\bar{g}_{\mathcal{R}} + s)(\mathbf{x}) - \varepsilon$. Then we show that

$$(\hat{g}_{\mathcal{R}} + s)(\mathbf{y}) \geq \left(1 - \frac{1}{e}\right) \max_{\mathbf{x} \in \mathcal{P}}(\hat{g}_{\mathcal{R}} + s)(\mathbf{x}) - \varepsilon.$$

From Proposition 1, we can know that

$$\left(1 - \frac{1}{e}\right) \bar{g}_{\mathcal{R}}(\mathbf{x}) \leq \hat{g}_{\mathcal{R}}(\mathbf{x}) \leq \bar{g}_{\mathcal{R}}(\mathbf{x}), \forall \mathbf{x} \in \mathcal{P},$$

and because $s(\mathbf{x})$ is non-negative on set \mathcal{P} , we have

$$\left(1 - \frac{1}{e}\right) (\bar{g}_{\mathcal{R}} + s)(\mathbf{y}) \leq (\hat{g}_{\mathcal{R}} + s)(\mathbf{y}) \leq (\bar{g}_{\mathcal{R}} + s)(\mathbf{y}).$$

Then we have

$$\begin{aligned}
(\hat{g}_{\mathcal{R}} + s)(\mathbf{y}) &\geq \left(1 - \frac{1}{e}\right) (\bar{g}_{\mathcal{R}} + s)(\mathbf{y}) \\
&\geq \left(1 - \frac{1}{e}\right) \left(\max_{\mathbf{x} \in \mathcal{P}} (\bar{g}_{\mathcal{R}} + s)(\mathbf{x}) - \varepsilon\right) \\
&\geq \left(1 - \frac{1}{e}\right) \max_{\mathbf{x} \in \mathcal{P}} (\hat{g}_{\mathcal{R}} + s)(\mathbf{x}) - \varepsilon.
\end{aligned}$$

□

B.5 Proof of Lemma 6

The following lemma is a more detailed version of Lemma 6.

Lemma 14 (Detailed version of Lemma 6). *Suppose that functions $q_{v,j}(x_j)$'s are L_q -Lipschitz, then function $g(\mathbf{x})$ is $n^2\sqrt{d}L_q$ -Lipschitz, and functions $\hat{g}_{\mathcal{R}}(\mathbf{x})$ and $\bar{g}_{\mathcal{R}}(\mathbf{x})$ are $\nu^{(1)}(\mathcal{R})n\sqrt{d}L_q$ -Lipschitz. The subgradient of the function $\bar{g}_{\mathcal{R}}(\mathbf{x})$ is*

$$\begin{aligned}
\partial \bar{g}_{\mathcal{R}}(\mathbf{x}) &= \frac{n}{\theta} \sum_{R \in \mathcal{R}} \partial \min \left\{ 1, \sum_{j \in [d], v \in R} q_{v,j}(x_j) \right\} \\
&= \frac{n}{\theta} \sum_{R \in \mathcal{R}} \begin{cases} 0, & \text{if } \sum_{j \in [d], v \in R} q_{v,j}(x_j) \geq 1 \\ \sum_{v \in R, j \in [d]} \nabla q_{v,j}(x_j), & \text{if } \sum_{j \in [d], v \in R} q_{v,j}(x_j) < 1 \end{cases}
\end{aligned}$$

and can be computed in time $O(T_q)$ if we assume that the gradient and function value of $q_{v,j}(\mathbf{x})$ can be generated in time $O(\sum_{R \in \mathcal{R}} |R|(1 + T_q))$.

Proof of Lemma 14. First recall that $h_v(\mathbf{x}) = 1 - \prod_{j \in [d]} (1 - q_{v,j}(x_j))$. Then since $q_{v,j}(\mathbf{x})$ is L_q -Lipschitz, it can be easily shown that $h_v(\mathbf{x})$ is $L_q\sqrt{d}$ -Lipschitz, since

$$\begin{aligned}
|h_v(\mathbf{x}) - h_v(\mathbf{y})| &= \left| 1 - \prod_{j \in [d]} (1 - q_{v,j}(x_j)) - 1 + \prod_{j \in [d]} (1 - q_{v,j}(y_j)) \right| \\
&\leq \left| \prod_{j=1}^d (1 - q_{v,j}(x_j)) - (1 - q_{v,j}(y_1)) \prod_{j=2}^d (1 - q_{v,j}(x_j)) \right| \\
&\quad + \left| (1 - q_{v,j}(y_1)) \prod_{j=2}^d (1 - q_{v,j}(y_j)) - \prod_{j=1}^2 (1 - q_{v,j}(y_j)) \prod_{j=3}^d (1 - q_{v,j}(x_j)) \right| \\
&\quad + \dots + \left| \prod_{j=1}^{d-1} (1 - q_{v,j}(y_j)) \cdot (1 - q_{v,j}(x_d)) - \prod_{j=1}^d (1 - q_{v,j}(y_j)) \right| \\
&\leq \sum_{j=1}^d |q_{v,j}(x_j) - q_{v,j}(y_j)| \\
&\leq L_q \cdot \|\mathbf{x} - \mathbf{y}\|_1 \\
&\leq L_q \sqrt{d} \|\mathbf{x} - \mathbf{y}\|_2.
\end{aligned}$$

Then from the previous lemma (Lemma 13), we can see that $g(\mathbf{x})$ is $n^2\sqrt{d}L_q$ -Lipschitz and $\hat{g}_{\mathcal{R}}(\mathbf{x})$ is all $\nu^{(1)}(\mathcal{R})n\sqrt{d}L_q$ -Lipschitz. We also have

$$\begin{aligned}
|\bar{g}_{\mathcal{R}}(\mathbf{x}) - \bar{g}_{\mathcal{R}}(\mathbf{y})| &\leq \frac{n}{\theta} \sum_{R \in \mathcal{R}} \left| \min \left\{ 1, \sum_{j \in [d], v \in R} q_{v,j}(x_j) \right\} - \min \left\{ 1, \sum_{j \in [d], v \in R} q_{v,j}(y_j) \right\} \right| \\
&\leq \frac{n}{\theta} \sum_{R \in \mathcal{R}} \left| \sum_{j \in [d], v \in R} q_{v,j}(x_j) - \sum_{j \in [d], v \in R} q_{v,j}(y_j) \right| \\
&\leq \frac{n}{\theta} \sum_{R \in \mathcal{R}} \sum_{j \in [d], v \in R} |q_{v,j}(x_j) - q_{v,j}(y_j)| \\
&\leq \frac{n}{\theta} \sum_{R \in \mathcal{R}} \sum_{j \in [d], v \in R} L_q |x_j - y_j| \\
&\leq \frac{n}{\theta} \sum_{R \in \mathcal{R}} \sum_{v \in R} L_q \|\mathbf{x} - \mathbf{y}\|_1 \\
&\leq \nu^{(1)}(\mathcal{R}) n L_q \sqrt{d} \|\mathbf{x} - \mathbf{y}\|_2.
\end{aligned}$$

As for the subgradient of $\bar{g}_{\mathcal{R}}(\mathbf{x})$ and the time complexity to generate the subgradient, it is trivial. Note that in the time complexity $\sum_{R \in \mathcal{R}} |R|(1 + T_q)$, the constant 1 is used for basic operations. \square

C Omitted Proofs for Time Complexity (Theorem 5)

In this section, we present our time complexity results for ProxGrad-RIS and UpperGrad-RIS. We divide this section into 2 parts. In the first part, we show the time complexity of ProxGrad-RIS and UpperGrad-RIS in a general form (Theorem 6 and 7), and then Theorem 5 will become a corollary. However, the time complexity bound in Theorem 6 and 7 is too conservative and cannot reflect the empirical running time in experiments. In order to close this gap, we give time complexity bounds based on the moments of RR-set size in the second part. Then we will give some statistics of the moments of RR-set size in the section describing our experiments.

C.1 Proof of Theorem 5

In this subsection, we prove Theorem 5. We actually prove the full version of the running times for the two algorithms in Theorems 6 and 7. Our proof follows from the original proof of the time complexity of IMM algorithm. First, we have to show a lemma (Lemma 15), which states that in Algorithm 2, with high probability, we output LB does not differ so much from the optimal value OPT_{g+s} . For convenience, all the notations follow from Algorithm 2. We use $\text{OPT}_{g+s} = (g+s)(\mathbf{x}_{g+s}^*)$ to denote the maximum value of $(g+s)$ in set \mathcal{P} .

Lemma 15. *For every $i = 1, 2, \dots, \lfloor \log_2(n + \lambda k) \rfloor - 1$, if $\text{OPT}_{g+s} \geq (1 + \varepsilon/3 + \varepsilon')^2 \cdot x_i / (\alpha - \varepsilon/3)$, then with probability at least $1 - \frac{1}{2n^\ell \cdot \mathcal{N}(\mathcal{P}, \frac{\varepsilon/3}{L_2} x_i) \cdot \log_2(n + \lambda k)}$, $(\hat{g}_{\mathcal{R}_i} + s)(\mathbf{y}_i) \geq (\alpha - \varepsilon/3) \text{OPT}_{g+s} / (1 + \varepsilon/3 + \varepsilon')$ and $(\hat{g}_{\mathcal{R}_i} + s)(\mathbf{y}_i) \geq (1 + \varepsilon/3 + \varepsilon') x_i$.*

Proof. First we know that $(\hat{g}_{\mathcal{R}_i} + s)(\mathbf{y}_i) \geq (\alpha - \varepsilon/3) \max_{\mathbf{y} \in \mathcal{P}} (\hat{g}_{\mathcal{R}_i} + s)(\mathbf{y})$. For any \mathcal{R} , let $X_i^{\mathcal{R}}(\mathbf{x}) = 1 - \prod_{v \in R_i} (1 - h_v(\mathbf{x}))$ where $\mathcal{R} = \{R_1, R_2, \dots, R_\theta\}$, then $X_i^{\mathcal{R}}(\mathbf{x}) \in [0, 1]$ and $\hat{g}_{\mathcal{R}}(\mathbf{x}) = \frac{n}{\theta} \sum_{i=1}^{\theta} X_i^{\mathcal{R}}(\mathbf{x})$. We also know that $X_i^{\mathcal{R}}(\mathbf{x})$ are independent. By Chernoff Bound (Proposition 2), we have

$$\begin{aligned}
&\Pr\{(\hat{g}_{\mathcal{R}_i} + s)(\mathbf{y}_i) \leq (\alpha - \varepsilon/3) \text{OPT}_{g+s} / (1 + \varepsilon/3 + \varepsilon')\} \\
&\leq \Pr\{(\alpha - \varepsilon/3) \max_{\mathbf{y} \in \mathcal{P}} (\hat{g}_{\mathcal{R}_i} + s)(\mathbf{y}) \leq (\alpha - \varepsilon/3) \text{OPT}_{g+s} / (1 + \varepsilon/3 + \varepsilon')\}
\end{aligned}$$

$$\begin{aligned}
&= \Pr\{\max_{\mathbf{y} \in \mathcal{P}}(\hat{g}_{\mathcal{R}_i} + s)(\mathbf{y}) \leq \text{OPT}_{g+s}/(1 + \varepsilon/3 + \varepsilon')\} \\
&\leq \Pr\{(\hat{g}_{\mathcal{R}_i} + s)(\mathbf{x}_{g+s}^*) \leq \text{OPT}_{g+s}/(1 + \varepsilon/3 + \varepsilon')\} \\
&= \Pr\{\hat{g}_{\mathcal{R}_i}(\mathbf{x}_{g+s}^*) \leq \text{OPT}_{g+s}/(1 + \varepsilon/3 + \varepsilon') - s(\mathbf{x}_{g+s}^*)\} \\
&= \Pr\left\{\sum_{j=1}^{\theta_i} X_j^{\mathcal{R}_i}(\mathbf{x}_{g+s}^*) \frac{n}{\theta_i} \leq \text{OPT}_{g+s}/(1 + \varepsilon/3 + \varepsilon') - s(\mathbf{x}_{g+s}^*)\right\} \\
&= \Pr\left\{\sum_{j=1}^{\theta_i} X_j^{\mathcal{R}_i}(\mathbf{x}_{g+s}^*) \leq \frac{\theta_i}{n} \text{OPT}_{g+s}/(1 + \varepsilon/3 + \varepsilon') - \frac{\theta_i}{n} s(\mathbf{x}_{g+s}^*)\right\} \\
&= \Pr\left\{\sum_{j=1}^{\theta_i} X_j^{\mathcal{R}_i}(\mathbf{x}_{g+s}^*) - \frac{\theta_i}{n} g(\mathbf{x}_{g+s}^*) \leq \frac{\theta_i}{n} \text{OPT}_{g+s}/(1 + \varepsilon/3 + \varepsilon') - \frac{\theta_i}{n} g(\mathbf{x}_{g+s}^*) - \frac{\theta_i}{n} s(\mathbf{x}_{g+s}^*)\right\} \\
&\leq \exp\left(-\frac{\varepsilon'^2}{2(1 + \varepsilon/3 + \varepsilon')^2} \frac{\theta_i}{n} \text{OPT}_{g+s}\right) \\
&\leq \exp\left(-\frac{(1 + \varepsilon/3 + \varepsilon')^2 \cdot x_i \varepsilon'^2 n \cdot (2 + \frac{2}{3}\varepsilon') \cdot (\ln \mathcal{N}(\mathcal{P}, \frac{\varepsilon/3}{L_2} x_i) + \ell \ln n + \ln 2 + \ln \log_2(n + \lambda k))}{2(\alpha - \varepsilon/3)(1 + \varepsilon/3 + \varepsilon')^2 n \varepsilon'^2 x_i}\right) \\
&\leq \frac{1}{2n^\ell \cdot \mathcal{N}(\mathcal{P}, \frac{\varepsilon/3}{L_2} x_i) \cdot \log_2(n + \lambda k)}.
\end{aligned}$$

Combined with the assumption that $\text{OPT}_{g+s} \geq (1 + \varepsilon/3 + \varepsilon')^2 \cdot x_i / (\alpha - \varepsilon/3)$, we have: if $(\hat{g}_{\mathcal{R}_i} + s)(\mathbf{y}_i) \geq (\alpha - \varepsilon/3)\text{OPT}_{g+s}/(1 + \varepsilon/3 + \varepsilon')$, then

$$(\hat{g}_{\mathcal{R}_i} + s)(\mathbf{y}_i) \geq (\alpha - \varepsilon/3)\text{OPT}_{g+s}/(1 + \varepsilon/3 + \varepsilon') \geq \frac{\alpha - \varepsilon/3}{1 + \varepsilon/3 + \varepsilon'} (1 + \varepsilon/3 + \varepsilon')^2 \cdot x_i / (\alpha - \varepsilon/3) = (1 + \varepsilon/3 + \varepsilon') x_i.$$

□

Then, with the previous high probability lemma, we can upper bound $\mathbb{E}[\theta^{(1)} + \theta^{(2)} + \theta_{i_{ret}}]$, where i_{ret} denote the index of Algorithm 2 that break from the for-loop. We use L_1, L_2 to denote the Lipschitz constant of the function $(g + s)$ and the function $(\hat{g}_{\mathcal{R}_i} + s)$.

Lemma 16. *Let i_{ret} denote the index of Algorithm 2 that break from the for-loop. If we have $n + \lambda k = O(n^\ell)$, then*

$$\mathbb{E}[\theta^{(1)} + \theta^{(2)} + \theta_{i_{ret}}] = O\left(\frac{n \cdot \ln\left(n^\ell \mathcal{N}(\mathcal{P}, \frac{\varepsilon}{L_1 + L_2} LB)\right)}{\varepsilon^2 \text{OPT}_{g+s}}\right).$$

Proof. Let x_{ret} denote the value $x_{i_{ret}}$ when Algorithm 2 break from the for-loop. We first prove that,

$$\mathbb{E}\left[\frac{1}{x_{ret}}\right] = O\left(\frac{1 + \frac{n + \lambda k}{n^\ell}}{\text{OPT}_{g+s}}\right).$$

Let i denote the smallest index such that $\text{OPT}_{g+s} \geq (1 + \varepsilon/3 + \varepsilon')^2 \cdot x_i / (\alpha - \varepsilon)$. From the previous lemma(Lemma 15), we know that with probability at least $1 - \frac{1}{2n^\ell \cdot \mathcal{N}(\mathcal{P}, \frac{\varepsilon}{L_2} x_i) \cdot \log_2(n + \lambda k)}$, we have

$$(\hat{g}_{\mathcal{R}_i} + s)(\mathbf{y}_i) \geq (1 + \varepsilon/3 + \varepsilon') x_i,$$

so the algorithm will break from for-loop and $x_{ret} \geq x_i$. Let \mathcal{A} denote the event that $(\hat{g}_{\mathcal{R}_i} + s)(\mathbf{y}_i) \geq (1 + \varepsilon/3 + \varepsilon') x_i$. Let $X = \frac{1}{x_{ret}}$ and we have

$$\mathbb{E}[X] = \mathbb{E}[X|\mathcal{A}] \Pr[\mathcal{A}] + \mathbb{E}[X|\neg\mathcal{A}] \Pr[\neg\mathcal{A}].$$

We know that $\mathbb{E}[X|\mathcal{A}] = O\left(\frac{1}{\text{OPT}_{g+s}}\right)$ and $\Pr[\mathcal{A}] \leq 1$. We also know that $\mathbb{E}[X|\neg\mathcal{A}] \leq 1$ since $x_j \geq 1$. Then

$$\Pr[\neg\mathcal{A}] \leq \frac{1}{2n^\ell \cdot \mathcal{N}(\mathcal{P}, \frac{\varepsilon/3}{L_2}x_i) \cdot \log_2(n + \lambda k)} = O\left(\frac{n + \lambda k}{n^\ell \text{OPT}_{g+s}}\right),$$

where we use the fact that $g(\mathbf{x}) \leq n$ and $s(\mathbf{x}) = \lambda(k - c(\mathbf{x})) \leq \lambda k$. Then we have

$$\mathbb{E}\left[\frac{1}{x_{ret}}\right] = O\left(\frac{1 + \frac{n + \lambda k}{n^\ell}}{\text{OPT}_{g+s}}\right) = O\left(\frac{1}{\text{OPT}_{g+s}}\right).$$

Recall that

$$\theta_i = \left\lceil \frac{n \cdot \left(2 + \frac{2}{3}\varepsilon'\right) \cdot \left(\ln \mathcal{N}(\mathcal{P}, \frac{\varepsilon/3}{L_2}x_i) + \ell \ln n + \ln 2 + \ln \log_2(n + \lambda k)\right)}{\varepsilon'^2 x_i} \right\rceil,$$

and

$$\theta^{(1)} = \frac{8n \cdot \ln(4n^\ell)}{\text{LB} \cdot (\alpha - \varepsilon/3)^2 \varepsilon^2 / 9}, \quad \theta^{(2)} = \frac{2\alpha' \cdot n \cdot \ln\left(4n^\ell \mathcal{N}(\mathcal{P}, \frac{\varepsilon/3}{L_1+L_2}\text{LB})\right)}{(\varepsilon/3 - \frac{1}{4}(\alpha - \varepsilon/3)^2 \varepsilon/3)^2 \text{LB}}.$$

We know that $\text{LB} \geq x_{ret}$, so we have

$$\mathbb{E}[\theta^{(1)} + \theta^{(2)} + \theta_{i_{ret}}] = O\left(\frac{n \cdot \ln\left(n^\ell \mathcal{N}(\mathcal{P}, \frac{\varepsilon/3}{L_1+L_2}\text{LB})\right)}{\varepsilon^2 / 9 \text{OPT}_{g+s}}\right) = O\left(\frac{n \cdot \ln\left(n^\ell \mathcal{N}(\mathcal{P}, \frac{\varepsilon}{L_1+L_2}\text{LB})\right)}{\varepsilon^2 \text{OPT}_{g+s}}\right).$$

□

The above lemma is the main lemma for the time complexity of Algorithm 1. Next, we show some existing propositions and lemmas, which will help to prove the time complexity. The following lemmas and assumptions comes from [23], and we use the lemmas and assumptions to prove Theorem 6 and 7.

Definition 4 (Martingale). *A series of random variables X_1, X_2, \dots is a martingale, if for all $i \geq 1$, $\mathbb{E}[|X_i|] < +\infty$ and $\mathbb{E}[X_{i+1}|X_1, \dots, X_i] = X_i$.*

Lemma 17 (Sufficient Condition for Martingale). *Suppose X_1, X_2, \dots, X_t are t random variables on $[0, 1]$, which satisfy that there exists a constant μ , $\mathbb{E}[X_i|X_1, \dots, X_{i-1}] = \mu$ for all $i \in [t]$. Let $Z_i = \sum_{j=1}^i (X_j - \mu)$, then Z_1, Z_2, \dots, Z_t is a martingale.*

Proposition 4 (Stopping Time of Martingale). *Let random variable τ denote the stopping time of a martingale $\{X_i\}_{i \geq 1}$. If there is a constant c that is independent to $\{X_i\}_{i \geq 1}$ and $\tau \leq c$, then $\mathbb{E}[X_\tau] = \mathbb{E}[X_1]$.*

Besides, we have the following assumptions.

Assumption 1. *In the triggering model, the time complexity to sample a triggering set T_v for any v is $O(|N^-(v)|)$, where $N^-(v)$ is the set of in-neighbors of node v .*

Assumption 2. *We assume that $\max_{\mathbf{x} \in \mathcal{P}} g(\mathbf{x}) \geq \max_{v \in V} \sigma(v)$, where $\sigma(v)$ denote the influence spread of node v . Besides, we have $\text{OPT}_{g+s} \geq \max_{v \in V} \sigma(v)$.*

Given a set $R \subseteq V$, let $\omega(R)$ denote the sum of in-degree of nodes in R . Based on Assumption 1, we know that generating a RR-set needs time $O(\omega(R) + 1)$. Next, we use $\text{EPT} = \mathbb{E}[\omega(R)]$ to denote the expectation of $\omega(R)$, and we have the following lemma, which also comes from [23].

Lemma 18.

$$\text{EPT} = \mathbb{E}[\omega(R)] = \frac{m}{n} \cdot \mathbb{E}_{\tilde{v}}[\sigma(\tilde{v})],$$

where $\sigma(\tilde{v})$ denote the influence spread of node \tilde{v} .

With the help of the previous lemma and proposition, we can prove Theorem 6 and Theorem 7, and their corresponding corollaries.

Theorem 6. *Under Assumption 1 and Assumption 2. Suppose that the proximal step can be finished in time T_{prox} . Besides, suppose that $c(\mathbf{x})$ is L_c -Lipschitz and β_c -smooth, $h_v(\mathbf{x})$ are L_h -Lipschitz and β_h -smooth for all v , and the gradient of $h_v(\mathbf{x})$ and $c(\mathbf{x})$ can be generated in time T_h and T_c . We also assume that the balance variable λ is also a constant and $n + \lambda k \leq n^l$. Then the expected running time of Algorithm 1 with proximal gradient descent oracle is bounded by*

$$O\left(\frac{\beta_h n^2 + L_h^2 n^3}{\varepsilon} \cdot \left((1 + T_h) \frac{(m+n) \cdot \ln\left(n^\ell \mathcal{N}\left(\mathcal{P}, \frac{\varepsilon}{2n^2 L_h + 2\lambda L_c}\right)\right)}{\varepsilon^2} + \log_2(n + \lambda k) T_c \right) + \frac{\log_2(n + \lambda k)(\beta_h n^2 + L_h^2 n^3) T_{prox}}{\varepsilon} \right).$$

Proof of Theorem 6. Let i_{ret} denote the index where Algorithm 2 break from the for-loop. First, note that in the sampling procedure, Algorithm 2 generate at most $2\theta_{i_{ret}} + \tilde{\theta}$ number of RR-sets, which is bounded by $O(\theta^{(1)} + \theta^{(2)} + \theta_{i_{ret}})$. We use $\tau = 2\theta_{i_{ret}} + \tilde{\theta}$ to denote the number of RR-sets. By Assumption 1, generating such number of RR-sets needs time $O(\sum_{j=1}^{\tau} (\omega(R_j) + 1))$. Let $W_i = \sum_{j=1}^i (\omega(R_j) - \text{EPT})$ for $i = 1, \dots, \tau$. Because the procedure to generate R_j is independent to the procedure that generates R_1, \dots, R_{j-1} , we have $\mathbb{E}[\omega(R_j) | \omega(R_1), \dots, \omega(R_{j-1})] = \mathbb{E}[\omega(R_j)] = \text{EPT}$. From Lemma 17, we know that $\{W_i\}_{i \leq \tau}$ is a martingale, and τ is a stopping time. Obviously, τ has an upper bound, which can be derived by setting $x_{ret} = 1$ and $LB = 1$. Then by the stopping time theorem (Proposition 4), the time complexity for generating RR-sets is $O(\mathbb{E}[\sum_{j=1}^{\tau} (\omega(R_j) + 1)]) = O(\mathbb{E}[W_\tau] + \mathbb{E}[\tau \cdot (\text{EPT} + 1)]) = O(\mathbb{E}[\theta^{(1)} + \theta^{(2)} + \theta_{i_{ret}}] \cdot (\text{EPT} + 1)) = O(\mathbb{E}[\tau] \cdot (\text{EPT} + 1))$.

Then, we count the total time complexity for calling the proximal gradient oracle. Each time we use the proximal gradient to optimize $\hat{g}_{\mathcal{R}_i} + s$ or $\hat{g}_{\mathcal{R}} + s$, we will use $O\left(\frac{\beta}{\varepsilon x_i}\right)$ or $O\left(\frac{\beta}{\varepsilon LB}\right)$ number of iterations, where β is the smoothness constant for $\hat{g}_{\mathcal{R}_i}$ and $\hat{g}_{\mathcal{R}}$. We can lower bound x_i and OB by 1, and by Lemma 13, we have $\hat{g}_{\mathcal{R}_i}$ and $\hat{g}_{\mathcal{R}}$ are all $(\beta_h n^2 + L_h^2 n^3)$ -smooth. So each time the number of iteration is upper bounded by $N = O\left(\frac{\beta_h n^2 + L_h^2 n^3}{\varepsilon}\right)$. Each time we need a proximal step, so the total proximal step time complexity is bounded by $O\left(\frac{\log_2(n + \lambda k)(\beta_h n^2 + L_h^2 n^3) T_{prox}}{\varepsilon}\right)$. Then we consider the time complexity that generate the gradient. By Lemma 13, the time complexity of generating gradient is $O(\sum_{R \in \mathcal{R}_i} |R|(1 + T_h) + T_c)$ or $O(\sum_{R \in \mathcal{R}} |R|(1 + T_h) + T_c)$, then the total time complexity for generating gradient is bounded by

$$O\left(N \cdot \left(2 \sum_{R \in \mathcal{R}_{i_{ret}}} |R|(1 + T_h) + \sum_{R \in \mathcal{R}} |R|(1 + T_h) + \log_2(n + \lambda k) T_c\right)\right).$$

By the martingale stopping time theorem (Proposition 4), we know that

$$\mathbb{E}\left[2 \sum_{R \in \mathcal{R}_{i_{ret}}} |R| + \sum_{R \in \mathcal{R}} |R|\right] = \mathbb{E}[\tau] \cdot \mathbb{E}[|R_1|],$$

and $|R_1| \leq \omega(R_1) + 1$ since each RR-set is weak connected. Then we have

$$\mathbb{E}\left[2 \sum_{R \in \mathcal{R}_{i_{ret}}} |R| + \sum_{R \in \mathcal{R}} |R|\right] \leq \mathbb{E}[\tau] \cdot (\text{EPT} + 1),$$

and the total expected time complexity of generating gradient is bounded by

$$O\left(N \cdot ((1 + T_h)(\text{EPT} + 1)\mathbb{E}[\tau] + \log_2(n + \lambda k) T_c)\right).$$

Then the total expected time complexity is bounded by (the time to generate RR-set is far less than the time to generate all of the gradient)

$$O\left(\frac{\beta_h n^2 + L_h^2 n^3}{\varepsilon} \cdot ((1 + T_h)(\text{EPT} + 1)\mathbb{E}[\tau] + \log_2(n + \lambda k)T_c) + \frac{\log_2(n + \lambda k)(\beta_h n^2 + L_h^2 n^3)T_{prox}}{\varepsilon}\right).$$

By Assumption 2, Lemma 13, Lemma 16 and Lemma 18, we have

$$\begin{aligned} (\text{EPT} + 1)\mathbb{E}[\tau] &= O\left(\frac{n \cdot \ln\left(n^\ell \mathcal{N}\left(\mathcal{P}, \frac{\varepsilon}{2n^2 L_h + 2\lambda L_c}\right)\text{LB}\right)}{\varepsilon^2 \text{OPT}_{g+s}} \left(\frac{m}{n} \cdot \mathbb{E}_{\tilde{v}}[\sigma(\tilde{v})] + 1\right)\right) \\ &\leq O\left(\frac{(m+n) \cdot \ln\left(n^\ell \mathcal{N}\left(\mathcal{P}, \frac{\varepsilon}{2n^2 L_h + 2\lambda L_c}\right)\right)}{\varepsilon^2}\right). \end{aligned}$$

The total time complexity is bounded by (let T_h and T_c be constants)

$$O\left(\frac{\beta_h n^2 + L_h^2 n^3}{\varepsilon} \cdot \left((1 + T_h) \frac{(m+n) \cdot \ln\left(n^\ell \mathcal{N}\left(\mathcal{P}, \frac{\varepsilon}{2n^2 L_h + 2\lambda L_c}\right)\right)}{\varepsilon^2} + \log_2(n + \lambda k)T_c\right) + \frac{\log_2(n + \lambda k)(\beta_h n^2 + L_h^2 n^3)T_{prox}}{\varepsilon}\right).$$

□

Theorem 7. *Under Assumption 1 and Assumption 2. Suppose that the projection step can be finished in time T_{proj} . Besides, suppose that $c(\mathbf{x})$ is L_c -Lipschitz and β_c -smooth, $q_{v,j}(x_j)$ are L_q -Lipschitz for all v, j , and the gradient and function value of $q_{v,i}(x_j)$ and $c(\mathbf{x})$ can be generated in time T_q and $T_c(L_c, L_q, \beta_c, T_q, T_c)$ are all constants. We also assume that the balance variable λ is also a constant and $n + \lambda k \leq n^l$. Then the expected running time of Algorithm 1 with proximal gradient descent oracle is bounded by*

$$O\left(\frac{(n^2 \sqrt{d} L_q + \lambda L_c)^2}{\varepsilon^2} \cdot \left((1 + T_h) \frac{(m+n) \cdot \ln\left(n^\ell \mathcal{N}\left(\mathcal{P}, \frac{\varepsilon}{2n^2 \sqrt{d} L_q + \lambda L_c}\right)\right)}{\varepsilon^2} + \log_2(n + \lambda k)T_c\right) + \frac{(n^2 \sqrt{d} L_q + \lambda L_c)^2 \log_2(n + \lambda k)T_{proj}}{\varepsilon^2}\right).$$

Proof of Theorem 7. The proof of this theorem is almost the same as the previous one. We only have to change the iteration step from $O\left(\frac{\beta_h n^2 + L_h^2 n^3}{\varepsilon}\right)$ to $O\left(\frac{(n^2 \sqrt{d} L_q + \lambda L_c)^2}{\varepsilon^2}\right)$ and the Lipschitz constant for $\hat{g}_{\mathcal{R}}$ and g from $n^2 L_h$ to $n^2 \sqrt{d} L_q$. □

The above 2 theorems are the main theorem of the the time complexity. First, we have the fact that when $c(\mathbf{x}) = \|\mathbf{x}\|_1$ or $c(\mathbf{x}) = \|\mathbf{x}\|_2$, then proximal step and the projection step can all be finished in time $\tilde{O}(d)$, and the time to generate the function value and gradient of $c(\mathbf{x})$ is also $O(d)$. Besides, as shown in [25], we know that the covering number $\mathcal{N}(\mathcal{P}, \varepsilon) \leq \mathcal{N}(\mathbb{B}_1(k), \varepsilon) \leq \mathcal{N}(\mathbb{B}_2(k), \varepsilon) \leq (3k/\varepsilon)^d$. Then since $\|\mathbf{x}\|_1 \geq \|\mathbf{x}_2\|$, we have $\mathbb{B}_1(k) \subseteq \mathbb{B}_2(k)$, and $\mathcal{N}(\mathbb{B}_1(k), \varepsilon) \leq (3k/\varepsilon)^d$. We can directly get the time complexity bounds in Theorem 5. $\mathcal{N}(\mathcal{P}, \varepsilon) \leq \mathcal{N}(\mathbb{B}_1(k), \varepsilon) \leq \mathcal{N}(\mathbb{B}_2(k), \varepsilon) \leq (3k/\varepsilon)^d$. Then since $\|\mathbf{x}\|_1 \geq \|\mathbf{x}_2\|$, we have $\mathbb{B}_1(k) \subseteq \mathbb{B}_2(k)$, and $\mathcal{N}(\mathbb{B}_1(k), \varepsilon) \leq (3k/\varepsilon)^d$. We can directly get the time complexity bounds in Theorem 5.

C.2 Time Complexity Bound Based on the Moments of RR Set Size

In this subsection, we give the proof of our time complexity bounds based on the moments of the size of the RR-sets and the optimal value OPT_{g+s} for ProxGrad-RIS and UpperGrad-RIS. We use $\nu^{(1)}$, $\nu^{(2)}$ and $\nu^{(3)}$ to denote the first, second and third moments of the mean size of a random generated RR-set. Formally, we have

$$\nu^{(1)} = \mathbb{E}_R[|R|], \quad \nu^{(2)} = E_R[|R|^2], \quad \nu^{(3)} = E_R[|R|^3].$$

Given the time complexity theorem in Theorem 8 and 9 and some statistics for the variable $\nu^{(1)}$, $\nu^{(2)}$, $\nu^{(3)}$, we can know why in experiments, ProxGrad-RIS and UpperGrad-RIS do not make so much difference with the heuristic greedy algorithm in terms of the running time.

To derive the time complexity bounds based on the moments of the size of RR-sets, we need to slightly revise our sampling algorithm. In Algorithm 2, we use the previous generated RR-sets and only generate $\theta_i - \theta_{i-1}$ RR-sets in round i (line 6 of Algorithm 2). However, in this subsection, we assume that we generate θ_i RR-sets in round i and we do not use the previous generated RR-sets. This is because we need to use the martingale stopping time for each round, and to construct a martingale, we need the empirical moments $\nu^{(1)}(\mathcal{R}_i), \nu^{(2)}(\mathcal{R}_i), \nu^{(3)}(\mathcal{R}_i)$ to be independent to \mathcal{R}_j for all $j < i$. The following theorem summarizes the time complexity of ProxGrad-RIS with the above resampling of RR sets adjustment.

Theorem 8. *Under Assumption 1 and Assumption 2. Suppose that the proximal step can be finished in time T_{prox} . Besides, suppose that $c(\mathbf{x})$ is L_c -Lipschitz and β_c -smooth, $h_v(\mathbf{x})$ are L_h -Lipschitz and β_h -smooth for all v , and the gradient of $h_v(\mathbf{x})$ and $c(\mathbf{x})$ can be generated in time T_h and T_c . We also assume that the balance variable λ is also a constant and $n + \lambda k \leq n^l$. Then the expected running time of ProxGrad-RIS (revised by resampling of RR sets in each sampling iteration) is bounded by*

$$O\left(\frac{n^2(\nu^{(2)}\beta_h + \nu^{(3)}L_h^2) \ln\left(n^\ell \mathcal{N}\left(\mathcal{P}, \frac{\varepsilon}{2n^2\sqrt{d}L_q + \lambda 2L_c}\right)\right)}{\varepsilon^3 \text{OPT}_{g+s}}(1 + T_h) + \frac{\nu^{(1)}n\beta_h + \nu^{(2)}nL_h^2}{\varepsilon} \log_2(n + \lambda k)(T_c + T_{\text{prox}}) + \frac{(m+n) \cdot \ln\left(n^\ell \mathcal{N}\left(\mathcal{P}, \frac{\varepsilon}{2n^2\sqrt{d}L_q + \lambda 2L_c}\right)\right)}{\varepsilon^2}\right).$$

Proof of Theorem 8. The first part is to compute the expected time to generate the RR-sets. This part is similar to the same part in the proof of Theorem 6 and 7. Let i_{ret} denote the index where Algorithm 2 break from the for-loop, and we know that Algorithm 2 generate at most $2\theta_{i_{\text{ret}}} + \tilde{\theta}$ number of RR-sets, which is bounded by $O(\theta^{(1)} + \theta^{(2)} + \theta_{i_{\text{ret}}})$. We use $\tau = 2\theta_{i_{\text{ret}}} + \tilde{\theta}$ to denote the number of RR-sets. By Assumption 1, generating such number of RR-sets needs time $O(\sum_{j=1}^{\tau}(\omega(R_j) + 1))$. Let $W_i = \sum_{j=1}^i(\omega(R_j) - \text{EPT})$ for $i = 1, \dots, \tau$. Because the procedure to generate R_j is independent to the procedure that generates R_1, \dots, R_{j-1} , we have $\mathbb{E}[\omega(R_j)|\omega(R_1), \dots, \omega(R_{j-1})] = \mathbb{E}[\omega(R_j)] = \text{EPT}$. From Lemma 17, we know that $\{W_i\}_{i \leq \tau}$ is a martingale, and τ is a stopping time. Obviously, τ has an upper bound, which can be derived by setting $x_{\text{ret}} = 1$ and $LB = 1$. Then by the stopping time theorem(Proposition 4), the expected time complexity for generating RR-sets is $O(\mathbb{E}[\sum_{j=1}^{\tau}(\omega(R_j) + 1)]) = O(\mathbb{E}[W_\tau] + \mathbb{E}[\tau \cdot (\text{EPT} + 1)]) = O(\mathbb{E}[\theta^{(1)} + \theta^{(2)} + \theta_{i_{\text{ret}}}] \cdot (\text{EPT} + 1)) = O(\mathbb{E}[\tau] \cdot (\text{EPT} + 1))$. Based on Assumption 1 and 2 and Lemma 16 and 18, we know that

$$\begin{aligned} O(\mathbb{E}[\tau] \cdot (\text{EPT} + 1)) &= O\left(\frac{n \cdot \ln\left(n^\ell \mathcal{N}\left(\mathcal{P}, \frac{\varepsilon}{L_1 + L_2} LB\right)\right)}{\varepsilon^2 \text{OPT}_{g+s}} \cdot \left(\frac{m}{n} \cdot \mathbb{E}_{\tilde{v}}[\sigma(\tilde{v})] + 1\right)\right) \\ &= O\left(\frac{(m+n) \cdot \ln\left(n^\ell \mathcal{N}\left(\mathcal{P}, \frac{\varepsilon}{2n^2\sqrt{d}L_q + \lambda 2L_c}\right)\right)}{\varepsilon^2}\right). \end{aligned}$$

Then we count the time to generate the gradients and the proximal steps. In round i , Algorithm 2 will use the RR-sets \mathcal{R}_i with $|\mathcal{R}_i| = \theta_i$. In ProxGrad-RIS, in the i -th round, the algorithm iterates for

$O\left(\frac{\nu^{(1)}(\mathcal{R}_i)n\beta_h + \nu^{(2)}(\mathcal{R}_i)nL_h^2}{\varepsilon}\right)$ times. Then, the total time to generate the gradient and to complete the proximal steps in round i are bounded by

$$\begin{aligned} & O\left(\frac{\nu^{(1)}(\mathcal{R}_i)n\beta_h + \nu^{(2)}(\mathcal{R}_i)nL_h^2}{\varepsilon} \left(\sum_{R \in \mathcal{R}_i} |R|(1 + T_h) + T_c\right)\right) \\ &= O\left(\frac{\nu^{(1)}(\mathcal{R}_i)n\beta_h + \nu^{(2)}(\mathcal{R}_i)nL_h^2}{\varepsilon} \left(\theta_i \nu^{(1)}(\mathcal{R}_i)(1 + T_h) + T_c\right)\right) \\ &\leq O\left(\frac{\nu^{(2)}(\mathcal{R}_i)n\beta_h + \nu^{(3)}(\mathcal{R}_i)nL_h^2}{\varepsilon} \theta_i(1 + T_h) + \frac{\nu^{(1)}(\mathcal{R}_i)n\beta_h + \nu^{(2)}(\mathcal{R}_i)nL_h^2}{\varepsilon} T_c\right), \end{aligned}$$

where the last line comes from the fact that $\nu^{(1)}(\mathcal{R}_i)^2 \leq \nu^{(2)}(\mathcal{R}_i)$ and $\nu^{(1)}(\mathcal{R}_i)\nu^{(2)}(\mathcal{R}_i) \leq \nu^{(3)}(\mathcal{R}_i)$. $\nu^{(1)}(\mathcal{R}_i)^2 \leq \nu^{(2)}(\mathcal{R}_i)$ comes directly from the Cauchy-Schwartz inequality, and $\nu^{(1)}(\mathcal{R}_i)\nu^{(2)}(\mathcal{R}_i) \leq \nu^{(3)}(\mathcal{R}_i)$ comes from the fact that for any $a, b \geq 0$, we have $a^3 + b^3 \geq ab^2 + a^2b$. Then we sum up all of them from $i = 1$ to $i = i_{ret}$ and bound the expectation. We want to bound

$$O\left(\mathbb{E}\left[\sum_{i=1}^{i_{ret}} \frac{\nu^{(2)}(\mathcal{R}_i)n\beta_h + \nu^{(3)}(\mathcal{R}_i)nL_h^2}{\varepsilon} \theta_i(1 + T_h) + \frac{\nu^{(1)}(\mathcal{R}_i)n\beta_h + \nu^{(2)}(\mathcal{R}_i)nL_h^2}{\varepsilon} T_c\right]\right).$$

First note that the second part is easy to compute, and it is just bounded by

$$O\left(\frac{\nu^{(1)}n\beta_h + \nu^{(2)}nL_h^2}{\varepsilon} \log_2(n + \lambda k) T_c\right),$$

since we know that $i_{ret} \leq \log_2(n + \lambda k)$. Then we bound the first term. For convenience, we use the following notations

$$Y_i = \frac{\nu^{(2)}(\mathcal{R}_i)n\beta_h + \nu^{(3)}(\mathcal{R}_i)nL_h^2}{\varepsilon} \theta_i(1 + T_h), \quad Z_i = \frac{\nu^{(2)}n\beta_h + \nu^{(3)}nL_h^2}{\varepsilon} \theta_i(1 + T_h).$$

It is easy to show that $\mathbb{E}Y_i = Z_i$, and let $W_j = \sum_{i=1}^j (Y_i - Z_i)$. We will show that W_j is a martingale. This is due to the fact that since we generate new RR-sets in each round, and each RR-set is independent to others, then $Y_i - Z_i$ is independent to all of the information before round i , then W_j is a martingale. Also notice that i_{ret} is a stopping time, since we only decide if the sampling procedure will stop based on the previous information. By the Martingale Stopping Time theorem(Proposition 4), we have

$$\mathbb{E} \sum_{i=1}^{i_{ret}} (Y_i - Z_i) = \mathbb{E}[Y_1 - Z_1] \mathbb{E}[i_{ret}] = 0.$$

Then we only have to bound $O(\mathbb{E} \sum_{i=1}^{i_{ret}} Z_i)$. We have the following

$$\sum_{i=1}^{i_{ret}} Z_i = \sum_{i=1}^{i_{ret}} \frac{\nu^{(2)}n\beta_h + \nu^{(3)}nL_h^2}{\varepsilon} \theta_i(1 + T_h) \leq 2 \frac{\nu^{(2)}n\beta_h + \nu^{(3)}nL_h^2}{\varepsilon} \theta_{i_{ret}}(1 + T_h).$$

Similarly, the total time to generate the gradient and to complete the proximal steps which are not counted in the sampling procedure(Algorithm 2) are bounded by

$$O\left(\frac{\nu^{(2)}(\mathcal{R})n\beta_h + \nu^{(3)}(\mathcal{R})nL_h^2}{\varepsilon} \tilde{\theta}(1 + T_h) + \frac{\nu^{(1)}(\mathcal{R})n\beta_h + \nu^{(2)}(\mathcal{R})nL_h^2}{\varepsilon} T_c\right).$$

Taking the expectation, because we generate the RR-sets independently and thus $\nu^{(2)}(\mathcal{R}), \nu^{(3)}(\mathcal{R})$ are independent to $\tilde{\theta}$, we have the following bound for the expected time

$$O\left(\frac{\nu^{(2)}n\beta_h + \nu^{(3)}nL_h^2}{\varepsilon} \mathbb{E}[\tilde{\theta}](1 + T_h) + \frac{\nu^{(1)}n\beta_h + \nu^{(2)}nL_h^2}{\varepsilon} T_c\right).$$

Then from Lemma 16, we have

$$\frac{\nu^{(2)}n\beta_h + \nu^{(3)}nL_h^2}{\varepsilon} \mathbb{E}[2\theta_{iret} + \tilde{\theta}](1 + T_h) \leq \frac{n^2(\nu^{(2)}\beta_h + \nu^{(3)}L_h^2) \ln \left(n^\ell \mathcal{N}(\mathcal{P}, \frac{\varepsilon}{2n^2\sqrt{d}L_q + \lambda 2L_c}) \right)}{\varepsilon^3 \text{OPT}_{g+s}} (1 + T_h).$$

Also, it is easy to show that the expected time for proximal steps is bounded by

$$\frac{\nu^{(1)}n\beta_h + \nu^{(2)}nL_h^2}{\varepsilon} \log_2(n + \lambda k) T_{prox},$$

since the algorithm applies the proximal steps in each iteration, which is the same as calling the gradient of function $c(\mathbf{x})$.

Combining all of them together (generating gradient, proximal steps, generating RR-sets), we know that the expected time complexity is

$$\begin{aligned} & O \left(\frac{n^2(\nu^{(2)}\beta_h + \nu^{(3)}L_h^2) \ln \left(n^\ell \mathcal{N}(\mathcal{P}, \frac{\varepsilon}{2n^2\sqrt{d}L_q + \lambda 2L_c}) \right)}{\varepsilon^3 \text{OPT}_{g+s}} (1 + T_h) \right. \\ & \left. + \frac{\nu^{(1)}n\beta_h + \nu^{(2)}nL_h^2}{\varepsilon} \log_2(n + \lambda k) (T_c + T_{prox}) + \frac{(m+n) \cdot \ln \left(n^\ell \mathcal{N}(\mathcal{P}, \frac{\varepsilon}{2n^2\sqrt{d}L_q + \lambda 2L_c}) \right)}{\varepsilon^2} \right). \end{aligned}$$

□

We remark that, when comparing Theorem 8 with Theorem 6, if we do the following relaxations for the corresponding terms in the bound of Theorem 8: $\nu^{(2)} \leq \nu^{(1)} \cdot n$, $\nu^{(3)} \leq \nu^{(1)} \cdot n^2$, $\nu^{(1)} \leq n$, and $\nu^{(2)} \leq n^2$, together with $\nu^{(1)}/\text{OPT}_{g+s} \leq m/n$ (Lemma 18), then we will have the bound given in Theorem 6. This indicates how loose is the bounds we give in Theorem 5 in the main text. In our experiments (Section ??), we will demonstrate how much this relaxation is numerically in our dataset.

The following theorem summarizes our result for the UpperGrad-RIS algorithm, revised with the resampling of RR sets in each iteration step as described before.

Theorem 9. *Under Assumption 1 and Assumption 2. Suppose that the projection step can be finished in time T_{proj} . Besides, suppose that $c(\mathbf{x})$ is L_c -Lipschitz and β_c -smooth, $q_{v,j}(x_j)$ are L_q -Lipschitz for all v, j , and the gradient and function value of $q_{v,1}(x_j)$ and $c(\mathbf{x})$ can be generated in time T_q and $T_c(L_c, L_q, \beta_c, T_q, T_c)$ are all constants). We also assume that the balance variable λ is also a constant and $n + \lambda k \leq n^\ell$. Then the expected running time of UpperGrad-RIS (revised by resampling of RR sets in each sampling iteration) is bounded by*

$$\begin{aligned} & O \left(\frac{n(n^2 d \nu^{(3)} L_q^2 + 2\lambda n \sqrt{d} \nu^{(2)} L_q L_c + \lambda^2 L_c^2) \ln \left(n^\ell \mathcal{N}(\mathcal{P}, \frac{\varepsilon}{2n^2\sqrt{d}L_q + \lambda 2L_c}) \right)}{\varepsilon^4 \text{OPT}_{g+s}} (1 + T_h) \right. \\ & \quad + \frac{(n^2 d \nu^{(3)} L_q^2 + 2\lambda n \sqrt{d} \nu^{(2)} L_q L_c + \lambda^2 L_c^2)}{\varepsilon^2} \log_2(n + \lambda k) (T_c + T_{proj}) \\ & \quad \left. + \frac{(m+n) \cdot \ln \left(n^\ell \mathcal{N}(\mathcal{P}, \frac{\varepsilon}{2n^2\sqrt{d}L_q + \lambda 2L_c}) \right)}{\varepsilon^2} \right). \end{aligned}$$

Proof of Theorem 9. The proof of this theorem is almost the same as the previous one. We only have to change the iteration step from $O \left(\frac{\nu^{(1)}(\mathcal{R})n\beta_h + \nu^{(2)}(\mathcal{R})nL_h^2}{\varepsilon} \right)$ to $O \left(\frac{(\nu^{(1)}(\mathcal{R})n\sqrt{d}L_q + \lambda L_c)^2}{\varepsilon^2} \right)$ and the Lipschitz constant for $\hat{g}_{\mathcal{R}}$ and g from n^2L_h to $n^2\sqrt{d}L_q$. □

D Properties of the Original Function $g(\mathbf{x})$

In this section, we discuss the properties of the original function $g(\mathbf{x})$. We will compute the gradient of $g(\mathbf{x})$ and show how to compute the stochastic gradient of $g(\mathbf{x})$. Then we give upper bounds for the smoothness constant of the function $g(\mathbf{x})$ and the variance of the stochastic gradient estimator $\widehat{\nabla}g(\mathbf{x})$ in terms of its L2-norm difference with the true gradient, defined as $\text{Var} = \mathbb{E}[|\widehat{\nabla}g(\mathbf{x}) - \nabla g(\mathbf{x})|_2^2]$. These bounds would lead to our settings of the step size and number of iterations for the stochastic gradient method on the original objective function $g(\mathbf{x})$.

The following lemma provides the exact gradient formula for $g(\mathbf{x})$.

Lemma 19. *The gradient of function $g(\mathbf{x})$ can be written as:*

$$\nabla g(\mathbf{x}) = \sum_{u' \in V} \left[\sum_{S: u' \notin S} (\sigma(S \cup \{u'\}) - \sigma(S)) \left(\prod_{u \in S} h_u(\mathbf{x}) \right) \left(\prod_{v \notin S \cup \{u'\}} (1 - h_v(\mathbf{x})) \right) \cdot \nabla h_{u'}(\mathbf{x}) \right].$$

Proof. We first recall the expression of the function $g(\mathbf{x})$,

$$g(\mathbf{x}) = \mathbb{E}_S[\sigma(S)] = \sum_{S \subseteq V} \left(\sigma(S) \left(\prod_{u \in S} h_u(\mathbf{x}) \right) \left(\prod_{v \notin S} (1 - h_v(\mathbf{x})) \right) \right).$$

The gradient of $g(\mathbf{x})$ is given as follow

$$\begin{aligned} \nabla g(\mathbf{x}) &= \nabla_{\mathbf{x}} \left(\sum_{S \subseteq V} \left(\sigma(S) \left(\prod_{u \in S} h_u(\mathbf{x}) \right) \left(\prod_{v \notin S} (1 - h_v(\mathbf{x})) \right) \right) \right) \\ &= \sum_{S \subseteq V} \sigma(S) \nabla_{\mathbf{x}} \left(\left(\prod_{u \in S} h_u(\mathbf{x}) \right) \left(\prod_{v \notin S} (1 - h_v(\mathbf{x})) \right) \right) \\ &= \sum_{S \subseteq V} \sigma(S) \left(\sum_{u' \in S} \nabla_{\mathbf{x}} h_{u'}(\mathbf{x}) \left(\prod_{u \in S, u \neq u'} h_u(\mathbf{x}) \right) \left(\prod_{v \notin S} (1 - h_v(\mathbf{x})) \right) \right. \\ &\quad \left. - \sum_{v' \notin S} \nabla_{\mathbf{x}} h_{v'}(\mathbf{x}) \left(\prod_{u \in S} h_u(\mathbf{x}) \right) \left(\prod_{v \notin S, v \neq v'} (1 - h_v(\mathbf{x})) \right) \right). \end{aligned}$$

We then rearrange the gradient term, we have

$$\begin{aligned} \nabla g(\mathbf{x}) &= \sum_{u' \in V} \left[\sum_{S: u' \in S} \sigma(S) \left(\prod_{u \in S, u \neq u'} h_u(\mathbf{x}) \right) \left(\prod_{v \notin S} (1 - h_v(\mathbf{x})) \right) \cdot \nabla h_{u'}(\mathbf{x}) \right. \\ &\quad \left. - \sum_{T: u' \notin T} \sigma(T) \left(\prod_{u \in T} h_u(\mathbf{x}) \right) \left(\prod_{v \notin T, v \neq u'} (1 - h_v(\mathbf{x})) \right) \cdot \nabla h_{u'}(\mathbf{x}) \right] \\ &= \sum_{u' \in V} f_{u'}(\mathbf{x}) \nabla h_{u'}(\mathbf{x}), \end{aligned}$$

where $f_{u'}(\mathbf{x})$ is defined as

$$f_{u'}(\mathbf{x}) := \sum_{S: u' \notin S} (\sigma(S \cup \{u'\}) - \sigma(S)) \left(\prod_{u \in S} h_u(\mathbf{x}) \right) \left(\prod_{v \notin S \cup \{u'\}} (1 - h_v(\mathbf{x})) \right). \quad (8)$$

□

Lemma 19 provides the exact gradient formula for $g(\mathbf{x})$, but it involves an exponentially large number of summation terms and cannot be efficiently computed. Instead, we use Lemma 19 to define a simple stochastic gradient estimator $\widehat{\nabla}g(\mathbf{x})$ as the unbiased estimator of $\nabla g(\mathbf{x})$.

Definition 5 (Stochastic Gradient Estimator). *For any vector $x \in \mathcal{D}$, we construct a stochastic gradient estimator $\widehat{\nabla}g(\mathbf{x})$ as follows: First, we sample a node $u' \in V$ uniformly at random from V . Then, we sample a subset $S \subseteq V \setminus \{u'\}$ according to $h_u(\mathbf{x})$ for $u \in V \setminus \{u'\}$, more specifically, for each $u \in V \setminus \{u'\}$, we include u in S with probability $h_u(\mathbf{x})$ and exclude u from S with probability $1 - h_u(\mathbf{x})$, and different u 's are sampled independently. Next we sample a live-edge graph L based on the triggering model, and compute the marginal gain of u' on S in graph L , denoted $\sigma_L(u'|S)$, which is the number of nodes in L that can be reached from u' but not from S . Finally, we set $\widehat{\nabla}g(\mathbf{x}) = n \cdot \sigma_L(u'|S) \nabla h_{u'}(\mathbf{x})$.*

It is straightforward to see that $\widehat{\nabla}g(\mathbf{x})$ is an unbiased estimator of $\nabla g(\mathbf{x})$, i.e. $\mathbb{E}[\widehat{\nabla}g(\mathbf{x})] = \nabla g(\mathbf{x})$.

The next lemma provides the bounds on the smoothness of $g(\mathbf{x})$ and the variance of its stochastic gradient estimator defined above.

Lemma 20. *Assuming that the function $h_u(\mathbf{x})$ is L_h -Lipschitz and β_h -smooth, then we have the following bound for the smoothness constant β_g of $g(\mathbf{x})$ and the variance of stochastic gradient estimator $\text{Var} = \mathbb{E}[|\widehat{\nabla}g(\mathbf{x}) - \nabla g(\mathbf{x})|_2^2]$.*

$$\beta_g \leq \beta_h n^2 + 2L_h^2 n^3, \quad (9)$$

$$\text{Var} \leq 4L_h^2 n^4. \quad (10)$$

Proof. In this proof, we use the same $f_{u'}(\mathbf{x})$ definition as given in Eq. (8). we have

$$\begin{aligned} \|\nabla g(\mathbf{x}) - \nabla g(\mathbf{y})\|_2 &= \left\| \sum_{u' \in V} f_{u'}(\mathbf{x}) \nabla h_{u'}(\mathbf{x}) - \sum_{u' \in V} f_{u'}(\mathbf{y}) \nabla h_{u'}(\mathbf{y}) \right\|_2 \\ &= \left\| \sum_{u' \in V} f_{u'}(\mathbf{x}) \nabla h_{u'}(\mathbf{x}) - \sum_{u' \in V} f_{u'}(\mathbf{x}) \nabla h_{u'}(\mathbf{y}) \right. \\ &\quad \left. + \sum_{u' \in V} f_{u'}(\mathbf{x}) \nabla h_{u'}(\mathbf{y}) - \sum_{u' \in V} f_{u'}(\mathbf{y}) \nabla h_{u'}(\mathbf{y}) \right\|_2 \\ &\leq \left\| \sum_{u' \in V} f_{u'}(\mathbf{x}) \nabla h_{u'}(\mathbf{x}) - \sum_{u' \in V} f_{u'}(\mathbf{x}) \nabla h_{u'}(\mathbf{y}) \right\|_2 \\ &\quad + \left\| \sum_{u' \in V} f_{u'}(\mathbf{x}) \nabla h_{u'}(\mathbf{y}) - \sum_{u' \in V} f_{u'}(\mathbf{y}) \nabla h_{u'}(\mathbf{y}) \right\|_2 \\ &\leq \sum_{u' \in V} |f_{u'}(\mathbf{x})| \cdot \|\nabla h_{u'}(\mathbf{x}) - \nabla h_{u'}(\mathbf{y})\|_2 \\ &\quad + \sum_{u' \in V} |f_{u'}(\mathbf{x}) - f_{u'}(\mathbf{y})| \cdot \|\nabla h_{u'}(\mathbf{y})\|_2 \\ &\leq \sum_{u' \in V} |f_{u'}(\mathbf{x})| \cdot \beta_h \|\mathbf{x} - \mathbf{y}\|_2 + \sum_{u' \in V} |f_{u'}(\mathbf{x}) - f_{u'}(\mathbf{y})| \cdot L_h, \end{aligned}$$

where the last inequality comes from the assumptions we made before.

Then, it is easy to see that

$$f_{u'}(\mathbf{x}) \leq \sum_{S: u' \notin S} n \cdot \left(\prod_{u \in S} h_u(\mathbf{x}) \right) \left(\prod_{v \notin S \cup \{u'\}} (1 - h_v(\mathbf{x})) \right) = n.$$

Then we bound $|f_{u'}(\mathbf{x}) - f_{u'}(\mathbf{y})|$. First, we define

$$r_{u'}(\mathbf{x}) = \sum_{S: u' \notin S} (\sigma(S \cup \{u'\}) - \sigma(S)) \left(\prod_{u \in S} x_u \right) \left(\prod_{v \notin S \cup \{u'\}} (1 - x_v) \right).$$

where the input for $r_{u'}(\cdot)$ is defined as $\mathbf{x} = (x_v)_{v \in V} \in \mathbb{R}^n$. Then we show that

$$|r_{u'}(\mathbf{x}) - r_{u'}(\mathbf{y})| \leq 4n \sum_{v \in V} |x_v - y_v|.$$

We first show that

$$|r_{u'}(\mathbf{x}_{-w}, x_w) - r_{u'}(\mathbf{x}_{-w}, x'_w)| \leq 4n|x_w - x'_w|.$$

First, if $w = u'$, then it is obvious that $r_{u'}(x_{-w}, x_w) - r_{u'}(x_{-w}, x'_w) = 0$, since the variable corresponding to w does not appear in the formula. Then assume that $u' \neq w$, we can get the following equation by rearranging the terms.

$$\begin{aligned} & r_{u'}(\mathbf{x}_{-w}, x_w) \\ = & x_w \sum_{S: u' \notin S, w \in S} (\sigma(S \cup \{u'\}) - \sigma(S)) \prod_{u \in S, u \neq w} x_u \prod_{v \notin S, v \neq u'} (1 - x_v) \\ & + (1 - x_w) \sum_{S: u', w \notin S} (\sigma(S \cup \{u'\}) - \sigma(S)) \prod_{u \in S} x_u \prod_{v \notin S, v \neq u', w} (1 - x_v) \end{aligned}$$

Then we have the following inequalities

$$\begin{aligned} & |r'_u(\mathbf{x}_{-w}, x_w) - r'_u(\mathbf{x}_{-w}, x'_w)| \\ = & \left| (x_w - x'_w) \sum_{S: u' \notin S, w \in S} (\sigma(S \cup \{u'\}) - \sigma(S)) \prod_{u \in S, u \neq w} x_u \prod_{v \notin S, v \neq u'} (1 - x_v) \right. \\ & \left. + (x'_w - x_w) \sum_{S: u', w \notin S} (\sigma(S \cup \{u'\}) - \sigma(S)) \prod_{u \in S} x_u \prod_{v \notin S, v \neq u', w} (1 - x_v) \right| \\ \leq & 2n|x_w - x'_w| \end{aligned}$$

The last inequality is due to the two sum terms are the expectations of marginal influence. The marginal influence is smaller than n due to submodularity.

Then given two vectors $\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$, we define $\mathbf{z}_0 = \mathbf{x}$, and $\mathbf{z}_i = (y_1, \dots, y_i, x_{i+1}, \dots, x_n)$ for $i = 1, 2, \dots, n-1$ and $\mathbf{z}_n = \mathbf{y}$. Then we have

$$\begin{aligned} |r_{u'}(\mathbf{x}) - r_{u'}(\mathbf{y})| &= \left| \sum_{i=0}^{n-1} (r_{u'}(\mathbf{z}_i) - r_{u'}(\mathbf{z}_{i+1})) \right| \\ &\leq \sum_{i=0}^{n-1} |r_{u'}(\mathbf{z}_i) - r_{u'}(\mathbf{z}_{i+1})| \\ &\leq \sum_{i=0}^{n-1} 2n|x_{i+1} - y_{i+1}| \\ &= 2n \sum_{v \in V} |x_v - y_v|. \end{aligned}$$

Then, let $h(\mathbf{x}) := (h_u(\mathbf{x}))_{u \in V}$ to be the vector for the activation probabilities, then we have $f_{u'}(\mathbf{x}) = r_{u'}(h(\mathbf{x}))$, and we have

$$|f_{u'}(\mathbf{x}) - f_{u'}(\mathbf{y})| = |r_{u'}(h(\mathbf{x})) - r_{u'}(h(\mathbf{y}))|$$

$$\begin{aligned}
&\leq 2n \sum_{v \in V} |h_v(\mathbf{x}) - h_v(\mathbf{y})| \\
&\leq 2n \sum_{v \in V} L_h \|\mathbf{x} - \mathbf{y}\|_2 \\
&= 2L_h n^2 \cdot \|\mathbf{x} - \mathbf{y}\|_2,
\end{aligned}$$

where the last inequality comes from the assumption that $h_v(\mathbf{x})$ is L_h -Lipschitz. Plug in all the terms, we have shown that

$$\begin{aligned}
\|\nabla g(\mathbf{x}) - \nabla g(\mathbf{y})\|_2 &\leq \sum_{u' \in V} |f_{u'}(\mathbf{x})| \cdot \beta_h \|\mathbf{x} - \mathbf{y}\|_2 + \sum_{u' \in V} |f_{u'}(\mathbf{x}) - f_{u'}(\mathbf{y})| \cdot L_h \\
&\leq \sum_{u' \in V} n \cdot \beta_h \|\mathbf{x} - \mathbf{y}\|_2 + \sum_{u' \in V} L_h \cdot w L_h n^2 \cdot \|\mathbf{x} - \mathbf{y}\|_2 \\
&= (\beta_h n^2 + 2n^3 L_h^2) \|\mathbf{x} - \mathbf{y}\|_2.
\end{aligned}$$

As shown above, the original function g is $(\beta_h n^2 + 2n^3 L_h^2)$ -smooth.

Next we bound $\mathbb{E}[\|\widehat{\nabla} g(\mathbf{x}) - \nabla g(\mathbf{x})\|_2^2]$. First, it is easy to show that $|f_{u'}(\mathbf{x})|$ is bounded by n , so we have

$$\|\widehat{\nabla} g(\mathbf{x})\|_2 \leq \sum_{u'} |f_{u'}(\mathbf{x})| \cdot \|\nabla h_{u'}(\mathbf{x})\|_2 \leq n^2 L_h.$$

Then similar to the argument above, we also have

$$\|\nabla g(\mathbf{x})\|_2 \leq n^2 L_h.$$

Then

$$\mathbb{E}[\|\widehat{\nabla} g(\mathbf{x}) - \nabla g(\mathbf{x})\|_2^2] \leq \mathbb{E}[(2n^2 L_h)^2] = 4L_h^2 n^4 = O(n^4).$$

□

The following corollary summarizes our stochastic gradient algorithm for the original objective function, which will be used in the experiment section. It is a direct consequence of Theorem 2 and Lemma 20, and the proof is omitted.

Corollary 1. *Assuming that the function $h_v(\mathbf{x})$ is L_h -Lipschitz and β_h -smooth. By choosing the unbiased stochastic gradient estimator shown in Definition 5, we run stochastic gradient algorithm which comes from [9]. The algorithm runs T rounds with gradient step size $\mu_t = \frac{1}{\beta_h n^2 + 2L_h^2 n^3 + \frac{2\sqrt{2}L_h n^2}{\Delta} \sqrt{t}}$, where $\Delta = \sup_{\mathbf{x}, \mathbf{y} \in \mathcal{P}} \|\mathbf{x} - \mathbf{y}\|_2$.*

Then

$$\mathbb{E} \left[\max_{t=0,1,2,\dots,T} g(\mathbf{x}_t) \right] \geq \frac{OPT}{2} - \left(\frac{\Delta^2 (\beta_h n^2 + 2L_h^2 n^3)}{4T} + \frac{\sqrt{2} \Delta L_h n^2}{\sqrt{T}} \right).$$

The above theorem shows that with the number of iterations $T = O\left(\frac{d^2 (\beta_h n^2 + L_h^2 n^3)}{\varepsilon} + \frac{d L_h n^2}{\varepsilon^2}\right)$, we obtain a solution whose objective value is at least $(\frac{OPT}{2} - \varepsilon)$.